


Universidad Técnica Federico Santa María
Departamento de Informática
Análisis Inteligente de Datos



Capítulo 2: Métodos Estadísticos II semestre de 2010

Héctor Allende: hallende@inf.utfsm.cl

Modelos Regresión (Metamodelo)

- El análisis de regresión consiste en determinar cómo una variable y , está relacionada (asociada, correlacionada o ligada) con una o más variables explicativas x_1, \dots, x_n

y	Respuesta Variable dependiente Salida
x_i	Regresores Variables explicativas Variable independiente Entrada

Profesor H.Allende

2

Modelo Regresión lineal Simple

Modelo de Regresión \equiv Modelo Explicativo Estático

$$y_{ij} = \beta_0 + \beta_1 x_i + u_{ij} \text{ (Hipótesis Estructural)}$$

Supuestos básicos

y_{ij}, u_{ij} : variables aleatorias dependiente ; β_0, β_1 : Parámetros y
 x_i : variable explicativa determinística (estocástica).

Supuestos distribucionales usuales.

- $E[u_{ij}] = 0$.
- $\text{Var}[u_{ij}] = \sigma^2$, Cte ; Perturbación independiente de x .
- $u_{ij} \sim N(0, \sigma^2)$.
- $E[u_{ij} u_{kh}] = 0, \forall (i,j) \neq (k,h) \Rightarrow$

Profesor H.Allende

3

Modelos Regresión Simple

1- $E[y_{ij}/x_i] = \beta_0 + \beta_1 x_i$

2- $\text{Var}[y_{ij}/x_i] = \sigma^2$

3- $f(y_{ij}/x_i)$ es normal

4- Las obs. son independientes entre si

1.2 Estimación de parámetros.

$$y_{ij} = \beta_0 + \beta_1 x_i + u_{ij} \Rightarrow \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \Leftrightarrow E[y_i/x_i] = \beta_0 + \beta_1 x_i$$

Residuo (e_{ij}) = Valor observado (y_{ij}) - Valor previsto (\hat{y}_i).

1.2.1 Método de Máxima Verosimilitud. (función de Verosimilitud).

$$\ell(\beta_0, \beta_1, \sigma^2, y_{ij}) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{1}{2\sigma^2} (y_{ij} - \beta_0 - \beta_1 x_i)^2 \right]$$

$$\Rightarrow L(\beta_0, \beta_1, \sigma^2) = \log \left\{ \prod_{i=1}^n \ell(\beta_0, \beta_1, \sigma^2, y_{ij}) \right\}$$

Profesor H.Allende

4

Modelos Regresión Simple

Derivando $L()$ con respecto a los parámetros :

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\text{Cov}(x, y)}{S_x^2} \quad y \quad \hat{\sigma}^2 = \frac{\sum \sum e_{ij}^2}{n}$$

Métodos: Mínimos Cuadrados y Máxima verosimilitud.

$$\text{Maximizar } L(\beta_0, \beta_1, \sigma^2) \Leftrightarrow \text{Min}_{(\beta_0, \beta_1)} \sum \sum (y_{ij} - \beta_0 - \beta_1 x_i)^2$$

$$M = \sum \sum (y_{ij} - \beta_0 - \beta_1 x_i)^2 = \sum \sum e_{ij}^2$$

$$\Rightarrow \hat{S}_R^2 = \frac{\sum \sum e_{ij}^2}{n-2}$$

Profesor H.Allende

5

Distribución de los Parámetros

$\hat{\beta}_0$ y $\hat{\beta}_1$ son variables aleatorias

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \approx N(\beta_0, V(\hat{\beta}_0))$$

$$\hat{\beta}_1 = \frac{\text{Cov}(x, y)}{S_x^2} \approx N(\beta_1, V(\hat{\beta}_1))$$

Profesor H.Allende

6

Propiedades de los estimadores.

Propiedades de la ley de probabilidades de β_1

$$E[\beta_1] = \beta_1 \quad \hat{\beta}_1 \approx N\left(\beta_1; \frac{\sigma^2}{nS_x^2}\right)$$

$$Var[\beta_1] = \frac{\sigma^2}{nS_x^2}$$

Intervalos de confianza

$$IC_{\gamma=(1-\alpha)} = \left[\hat{\beta}_1 \pm t_{1-\alpha/2} \frac{\hat{S}_R}{\sqrt{nS_x^2}} \right]$$

Profesor H.Allende

7

Propiedades de los estimadores.

Propiedades de la ley de probabilidades de β_0

$$E[\hat{\beta}_0] = \beta_0 \quad \hat{\beta}_0 \approx N\left(\beta_0; \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{nS_x^2} \right)\right)$$

$$Var[\hat{\beta}_0] = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{nS_x^2} \right)$$

Intervalos de confianza

$$IC_{\gamma=(1-\alpha)} = \left[\hat{\beta}_0 \pm t_{1-\alpha/2} \frac{\hat{S}_R}{\sqrt{n}} \sqrt{1 + \frac{\bar{x}^2}{S_x^2}} \right]$$

Profesor H.Allende

8

Propiedades de los estimadores.

Propiedades de la ley de probabilidades \hat{S}_R^2

$$E[\hat{S}_R^2] = \sigma^2 \quad \sum \sum \frac{e_{ij}^2}{\sigma^2} \approx \chi^2_{(n-2)}$$

$$Var[\hat{S}_R^2] = \frac{2\sigma^4}{n-2}$$

Intervalos de confianza para σ^2

$$IC(\sigma^2)_{\gamma=(1-\alpha)} = \left[\frac{(n-2)\hat{S}_R^2}{\chi^2_{1-\alpha/2}}; \frac{(n-2)\hat{S}_R^2}{\chi^2_{\alpha/2}} \right]$$

Profesor H.Allende

9

Contraste de regresión.

Prueba de hipótesis. $H_0: \beta_1 = 0$ v/s $H_1: \beta_1 \neq 0$

Sea Variación Total = $VT = \sum_i \sum_j (y_{ij} - \bar{y})^2$

Variación no Explicada = $VNE = \sum_i \sum_j (y_{ij} - \hat{y}_i)^2$

Variación Explicada = $VE = \sum_i n_i (\hat{y}_i - \bar{y})^2$

Estadístico $F_0 = \frac{\hat{\beta}_1^2 n S_x^2}{\hat{S}_R^2} = \frac{VE}{\hat{S}_R^2} \approx F_{(1, n-2)}$

Para la región crítica $P(F_{(1, n-2)} \leq C) = 1-\alpha$, se rechaza H_0 para $F_0 > C$.

Profesor H.Allende

10

Predicción de las medias condicionales.

$$\hat{y}_h = \bar{y} + \hat{\beta}_1 (x_h - \bar{x})$$

$$\Rightarrow E[\hat{y}_h] = \beta_0 + \beta_1 x_h = m_h \quad (E[y/x_h])$$

$$Var[\hat{y}_h] = \left\{ \frac{1}{n} + \frac{(x_h - \bar{x})^2}{\sum n_i (x_i - \bar{x})^2} \right\} \sigma^2 = v_{hh} \sigma^2$$

Intervalo de confianza para las medias. $IC_y(m_h) = [\hat{y}_h \pm t_{\alpha/2} \hat{S}_R \sqrt{v_{hh}}]$

Se desea prever el valor de y para $x = x_h$. Intuitivamente

$$\hat{y}_h = \hat{\beta}_0 + \hat{\beta}_1 x_h$$

Corresponde a la Estimaciones de las medias condicionales.

Criterio de predicción. Error cuadrático medio mínimo

$$E[(y - \tilde{y}/x)^2] = Var(y/x) + (E[y/x] - \tilde{y})^2$$

Profesor H.Allende

11

Intervalo de Confianza

Intervalo de Confianza para las observaciones $IC_\gamma(y_h) = [\hat{y}_h \pm t_{\alpha/2} \hat{S}_R \sqrt{1 + v_{hh}}]$

Coefficiente de Correlación.

Coefficiente de determinación: $R^2 = \frac{VE}{VT} = \frac{\sum (y_i - \bar{y})^2}{\sum (y_i - \bar{y})^2}$

Coefficiente de correlación: $r = \frac{Cov(x, y)}{S_x S_y}$

Relación entre \hat{S}_R^2 y r : $r = \left\{ 1 - \frac{(n-2)\hat{S}_R^2}{nS_y^2} \right\}^{1/2}$

Profesor H.Allende

12

Contraste de linealidad.

Hipótesis.

$$H_0: E[y_{ij}/x_i] = \beta_0 + \beta_1 x_i \quad \text{v/s} \quad H_1: E[y_{ij}/x_i] \neq \beta_0 + \beta_1 x_i$$

$$\text{Varianza entre las medias y las rectas: } \hat{S}_{12}^2 = \frac{\sum n_i (\bar{y}_i - \bar{y})^2}{d-2}$$

$$\text{Varianza de la perturbación sin la linealidad: } \hat{S}_2^2 = \frac{\sum \sum (y_{ij} - \bar{y}_i)^2}{n-d}$$

$$\text{Estadístico de Prueba: } F_0 = \frac{\hat{S}_{12}^2}{\hat{S}_2^2} \approx F_{(d-2, n-d)}$$

Región crítica $P(F_{(d-2, n-d)} \leq C) = 1 - \alpha$, se rechaza H_0 para $F_0 > C$.

Profesor H.Allende

13

Análisis de residuos.

Tiene por objeto contrastar a posteriori las hipótesis de linealidad del modelo. Es especialmente importante cuando se tiene un solo valor de la variable y para cada valor de la variable de control x .

El análisis de los residuos se utiliza para verificar:

- Si su distribución es aproximadamente normal.
- Si su variabilidad es constante y no depende de x .
- Si presentan evidencia de una relación no lineal.
- Si existen observaciones atípicas o heterogéneas.

Profesor H.Allende

14

Regresión Lineal

Modelo Lineal General:

$$y_j = \beta_0 + \beta_1 x_{1j} + \beta_2 x_{2j} + \dots + \beta_k x_{kj} + u_j \quad j = 1..n$$

$$u_j \sim N(0, \sigma^2) \quad \text{i.i.d}$$

Valor esperado:

$$E[y_j / x_1, x_2, \dots, x_k] = \beta_0 + \sum_{i=1}^k \beta_i x_{ij}$$

Profesor H.Allende

15

REGRESION GENERAL

La variable de respuesta " y " depende de muchas variables x_1, x_2, \dots, x_n , aunque algunas de estas son no observables.

El modelo de regresión pretende develar efecto de las variables explicativas más importantes y representa las restantes mediante una v.a. la perturbación.

$$\text{Es decir: } y = f(x_1, x_2, \dots, x_k) + \underbrace{g(x_{k+1}, \dots, x_n)}_{\mu}$$

Suponga que en el rango de interés, la función f admite una aproximación lineal:

$$f(x_1, x_2, \dots, x_k) = \sum_{j=0}^k \beta_j x_j, \text{ donde } x_0 = 1.$$

En tal caso

$$y = \sum_{j=0}^k \beta_j x_j + \mu.$$

Ejemplo: Modelo para predecir el valor de las propiedades en Viña del Mar de sus características físicas, geográficas etc.

Profesor H.Allende

16

REGRESION GENERAL

Se hacen las siguientes hipótesis sobre la distribución de las variables:

-Para cada conjunto fijo de las x , la distribución de y es normal

$$E[y/x_1, \dots, x_k] = \sum_{j=0}^k \beta_j x_j \quad \text{Var}[y/x_1, \dots, x_k] = \sigma^2 = cte.$$

Las variables y_i son independientes entre sí.

-El n° de variables explicativas es menor que el n° de observaciones.

-Las x 's son realmente distintas y no existen entre ellas relaciones lineales exactas.

$$\text{Luego } y_i = \beta_0 + \sum_{j=0}^k \beta_j x_{ji} + \mu_i.$$

Donde cada coeficiente β_j mide el efecto marginal sobre la respuesta de un aumento unitario en x_j .

μ_i : perturbación aleatoria ; $\mu_i \sim N[0, \sigma^2]$, $\forall i=1, \dots, n$.

Var $[\mu_i] = \sigma^2 = cte$, $\forall i=1, \dots, n$; $E[\mu_i \mu_j] = 0$, si $i \neq j$

Profesor H.Allende

17

Estimación de Parámetros

$$\text{Sea } y_i = \sum_{j=0}^k \beta_j x_{ji} + \mu_i \quad j=1, \dots, n; x_0=1 \quad \text{y sea } Q = \sum_{i=1}^n (y_i - \sum_{j=0}^k \beta_j x_{ji})^2$$

Bajo el supuesto de normalidad de la variable aleatoria y se sabe que

$$\underset{\underline{\beta}}{\text{Min}} Q \Leftrightarrow \underset{\underline{\beta}}{\text{Max}} \ln \prod_{i=1}^n f(y_i, \underline{\beta}, \sigma^2)$$

Derivando con respecto a β_0 y a β_j se obtiene las siguientes ecuaciones

$$\text{notación matricial: } X'Y = X'X\hat{\beta}$$

$$\text{donde } Y' = (y_1, \dots, y_n) \wedge X = \begin{pmatrix} 1 & x_{11} & \dots & x_{k1} \\ \vdots & \vdots & & \vdots \\ 1 & x_{1n} & \dots & x_{kn} \end{pmatrix}$$

Como por hipótesis $X'X$ no es singular se tiene que $\hat{\beta} = (X'X)^{-1} X'Y$

Profesor H.Allende

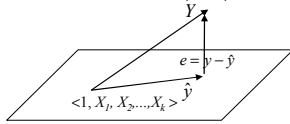
18

Interpretación geométrica.

Considere los vectores de \mathbb{R}^n : $1, X_1, X_2, \dots, X_k$ que forman las columnas de la matriz de diseño X . El objetivo de la estimación es determinar $\hat{\beta}$, como CL de \hat{Y}

$$\hat{Y} = \beta_0 1 + \beta_1 X_1 + \dots + \beta_k X_k$$

\hat{Y} está contenido en el subespacio generado por los vectores $\langle 1, X_1, X_2, \dots, X_k \rangle$. El criterio de mínimos cuadrados, impone que el norma del vector $\|e\| = \|Y - \hat{Y}\|$ sea mínima.



Del teorema de proyección se tiene que: $e \perp 1, X_1, \dots, X_k$ e \hat{Y}

Es decir $1'e = X_1'e = \dots = X_k'e = \hat{Y}'e = 0$
 $X'e = X'(Y - \hat{Y}) = X'(Y - X\hat{\beta}) = 0 \Rightarrow X'Y = X'X\hat{\beta}$

Profesor H.Allende

19

Interpretación geométrica.

Por lo tanto $\hat{Y} = Y + e$
 $\hat{Y} = X\hat{\beta} = X(X'X)^{-1}X'Y$
 $\hat{Y} = V'Y$

Siendo V' la matriz de proyección (simétrica e idempotente).

$$V' = V \quad y \quad V^2 = V$$

Esta matriz juega un rol importante en la etapa de diagnóstico.

$$e = Y - \hat{Y} = Y - V'Y = (I - V')Y$$

Profesor H.Allende

20

EJEMPLO 1

Ejemplo:

Los siguientes datos muestran el indicador global de desarrollo regional y, en términos del número de automóviles por mil habitantes (x_1) y el número de teléfonos por mil habitantes (x_2) en ocho de las 13 regiones del país.

Profesor H.Allende

21

Región	Indicador y	Automóviles x_1	Teléfonos x_2
I	64	58	111
II	78	84	131
III	83	78	158
IV	88	81	147
V	89	82	121
VI	99	102	165
VII	101	85	174
VIII	102	102	169

$$\hat{Y} = \beta_0 1 + \beta_1 \underline{x}_1 + \beta_2 \underline{x}_2 + u \quad Y = X\hat{\beta} + u \quad \text{con} \quad X = \begin{pmatrix} 1 & 58 & 111 \\ 1 & 84 & 131 \\ 1 & 78 & 158 \\ 1 & 81 & 147 \\ 1 & 82 & 121 \\ 1 & 102 & 165 \\ 1 & 85 & 174 \\ 1 & 102 & 169 \end{pmatrix}$$

$$X'X = \begin{pmatrix} 8 & 672 & 1176 \\ 672 & 57822 & 100453 \\ 1176 & 100453 & 176785 \end{pmatrix} \quad X'Y = \begin{pmatrix} 704 \\ 60251 \\ 105288 \end{pmatrix}$$

Resolviendo la ecuación matricial $X'Y = X'X\hat{\beta}$ se obtiene:

$$\hat{\beta}' = (9,05; 0,52; 0,24) = (\beta_0, \beta_1, \beta_2)$$

Profesor H.Allende

22

Conclusiones

Conclusiones.

1. Cualquier coeficiente de regresión estimado $\hat{\beta}_i$; puede interpretarse como la pendiente de la recta de regresión de los residuos de una regresión y respecto a todas las otras variables (parte de y no explicada por el resto de las x) con la contribución diferencial de x_i .
2. El coeficiente de regresión $\hat{\beta}_i$; tiene que interpretarse como el efecto diferencial de la variable x_i , eliminando los efectos de las otras variables explicativas.
3. El efecto sobre los coeficientes de regresión de las variables relevantes para explicar y , es distinto cuando las variables incluidas son independientes de las excluidas que cuando no lo son: en el primer caso no afectarán a los coeficientes $\hat{\beta}_i$, pero en el segundo pueden distorsionarlos apreciablemente.

Profesor H.Allende

23

Propiedades de los estimadores

2.3.1 Esperanza.

Sea: $C = (X'X)^{-1}X'$

Se puede demostrar que: $\hat{\beta} = \beta + Cu$

Luego, $E(\hat{\beta}) = E(\beta + Cu) = \beta + CE(u) = \beta$.

2.3.1 Covarianzas.

Sea $\hat{\beta} - \beta = Cu$

Se puede demostrar que: $\Sigma = E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)'] = \sigma^2(X'X)^{-1}$

Llamando q_{ij} a los elementos de la matriz Σ , se concluye que: $\hat{\beta}_j \sim N(\beta_j; \sigma^2 q_{jj})$

La matriz $(X'X)^{-1}$ en general no es diagonal, por lo tanto, su inversa tampoco lo será y los coeficientes $\hat{\beta}$ no serán independientes al no tener covarianzas nulas.

Profesor H.Allende

24

El Teorema de Gauss-Markov

El teorema de Gauss-Markov : Fundamento teórico principal del método de mínimos cuadrados en modelos lineales y establece que si las siguientes hipótesis son ciertas:

- Todos los valores de la variable aleatoria dependiente están generados por el modelo lineal: $Y = X\beta + U$
- Las perturbaciones u_i son no correlacionadas.
- Todas las perturbaciones tienen la misma varianza.
- Las perturbaciones son independientes de las v. a. X_i .
- Las variables X_i se obtienen sin errores de medida.
- Se desea encontrar el estimador (óptimo ECM), dentro de la clase de estimadores insesgados (centrados), que sean funciones lineales de Y . El estimador óptimo insesgado tendrá varianza mínima.

Entonces: Gauss-Markov aseguran que los estimadores mínimo cuadráticos son "óptimos" en el sentido restringido dado por f) - g), independiente de la distribución de U .

Profesor H.Allende

25

Estimación de la Varianza.

El modelo de regresión múltiple quedará especificado al estimar β y la varianza σ^2 de la perturbación $e = Y - \hat{Y} = (I - V)Y$

V es una matriz idempotente, luego $(I - V)$ también lo es. $\hat{\sigma}^2 = \frac{1}{n} e'e = U'(I - V)U$

La expresión es una forma cuadrática de variables aleatorias normales $N(0, \sigma^2)$ e independientes. Luego, $\frac{1}{\sigma^2} e'e \sim \chi^2_{rang(I - V), gl}$.

Como $(I - V)$ proyecta a Y sobre el complemento ortogonal al espacio definido por X , tendrá rango $n - k - 1$. $\Rightarrow \frac{1}{\sigma^2} e'e \sim \chi^2_{(n - k - 1)}$.

Finalmente, el estimador insesgado para σ^2 , llamado varianza residual es \hat{S}_R^2 :

$$\hat{S}_R^2 = \frac{\sum_{i=1}^n e_i^2}{n - k - 1}$$

Profesor H.Allende

26

Intervalos de Confianza y Pruebas de Hipótesis

Intervalos de confianza

Si se verifica que $\hat{\beta}_i$ y \hat{S}_R^2 son independientes, entonces $\frac{\hat{\beta}_i - \beta_i}{\hat{S}_R \sqrt{q_{ii}}} \sim t_{(n - k - 1), gl}$.

Luego, un intervalo de confianza para β_i de nivel $\gamma = 1 - \alpha$

$$IC_\gamma = \left[\hat{\beta}_i \pm t_{(n - k - 1)(\alpha/2), gl} \hat{S}_R \sqrt{q_{ii}} \right]$$

Pruebas o contrastes.

Se desea contrastar que la variable aleatoria $\hat{\beta}_i$ tiene media β_i^* . El test se realiza basado en la estadística t : siendo $E[\hat{\beta}_i] = \beta_i^*$.

$$t = \frac{\hat{\beta}_i - \beta_i^*}{\hat{S}_R \sqrt{q_{ii}}} \approx t_{(n - k - 1), gl}$$

Profesor H.Allende

27

Intervalos de Confianza y Pruebas de Hipótesis

Una prueba importante es $H_0: \beta_i^* = 0$. Bajo H_0 $\left(t_0 = \frac{\hat{\beta}_i}{\hat{S}_R \sqrt{q_{ii}}} \right) \sim t_{(n - k - 1), gl}$.

Rechazándose H_0 para $t_0 > c$ (valor crítico).

Regiones de confianza para conjuntos de coeficientes.

Como los coeficientes $\hat{\beta}_i$ son dependientes, Los intervalos de confianza individuales pueden dar una imagen errónea de sus valores conjuntos.

Sea $F = \frac{(\hat{\beta} - \beta)'(X'X)(\hat{\beta} - \beta)}{(k + 1)\hat{S}_R^2} \sim F_{(k + 1, n - k - 1), gl}$.

Luego, la región de confianza de nivel $(1 - \alpha)$ se obtiene calculando un valor crítico de la tabla F : $P(F_{(k + 1, n - k - 1)} \geq F_c) = \alpha$

Entonces, el elipsoide confidencial contendrá aquellos valores β tales que:

$$(\hat{\beta} - \beta)'(X'X)(\hat{\beta} - \beta) \leq F_c \hat{S}_R^2 (k + 1)$$

Profesor H.Allende

28

Intervalos de confianza para la varianza

Intervalos de confianza para la varianza.

Un intervalo de confianza de nivel $\gamma = 1 - \alpha$ para σ^2 es:

$$IC_\gamma = \left[\frac{(n - k - 1)\hat{S}_R^2}{\chi^2_{(1 - \alpha/2)(n - k - 1), gl}}, \frac{(n - k - 1)\hat{S}_R^2}{\chi^2_{\alpha/2(n - k - 1), gl}} \right]$$

Para intervalos de confianza de una cola:

$$IC_\gamma = \left[0; \frac{(n - k - 1)\hat{S}_R^2}{\chi^2_{(\alpha, n - k - 1)}} \right]$$

Profesor H.Allende

29

Contraste de regresión

El contraste de regresión para coeficientes individuales.

$H_0: \beta_h = 0$ v/s $H_1: \beta_h \neq 0$, Estadística $t = \left(\hat{\beta}_h / (\hat{S}_R \sqrt{q_{hh}}) \right) \sim t_{(n - k - 1), gl}$.

Usando ANDEVA.

VE(k): Variación explicada por el modelo completo.

VE(k-1): Variación explicada por el modelo sin la variable x_h .

$$\Delta VE = VE(k) - VE(k-1) = (\hat{\beta}_h / q_{hh})$$

Si $\beta_h = 0$, ΔVE depende solo del error experimental.

Luego, una estadística $F = \frac{\Delta VE}{\hat{S}_R^2} = F_{(1, n - k - 1)}$

Profesor H.Allende

30

Contraste de regresión

El contraste de regresión para grupos de coeficientes.

$$H_0: \beta_1 = \dots = \beta_k = 0 \quad \text{v/s} \quad H_1: \text{algún } \beta_i \neq 0, \quad i=1, \dots, k$$

Sea $\hat{\beta}^*$ el vector de coeficientes que no incluye a la componente β_0

$$F = \frac{\hat{\beta}^{*'} (X^{*'} X^*)^{-1} \hat{\beta}^*}{k \hat{\sigma}_R^2} \approx F_{(k, n-k-1)}$$

Descomposición de la varianza.

$$\text{Pitágoras Generalizado: } \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Tabla de ANDEVA.

Fuente	Suma de Cuadrados	g.l	Varianza	Contraste
VE	$\sum (y_i - \bar{y})^2$	k	$\hat{\sigma}_e^2$	
VNE	$\sum (y_i - \hat{y}_i)^2$	n-k-1	$\hat{\sigma}_R^2$	$F = \frac{\hat{\sigma}_e^2}{\hat{\sigma}_R^2}$
VT	$\sum (y_i - \bar{y})^2$	n-1	$\hat{\sigma}_y^2$	

Profesor H.Allende

31

Correlación en Regresión Múltiple

El contraste de regresión establece que la VE es significativamente mayor que

$$\text{VNE. Bajo } H_0, \quad F_{(k, n-k-1)gl} \sim F = \frac{VE}{k \hat{\sigma}_R^2}$$

El coeficiente de determinación.

Es una medida descriptiva global del ajuste de un modelo: $R^2 = \frac{VE}{VT} = \frac{\sum (y_i - \bar{y})^2}{\sum (y_i - \bar{y})^2}$

Al valor R se le denomina coeficiente de correlación múltiple.

Observaciones.

- Desde un punto de vista estricta la correlación se define solo para v.a., al ser X variables fijas el nombre no es totalmente correcto.
- R² aumenta cuando k aumenta.
- R² es muy sensible con respecto a la formulación del modelo y a la elección de la variable dependiente "y".

Profesor H.Allende

32

El coeficiente de determinación corregido.

Para evitar que R² aumente cuando k aumenta, se define un R²-corregido como:

$$\bar{R}^2 = 1 - \frac{\sum e_i^2 / (n-k-1)}{\sum (y_i - \bar{y})^2 / (n-1)}$$

Donde se verifica: 1) $\bar{R}^2 = 1 - (1 - R^2) \left(\frac{n-1}{n-k-1} \right)$

$$2) \quad \hat{\sigma}_R^2 = \hat{\sigma}_e^2 \left(1 - \bar{R}^2 \right)$$

R² y el Test de F Regresión.

Una forma alternativa para contrastar la hipótesis de que todos los coeficiente de regresión son cero es:

$$F = \frac{VE}{k} \cdot \frac{1}{\hat{\sigma}_R^2} = \frac{VE}{k} \cdot \frac{1}{\hat{\sigma}_R^2} = \frac{VE}{k} \cdot \frac{1}{\hat{\sigma}_R^2}$$

Mientras $1 - R^2 = \frac{VNE}{VT}$

Luego, el contraste F de regresión puede escribirse: $F_{(k, n-k-1)} = \left(\frac{R^2}{1 - R^2} \right) \left(\frac{n-k-1}{k} \right)$

Profesor H.Allende

33

Correlación Parcial Múltiple

Correlación Parcial.

Dado un conjunto de variables (x_1, \dots, x_p) , el coeficiente de correlación parcial entre dos de ellas, algún x_i y x_j , es una medida adimensional de su relación lineal, cuando se eliminan de ambas los efectos debidos al resto de las variables.

Definición:

Consideremos k regresores (x_1, \dots, x_k) . Entonces el coeficiente de correlación parcial entre x_1 y x_2 se define como el coeficiente de correlación Lineal de Pearson entre x_1 y x_2 .

$$r_{12.34\dots k} = \text{cov}(e_{1.34\dots k}, e_{2.34\dots k}) / \hat{\sigma}_{e_{1.34\dots k}} \hat{\sigma}_{e_{2.34\dots k}}$$

Es decir $r_{12.34\dots k}$ es el coeficiente del modelo $e_{1.34\dots k} = (r_{12.34\dots k}) * e_{2.34\dots k} + u$

Donde $e_{1.34\dots k}$ y $e_{2.34\dots k}$ son los residuos de la regresión múltiple de x_1 y x_2 con respecto al resto de las variables de control (x_3, x_4, \dots, x_k) .

$$x_1 = \beta_0 + \sum_{j=3}^k \beta_j x_j + e_{1.34\dots k}$$

$$x_2 = \beta_0 + \sum_{j=3}^k \beta_j x_j + e_{2.34\dots k}$$

Profesor H.Allende

34

Al estar los residuos depurados de los efectos de las restantes variables, el $r_{12.34\dots k}$ representa la relación entre x_1 y x_2 que no pueden explicarse por las variables restantes.

El coeficiente de correlación parcial entre las variables de respuesta y un regresor x_i (notación: $r_{y_i.R}$) se obtiene fácilmente a partir de la estadística "t": $t_i = \hat{\beta}_i / \hat{\sigma}(\hat{\beta}_i)$

$$\text{Entonces } r_{y_i.R}^2 = \frac{t_i^2}{t_i^2 + n - k - 1}$$

Regresión con variables ortogonales

Es un caso especial de regresión múltiple donde todas las variables explicativas satisfacen

$$\begin{aligned} \sum_{i=1}^k (x_{ij} - \bar{x}_j) (x_{ih} - \bar{x}_h) &= 0, \quad \forall j, h \\ \Rightarrow X'X &= \text{Diagonal}(\sum (x_{1j} - \bar{x}_j)^2, \dots, \sum (x_{kj} - \bar{x}_j)^2) \\ \Rightarrow \hat{\beta}_j &= \frac{\sum (x_{ij} - \bar{x}_j) (y_i - \bar{y})}{\sum (x_{ij} - \bar{x}_j)^2} = \frac{\sum \tilde{x}_{ij} \tilde{y}_i}{\sum \tilde{x}_{ij}^2} \end{aligned}$$

Profesor H.Allende

35

Predicción

Predicción del valor medio.

La predicción del valor medio de la respuesta para ciertos valores concretos de las variables explicativas $x_h' = (1, x_{1h}, \dots, x_{kh})$ será: $\hat{y}_h = x_h' \hat{\beta}$

$$E(\hat{y}_h) = x_h' \beta = m_h$$

$$\text{Var}(\hat{y}_h) = \sigma^2 x_h' (X'X)^{-1} x_h = \sigma^2 v_{hh} = \frac{\sigma^2}{\hat{v}_{hh}}$$

Intervalo de confianza para m_h .

Un intervalo de confianza para m_h de nivel $\gamma = 1 - \alpha$ es:

$$IC_\gamma = \left[\hat{y}_h \pm t_{\alpha/2} \hat{\sigma}_R \sqrt{v_{hh}} \right]$$

Profesor H.Allende

36

Predicción de una observación.

Predicción de una observación.

La predicción de una observación y_h no observada se efectúa mediante la media de la distribución condicionada, \hat{y}_h dado x_h

$$\hat{y}_h = x_h' \hat{\beta} \quad E(y_h) = m_h$$

Error cuadrático medio de la predicción.

$$E[(y_h - \hat{y}_h)^2] = \sigma^2(1 + v_{hh})$$

Intervalo de confianza para m_h .

Un intervalo de confianza para y_h de nivel $\gamma = 1 - \alpha$ está dado por:

$$IC_\gamma = [\hat{y}_h \pm t_{\alpha/2} \hat{S}_R \sqrt{1 + v_{hh}}]$$

En la siguiente sección se describen los problemas principales que surgen al construir un modelo de regresión.

Profesor H.Allende

37

Diagnóstico y validación de los modelos de regresión múltiple.

Hipótesis	Realidad
Las variables X toman valores distintos en la muestra.	Multicolinealidad: Las variables X toman valores semejantes en la muestra.
$E[y] = \beta' X$	Error de especificación, $E[y] \neq \beta' X$
La distribución de u es normal.	Falta de normalidad: u no es normal.
$Var[u] = cte.$ Homocedasticidad.	$Var[u] \neq cte.$ Heterocedasticidad.
u independientes entre si.	Autocorrelación: u dependientes.

Profesor H.Allende

38

Modelos Regresión No-Lineal

Regresión No-Lineal:

Modelo Allométrico:

$$y_j = \beta_0 x_{1j}^{\beta_1} + \varepsilon_j \quad j = 1..m$$

Usado para representar la relación existente entre el peso de una parte de una planta respecto a toda la planta

Modelo Mitscherlich:

$$y_j = \beta_0 \left(1 - e^{-\beta_1(x_{1j} + \beta_2)}\right) + \varepsilon_j \quad j = 1..m$$

y_j Rendimiento del cultivo
 x_{1j} Cantidad de fertilizante utilizado
 β_0 Límite superior del cultivo
 β_1 Cantidad de fertilizante que ya hay en la tierra
 β_2 Rapidez de crecimiento del cultivo

Relaciona el rendimiento de una cosecha con la cantidad de fertilizante utilizado

Profesor H.Allende

39

Modelos Lineales Generalizados

Generalización:

- La distribución de la variable de salida no necesariamente tiene que ser "Normal".
- El valor esperado de la salida viene dado por:

$$g(E(y_j)) = \beta_0 \sum_{i=1}^k \beta_i x_{ij}$$

$g(\cdot)$ es una función monótona y diferenciable y se conoce como función de enlace.

Profesor H.Allende

40

Modelos Aditivos Generales

Modelo:

$$g(E(y_j)) = \beta_0 \sum_{i=1}^k s_i(x_{ij})$$

donde s_i es una función arbitraria, usualmente suaves.

Ej: B-splines.

Profesor H.Allende

41

Análisis de Variancia (ANOVA)

- Método usado para identificar que parámetros son significativamente distintos de cero en el modelo lineal.
- La Técnica desarrollada por R. Fisher ha sido ampliamente usada en el análisis de experimentos en la agricultura y en la industria.

Profesor H.Allende

42

Ejemplo "Turnips for winter fodder"

- Ejemplo: "Nabos para el forraje invernal"
Experimento para investigar el crecimiento de los nabos.

Variety	Treatments			Label	Blocks			
	Date	Density			I	II	III	IV
Barkant	21-08-1990	1 Kg/ha	A	2.7	1.4	1.3	3.8	
		2 Kg/ha	B	7.3	3.8	3.0	1.2	
		4 Kg/ha	C	6.5	4.6	4.7	0.8	
		8 Kg/ha	D	8.2	4.0	6.0	2.5	
	28-08-1990	1 Kg/ha	E	4.4	0.4	6.5	3.1	
		2 Kg/ha	F	2.6	7.1	7.0	3.2	
		4 Kg/ha	G	24.0	14.9	14.6	2.6	
		8 Kg/ha	H	12.2	18.9	15.6	9.9	
Marco	21-08-1990	1 Kg/ha	J	1.2	1.3	1.5	1.0	
		2 Kg/ha	K	2.2	2.0	2.1	2.5	
		4 Kg/ha	L	2.2	6.2	5.7	0.6	
		8 Kg/ha	M	4.0	2.8	10.8	3.1	
	28-08-1990	1 Kg/ha	N	2.5	1.6	1.3	0.3	
		2 Kg/ha	P	5.5	1.2	2.0	0.9	
		4 Kg/ha	Q	4.7	13.2	9.0	2.9	
		8 Kg/ha	R	14.9	13.3	9.3	3.6	

Profesor H.Allende

43

Ejemplo "Turnips for winter fodder"

- Modelo lineal utilizado

$$y_j = \beta_0 + \beta_B x_{Bj} + \beta_C x_{Cj} + \dots + \beta_R x_{Rj} + \beta_{II} x_{II,j} + \beta_{III} x_{III,j} + \beta_{IV} x_{IV,j} + \varepsilon_j \quad j = 1..64$$

$$x_{ij} = \begin{cases} 1 & \text{si el punto } j \text{ coincide con la bloque/tratamiento } i \\ 0 & \text{en todo otro caso} \end{cases}$$

- ¿Puede un cambio en el tratamiento cambiar el crecimiento de los nabos?
- ¿Existe alguna constante distinta de cero?
- ANOVA

Profesor H.Allende

44

ANOVA

- Considerando el modelo de regresión lineal general:

$$y_j = \beta_0 + \beta_1 x_{1j} + \beta_2 x_{2j} + \dots + \beta_k x_{kj} + u_j \quad j = 1..n$$

$$u_j \sim N(0, \sigma^2) \quad \text{i.i.d}$$

- Se estiman los parámetros a partir de los datos:

$$\hat{y}_j = \hat{\beta}_0 + \sum_{i=1}^k \hat{\beta}_i x_{ij}$$

Profesor H.Allende

45

ANOVA

- Sean los residuos: $e_j = Y_j - Y'_j$

- El tamaño de los residuos está relacionado con σ^2 , y corresponde a la varianza de los u's.

- Estimando σ^2

$$S^2 = \frac{\sum_{j=1}^m (y_j - \hat{y}_j)^2}{m - (k + 1)}$$

Suma al cuadrado de los residuos

Grados de libertad de los residuos

Profesor H.Allende

46

ANOVA

- Características de S^2 para comparar modelos:

- Si el modelo ajustado es el adecuado, entonces S^2 es una buena estimación de σ^2
- Si el modelo ajustado incluye términos redundantes, entonces S^2 es aún una buena estimación.
- Si el modelo ajustado no incluye una o más entradas que deberían estar, entonces S^2 tendería a ser más grande que σ^2

Profesor H.Allende

47

ANOVA

- Sea Ω_1 un modelo lineal ajustado al conjunto de datos y sea

$$S_1^2 = \frac{RSS_1}{v_1}$$

la estimación de la varianza.

- Sea Ω_0 un modelo lineal ajustado al conjunto de datos, pero con algunos parámetros en 0, y

$$S_0^2 = \frac{RSS_0}{v_0}$$

es la estimación de la varianza.

- Si S_0^2 es grande comparado con S_1^2 , entonces Ω_0 es un modelo inadecuado.
- Si son similares, entonces Ω_0 es un modelo satisfactorio.

Profesor H.Allende

48

Anova

Sean

- Suma de cuadrados extra $ESS = RSS_0 - RSS_1$
- Grados extra de libertad $v_E = v_0 - v_1$
- Estimación de la varianza $S_E^2 = \frac{ESS}{v_E}$

Test de hipótesis:

- Estadístico de test: $F = \frac{S_E^2}{S_1^2} = \frac{ESS / v_E}{RSS_1 / v_1}$
- F tiene una distribución F con parámetros v_E y v_1
- Si F es aproximadamente 1 entonces Ω_0 es adecuado.
- Si F es grande, entonces el modelo es inadecuado.

$$H_0 : \Omega_0 \text{ es el modelo correcto}$$

$$H_1 : \Omega_1 \text{ es el modelo correcto}$$

Profesor H.Allende

49

Ejemplo "Turnips for winter fodder"

Modelo lineal utilizado

$$\Omega_0 : y_j = \beta_0 + \beta_{II}x_{II,j} + \beta_{III}x_{III,j} + \beta_{IV}x_{IV,j} + \varepsilon_j \quad j = 1..64$$

$$\Omega_1 : y_j = \beta_0 + \beta_Bx_{Bj} + \beta_Cx_{Cj} + \dots + \beta_Rx_{Rj} + \beta_{II}x_{II,j} + \beta_{III}x_{III,j} + \beta_{IV}x_{IV,j} + \varepsilon_j \quad j = 1..64$$

Tablas ANOVA

	Df	Sum of Sq.	Mean Sq.	F. Value	Pr(F)
Block	3	163,73700	54,57891	2,27802	0,08867543
Residuals	60	1,437,53800	23,95897		

	Df	Sum of Sq.	Mean Sq.	F. Value	Pr(F)
Block	3	163,73700	54,57891	5,69043	0,00216381
Treat	15	1,005,92700	67,06182	6,99191	0,00000017
Residuals	45	431,61100	9,59135		

Profesor H.Allende

50

Ejemplo "Turnips for winter fodder"

Donde la fila de residuos viene dado por:

$$RSS_0 = 1437.5$$

$$v_0 = 60$$

$$s_0^2 = \frac{1437.5}{60} = 23.96$$

$$RSS_1 = 431.6$$

$$v_1 = 45$$

$$s_1^2 = \frac{431.6}{45} = 9.59$$

- Como las estimaciones de σ^2 son muy diferentes, entonces algunas de las entradas eliminadas de Ω_0 son necesarias

La fila de tratamiento:

$$ESS = RSS_0 - RSS_1 = 1437.5 - 431.6 = 1005.9$$

$$v_E = v_0 - v_1 = 60 - 45 = 15$$

$$s_E^2 = \frac{1005.9}{15} = 67.06$$

Profesor H.Allende

51

Ejemplo "Turnips for winter fodder"

Obteniendo el estadístico F:

$$\frac{s_E^2}{s_1^2} = \frac{67.06}{9.59} = 6.99$$

- Nivel de significancia se obtiene de la columna Pr(F)
- La fila "block" proviene al comparar el modelo Ω_0 con:

$$y_j = \beta_0 + \varepsilon_j \quad j = 1, 2, \dots, 64$$

Profesor H.Allende

52

Ejemplo "Turnips for winter fodder"

	Df	Sum of Sq.	Mean Sq.	F. Value	Pr(F)
Block	3	163,7367	54,5789	5,6904	0,0021638
Variety	1	83,9514	83,9514	8,7528	0,0049136
sowing	1	233,7077	233,7077	24,3665	0,0000114
density	3	470,3780	156,7927	16,3473	0,0000003
variety: sowing	1	36,4514	36,4514	3,8005	0,0574875
variety: density	3	8,6467	2,8822	0,3005	0,8248459
sowing: density	3	154,7930	51,5977	5,3796	0,0029884
variety: sowing: density	3	17,9992	5,9997	0,6256	6,6022439
Residuals	45	1,437,5380	23,9590		

Profesor H.Allende

53

Modelos Log-lineales

- Es una forma de investigar las relaciones entre variables categóricas.
- Es un tipo de modelo GLM:

$$Y_j \sim \text{Pois}(\mu_j)$$

$$\log(\mu_j) = \beta_0 + \beta_1x_{1j} + \beta_2x_{2j} + \dots + \beta_nx_{nj}$$

- Las asociaciones corresponden a los términos de interacción en el modelo.
 - Problema: Determinar qué parámetros β son nulos \rightarrow ANOVA
 - Considerar una cantidad llamada *desviación* cuando se desea comparar dos modelos.

Profesor H.Allende

54

Ejemplo: "Breast cancer"

Centre	Age	Survived	State of Tumour			
			Minimal Malignant Appearance	Inflammation Benign Appearance	Greater Malignant Appearance	Inflammation Benign Appearance
Tokyo	Under 50	No	9	7	4	3
		Yes	26	68	25	9
	50-60	No	9	9	11	2
		Yes	20	46	18	5
	70 or over	No	2	3	1	0
		Yes	1	6	5	1
Boston	Under 50	No	6	7	6	0
		Yes	11	24	4	0
	50-60	No	8	20	3	2
		Yes	18	58	10	3
	70 or over	No	9	18	3	0
		Yes	15	26	1	1
Glamorgan	Under 50	No	16	7	3	0
		Yes	16	20	8	1
	50-60	No	14	12	3	0
		Yes	27	39	10	4
	70 or over	No	3	7	3	0
		Yes	12,0	11	4	1

Profesor H.Allende

55

Ejemplo: "Breast cancer"

	Df	Deviance	Resid. Df	Resid. Dev	Pr(>F)
NULL			71		
centre	2	9,36190	69		
age	2	105,535	67		
survived	1	160,6009	66		
inflam	1	291,1986	65		
appear	1	7,5727	64		
centre:age	4	76,9628	60		
centre:survived	2	11,2698	58		
centre:inflam	2	23,2484	56		
centre:appear	2	13,3323	54		
age:survived	2	3,5257	52		
age:inflam	2	0,293	50		
age:appear	2	1,2082	48		
survived:inflam	1	0,9645	47		
survived:appear	1	9,6709	46		
inflam:appear	1	95,4381	45		

Profesor H.Allende

56

Construcción de modelos de regresión

- Una de las mayores dificultades en regresión múltiple es la construcción de un buen modelo, debido a que existen dos problemas
- Se requiere incorporar la mayor cantidad de variables para alcanzar la mayor explicación de Y.
 - Se requiere tener la menor cantidad de variables para no aumentar la varianza.

Profesor H.Allende

57

Construcción de modelos de regresión

- Existen diversas maneras de seleccionar un buen modelo de regresión, una de ellas es crear todos los modelos posibles en base a las variables existentes y seleccionar el mejor a través de alguna medida de calidad.
- Otra manera es construir un modelo en forma dinámica, mediante:
 - Eliminación Progresiva.
 - Introducción Progresiva
 - Regresión paso a paso

Profesor H.Allende

58

Construcción de modelos de regresión

- Eliminación Progresiva
- El modelo se construye considerando todas las variables.
 - Seleccionar la variable cuyo estadístico F sea mínimo de todos (X_j)
 - Si el valor del estadístico es menor que un umbral de eliminación (F_{OUT}) reconstruir el modelo sin la variable e ir al paso 2.
 - Fin del algoritmo.

Profesor H.Allende

59

Construcción de modelos de regresión

- Introducción Progresiva
- Selección de la variable más correlacionada con Y.
 - Construcción del modelo de regresión simple $Y = \beta_0 + \beta_1 X_j$.
 - Verificar a través del estadístico F si X_j aporta significativamente al modelo.
 - Selección de la variable con mayor correlación a Y dado el modelo existente, lo que equivale a seleccionar la variable con mayor estadístico F.
 - Si el valor del estadístico F de la variable seleccionada es mayor a F_{IN} , incorporar la variable al modelo e ir a 4.
 - Fin del algoritmo.

Profesor H.Allende

60

Construcción de modelos de regresión

- Regresión paso a paso.
 - Este método combina los dos algoritmos descritos anteriormente, para poder incorporar y eliminar variables a medida que construye el modelo.
 - Este método se basa en introducción progresiva, pero cuando se incorpora una nueva variable al modelo, se procede a eliminar las variables que ya pertenecían al modelo pero que ya no aportan en forma significativa a este.

Profesor H.Allende

61

Construcción de modelos de regresión

Indices de calidad para evaluar los modelos

Coef. de determinación múltiple: $R^2 = \frac{SS_R}{SS_{yy}} = 1 - \frac{SS_E}{SS_{yy}}$

Coef. de correlación corregido: $\bar{R}^2 = 1 - \frac{SS_E / (n - p)}{SS_{yy} / (n - 1)}$

Varianza residual: $MS_E = (1 - \bar{R}^2) \frac{SS_{yy}}{n - 1}$

Estadístico Cp de Mallows: $C_p = \frac{SS_E(p)}{\hat{\sigma}^2} - n + 2p$

Criterio de Akaike: $AIC = n \ln(\hat{\sigma}^2) + n + 2p$

Profesor H.Allende

62

Ejemplo de Regresión Simple

t	0	1	2	3	4	5	6
V(t)	30	60	46	32	10	4	17
	20	40		26	14	8	
		20			12		
$\bar{V}(t)$	25	40	46	29	12	6	17

Sea $x_t = \sin t$ $y_t = \bar{V}(t)$

Luego $y(t) = a + b x_t + u_t$

$$\min_{a,b} Q(a,b) = \min_{a,b} \sum_t (y_t - a - b x_t)^2$$

Profesor H.Allende

63

$$\hat{a} = \bar{y} - \hat{b}\bar{x} = 25,3 \quad \hat{b} = \frac{\text{cov}(x,y)}{S_x^2} = 20$$

$$S_y^2 = 1276 \quad \sum (y_t - \hat{y}_t)^2 = 22,45$$

% de Ajuste del Modelo =

$$1 - \frac{\sum \hat{e}_t^2}{S_y^2} = 0,98 * 100\% = 98\%$$

Profesor H.Allende

64