

# ECONOMETRÍA

Iker Caballero Bragagnini  
Undergraduate Level

## Tabla de contenido

INTRODUCCIÓN A LA ECONOMETRÍA.....	2
EL MODELO DE REGRESIÓN LINEAL.....	6
LOS MÍNIMOS CUADRADOS.....	14
EL ESTIMADOR DE MÍNIMOS CUADRADOS .....	26
LOS CONTRASTES DE HIPÓTESIS Y LA SELECCIÓN DEL MODELO .....	38
LAS VARIABLES BINARIAS O CATEGÓRICAS Y LOS EXPERIMENTOS.....	47
LA NO LINEALIDAD EN LAS VARIABLES.....	52
LA EVALUACIÓN DE LOS EFECTOS DE TRATAMIENTOS .....	52
EL MODELO DE REGRESIÓN GENERALIZADO Y LA HETEROSCEDASTICIDAD .....	ERROR! BOOKMARK NOT DEFINED.
LOS SISTEMAS DE ECUACIONES .....	ERROR! BOOKMARK NOT DEFINED.
LOS MODELOS PARA DATOS DE PANEL.....	60
LOS RESULTADOS BINARIOS Y LA ELECCIÓN DISCRETA .....	91
LOS MODELOS DE ELECCIÓN DISCRETA MULTINOMIAL DESORDENADA .....	101
LOS MODELOS DE ELECCIÓN DISCRETA MULTINOMIAL ORDENADA .....	117
LOS MODELOS PARA EL CONTEO DE EVENTOS.....	125
LOS MODELOS DE TRUNCACIÓN Y CENSURA .....	133
LOS MODELOS DE DOS ECUACIONES.....	148
LOS MODELOS DE DURACIÓN .....	152
LA CORRELACIÓN SERIAL .....	163
LOS DATOS NO ESTACIONARIOS .....	175

## Introducción a la econometría

- La econometría se basa en la aplicación de la estadística matemática y las herramientas de la inferencia estadística para la medición empírica de las relaciones postuladas por la teoría económica
  - Uno de los objetivos principales es hacer un análisis causal: analizar cualitativa y cuantitativamente cómo ciertos factores afectan a una variable asociada a un fenómeno económico de interés
    - Determinar los efectos de ciertas políticas, caracterizar y cuantificar la relación de comportamiento entre variables económicas y simular los efectos de políticas alternativas son algunas de las cosas que permiten los modelos econométricos
  - En econometría se suelen usar datos no experimentales u observacionales, de modo que los valores de la variable objeto de estudio como de los factores susceptibles de afectarla están fuera del control del economista
    - Por eso no se puede interpretar una correlación entre variables como un efecto causal
    - Los datos experimentales, en cambio, se determinan en entornos controlados (se controlan los valores de aquellos factores susceptibles de afectar al fenómeno de estudio)
- Dentro de la econometría es posible hacer distinciones entre diversas aplicaciones y áreas
  - La conexión entre los modelos de comportamiento subyacentes y la práctica moderna de la econometría ha incrementado mucho y es útil hacer una distinción entre microeconometría y macroeconometría
    - La microeconometría se caracteriza por el análisis de datos transversales y de panel centrado en los consumidores individuales, las empresas y los decisores a nivel microeconómico. Los practicantes se apoyan en las herramientas teóricas de la microeconomía y los análisis se centran en cuestiones difíciles que suelen requerir formulaciones complejas
    - La macroeconometría se centra en el análisis de series temporales, normalmente de variables agregadas macroeconómicas

- Otra distinción útil es la distinción entre la econometría teórica y la econometría aplicada
  - Los econometristas teóricos desarrollan nuevas técnicas para la estimación y los contrastes de hipótesis y analizan las consecuencias de aplicar métodos particulares cuando las suposiciones que justifican esos métodos no se dan
  - Los econometristas aplicados son los usuarios de estas técnicas y loss analistas de datos reales y simulados
- Los datos económicos se suelen presentar en estructuras variadas, lo cual hace que a veces se tenga que tener en cuenta sus características con tal de poder hacer un análisis econométrico adecuado. Se trabaja principalmente con cuatro tipos de bases de datos: de datos transversales, de secciones transversales agrupadas, de series de tiempo y de datos longitudinales
  - Una base de datos transversales o *cross-sectional data* consiste en una muestra de individuos tomada en un punto en el tiempo

obsno	wage	educ	exper	female	married
1	3.10	11	2	1	0
2	3.24	12	22	1	1
3	3.00	11	2	0	0
4	6.00	8	44	0	1
5	5.30	12	7	0	1
.	.	.	.	.	.
.	.	.	.	.	.
.	.	.	.	.	.
525	11.56	16	5	0	1
526	3.50	14	5	1	0

- Una importante característica es que comúnmente se puede asumir que estos datos se han obtenido con muestreo aleatorio de la población, aunque a veces esta suposición no es posible
- Una base de datos de series temporales consiste en observaciones de una o más variables a lo largo del tiempo. Como el pasado, el presente y el futuro suelen estar conectados, el tiempo es un aspecto importante a tener en cuenta en este tipo de bases de datos

obsno	year	avgmin	avgcov	prunemp	prgnp
1	1950	0.20	20.1	15.4	878.7
2	1951	0.21	20.7	16.0	925.0
3	1952	0.23	22.6	14.8	1015.9
.	.	.	.	.	.
.	.	.	.	.	.
.	.	.	.	.	.
37	1986	3.35	58.1	18.9	4281.6
38	1987	3.35	58.2	16.8	4496.7

- Una dificultad que se encuentra es que no se puede asumir que las observaciones son independientes en el tiempo. Por ello se han desarrollado modificaciones a las técnicas econométricas estándar para tener en cuenta y explotar la naturaleza dependiente de las series de tiempo económicas y para abordar otras cuestiones, como el hecho de que algunas variables económicas tienden a mostrar tendencias claras a lo largo del tiempo
- Otro aspecto importante es la frecuencia con la que se recolectan los datos, dado que muchos datos económicos muestran patrones estacionales que tienen que tenerse en cuenta en el análisis
- Una base de datos de secciones transversales agrupadas es una base de datos que combina datos transversales y series temporales, de modo que se aumenta la muestra con observaciones de diferentes puntos en el tiempo

obsno	year	hprice	proptax	sqrft	bdrms	bthrms
1	1993	85500	42	1600	3	2.0
2	1993	67300	36	1440	3	2.5
3	1993	134000	38	2000	4	2.5
.	.	.	.	.	.	.
.	.	.	.	.	.	.
.	.	.	.	.	.	.
250	1993	243600	41	2600	4	3.0
251	1995	65000	16	1250	2	1.0
252	1995	182400	20	2200	4	2.0
253	1995	97500	15	1540	3	2.0
.	.	.	.	.	.	.
.	.	.	.	.	.	.
.	.	.	.	.	.	.
520	1995	57200	16	1100	2	1.5

- De este modo, los individuos de la muestra pueden ser observados en periodos diferentes, pero no se tiene una serie

temporal por cada individuo, si no que un individuo puede ser seleccionado en un punto del tiempo diferente a otro

- Se suele analizar como una base de datos transversales, pero se tiene en cuenta también la dimensión temporal
- Una base de datos de panel o datos longitudinales es una base de datos para la cual se tiene una serie temporal por cada individuo transversal

obsno	city	year	murders	population	unem	police
1	1	1986	5	350000	8.7	440
2	1	1990	8	359200	7.2	471
3	2	1986	2	64300	5.4	75
4	2	1990	1	65100	5.5	75
.	.	.	.	.	.	.
.	.	.	.	.	.	.
.	.	.	.	.	.	.
297	149	1986	10	260700	9.6	286
298	149	1990	6	245000	9.8	334
299	150	1986	25	543000	4.3	520
300	150	1990	32	546200	5.2	493

- La diferencia entre esta base y la anterior es que para los mismos individuos de la muestra se tienen observaciones de diferentes momentos en el tiempo
- El objetivo principal de hacer un análisis econométrico empírico para comprobar la teoría económica y para evaluar políticas públicas es inferenciar que una variable tiene efecto causal sobre otra
  - La noción de *Ceteris Paribus* (manteniendo todos los otros factores relevantes constantes) juega un papel importante en el análisis causal, dado que la mayoría de cuestiones económicas son por naturaleza *Ceteris Paribus*
    - Si todos los otros factores no se mantienen fijos, no se puede saber el efecto causal de una variable, dado que las otras pueden influir en el resultado
    - Casi nunca es posible mantener absolutamente todos los otros factores constantes, por lo que lo importante suele ser saber si se mantienen fijos suficientes factores como para inferenciar el efecto causal
  - Existen muchas dificultades para poder inferenciar el efecto causal, debido a que se tratan con datos no experimentales u observacionales

## El modelo de regresión lineal

- La econometría se centra principalmente en la creación de modelos, por lo tanto, es importante entender cómo se suelen desarrollar estos modelos para poder establecer un marco práctico
  - Un modelo suele comenzar por la observación o proposición de que una variable es causada por o varía junto a otra, o por una proposición cualitativa sobre la relación entre una variable y una o más covariables que se espera que estén relacionadas a la variable interesante en cuestión
    - El modelo suele hacer una afirmación general sobre el comportamiento de los individuos
    - De este modo, un modelo econométrico no nace como un conjunto de ecuaciones, sino como una idea de algún tipo de relación que se tiene que traducir a un conjunto de ecuaciones que permitan responder preguntas importantes de la variable de interés
  - Desde un punto de vista estadístico, se tiene en mente una variable  $y$  y un vector de covariables  $x$  junto con una distribución de probabilidad conjunta entre ellas  $p(y, x)$ 
    - Por lo tanto, la relación entre variables sería el proceso estadístico por el cual se producen las variables
    - Como  $x$  es un vector de varias variables aleatorias, se puede descomponer la distribución conjunta en la distribución condicional de  $y$  sobre  $x$  y en la distribución conjunta de  $x$ 
$$p(y, x) = p(y|x)p(x)$$
    - Normalmente uno se interesa en la variación condicional  $p(y|x)$  de una de las variables en relación a otras más que en la variación conjunta  $p(y, x)$  de todas las variables
  - En consecuencia, la idea de una distribución conjunta proporciona un buen comienzo para pensar sobre la relación entre la variable de interés  $y$  y las covariables  $x$ 
    - Pensando en términos de la distribución condicional, uno se suele centrar en el su valor esperado  $E(y|x)$ , llamado función de regresión

- Aunque hay otros estadísticos útiles, es necesario establecer como analizar la función de regresión para poder analizar los otros estadísticos
- Aunque existen herramientas más complejas, la herramienta más simple y útil en econometría es el modelo de regresión lineal, dado que suele ser el modelo inicial para la investigación y permite establecer un marco de referencia para las relaciones entre variables
  - El modelo de regresión lineal múltiple se utiliza para estudiar la relación entre una variable dependiente o explicada  $y$  y una o más variables independientes o explicativas  $x_1, x_2, \dots, x_K$ , cuya forma genérica es la siguiente:

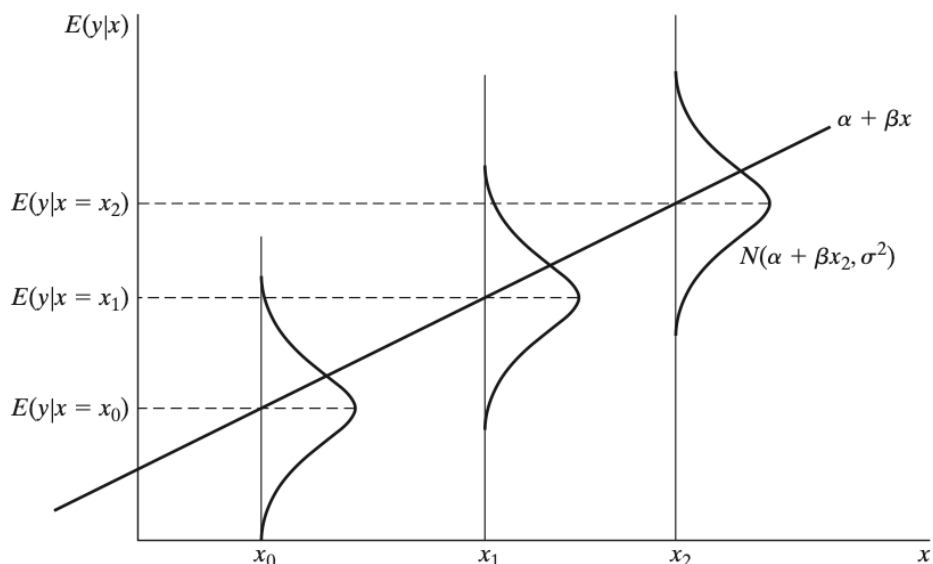
$$y = f(x_1, x_2, \dots, x_K) + \varepsilon = x_1\beta_1 + x_2\beta_2 + \dots + x_K\beta_K + \varepsilon$$

- La teoría económica especifica  $f(x_1, x_2, \dots, x_K)$ , la cual se conoce como la ecuación de regresión poblacional de  $y$  sobre  $x_1, x_2, \dots, x_K$
- Uno de los aspectos más útiles de este modelo es que permite identificar los efectos independientes de un conjunto de variables sobre la variable dependiente
- En este modelo,  $y$  es la variable regresora y  $x_1, x_2, \dots, x_K$  son los regresores o las covariables. Estas variables también las suele especificar la teoría económica
  - Para propósitos de modelaje, normalmente es útil pensar en términos de variación autónoma, de modo que se concibe el movimiento de las variables independientes fuera de las relaciones definidas por el modelo mientras que el movimiento de las variables independientes se considera respuesta a algún estímulo independiente o exógeno
- El término  $\varepsilon$  es una perturbación aleatoria o error (porque representa perturbaciones a la relación estable planteada), el cual nace de no poder capturar todas las influencias en una variable económica del modelo
  - El efecto neto de estos factores omitidos puede ser positivo o negativo y se captura en este término de error
  - No obstante, este término de error también se puede deber a errores de medición u otros problemas



- Se asume que cada observación en una muestra  $(y_i, x_{i1}, x_{i2}, \dots, x_{iK})$  para  $i = 1, 2, \dots, n$  se genera por un proceso subyacente descrito por la siguiente ecuación:

$$y_i = x_{i1}\beta_1 + x_{i2}\beta_2 + \dots + x_{iK}\beta_K + \varepsilon_i$$



- El valor observado de  $y_i$  es la suma de dos partes: una parte determinística y una parte aleatoria  $\varepsilon_i$
- De este modo, el objetivo es estimar los parámetros desconocidos del modelo, usar los datos para estudiar la validez de las proposiciones teóricas y posiblemente utilizar el modelo para predecir  $y$
- Siendo el vector columna  $\mathbf{x}_k$  que contiene las  $n$  observaciones en la variable  $x_k$  para  $k = 1, 2, \dots, K$ ,  $\mathbf{y}$  el vector columna que contiene las  $n$  observaciones  $y_1, y_2, \dots, y_K$ ,  $\boldsymbol{\varepsilon}$  el vector columna que contiene las  $n$  perturbaciones y  $\boldsymbol{\beta}$  el vector columna de parámetros, es posible reformular el modelo de manera más conveniente
- Uniendo estos datos en una matriz de datos  $\mathbf{X}$  de tamaño  $n \times K$  y asumiendo que la primera columna de  $\mathbf{X}$  es una columna de unos para que  $\beta_1$  sea el término constante del modelo, se puede reformular el modelo de regresión lineal múltiple de la siguiente manera:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

- Para evitar confusiones de notación, se usa  $\mathbf{x}$  para denotar una columna de  $\mathbf{X}$ , de modo que  $\mathbf{x}_k$  es la columna  $k$  de  $\mathbf{X}$ , y los suscritos  $j$  y  $k$  se usan para denotar columnas (variables),

mientras que  $i$  y  $t$  se utilizan para denotar filas (observaciones). Por lo tanto, se puede reescribir el proceso para un individuo:

$$\mathbf{x}'_i = (x_{1i}, x_{2i}, \dots, x_{Ki}) \Rightarrow y_i = \mathbf{x}'_i \boldsymbol{\beta} + \varepsilon_i$$

- En este caso, la matriz  $\mathbf{X}$  se puede expresar como un vector fila en el que cada componente es  $\mathbf{x}_j$  (el vector columna de observaciones para la variable  $j$ ), dado que es equivalente a una matriz adjunta y contiene los mismos datos. También se suele utilizar  $\mathbf{x}$  para denotar el vector de variables aleatorias

$$\mathbf{X} = [\mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_K] = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K)$$

- El interés principal reside en estimar e inferenciar sobre el vector de parámetros  $\boldsymbol{\beta}$
- El modelo de regresión lineal consiste en un conjunto de suposiciones sobre como la base de datos ha sido producida y el proceso generador de datos subyacente
  - Las suposiciones describen la forma del modelo y las relaciones entre las partes, además de sugerir unos procedimientos de estimación e inferencia apropiados. Estas suposiciones son las siguientes:

**A1. Linearity:**  $y_i = x_{i1}\beta_1 + x_{i2}\beta_2 + \dots + x_{iK}\beta_K + \varepsilon_i$ . The model specifies a linear relationship between  $y$  and  $x_1, \dots, x_K$ .

**A2. Full rank:** There is no exact linear relationship among any of the independent variables in the model. This assumption will be necessary for estimation of the parameters of the model.

**A3. Exogeneity of the independent variables:**  $E[\varepsilon_i | x_{j1}, x_{j2}, \dots, x_{jK}] = 0$ . This states that the expected value of the disturbance at observation  $i$  in the sample is not a function of the independent variables observed at any observation, including this one. This means that the independent variables will not carry useful information for prediction of  $\varepsilon_i$ .

**A4. Homoscedasticity and nonautocorrelation:** Each disturbance,  $\varepsilon_i$  has the same finite variance,  $\sigma^2$ , and is uncorrelated with every other disturbance,  $\varepsilon_j$ . This assumption limits the generality of the model, and we will want to examine how to relax it in the chapters to follow.

**A5. Data generation:** The data in  $(x_{j1}, x_{j2}, \dots, x_{jK})$  may be any mixture of constants and random variables. The crucial elements for present purposes are the strict mean independence assumption A3 and the implicit variance independence assumption in A4. Analysis will be done conditionally on the observed  $\mathbf{X}$ , so whether the elements in  $\mathbf{X}$  are fixed constants or random draws from a stochastic process will not influence the results. In later, more advanced treatments, we will want to be more specific about the possible relationship between  $\varepsilon_i$  and  $\mathbf{x}_j$ .

**A6. Normal distribution:** The disturbances are normally distributed. Once again, this is a convenience that we will dispense with after some analysis of its implications.

- Estas son suposiciones que se suelen utilizar en el modelo de regresión lineal, pero no quiere decir que sean todas necesarias para derivar algunos resultados (como el teorema de Gauss-Markov)
- Las primeras tres suposiciones conforman el modelo de regresión lineal, mientras que las otras especifican las características de las perturbaciones del modelo y las

condiciones para las cuales las observaciones muestrales de  $\mathbf{x}$  se pueden obtener

- La suposición de linealidad del modelo de regresión incluye un término de error aditivo, de modo que para que un modelo sea lineal en este sentido, es necesario que sea de la forma vista anteriormente con las variables originales o con una transformación adecuada
  - En el contexto de regresión, la linealidad se refiere a la manera en que los parámetros y el término de error entran en la ecuación, y no necesariamente a la relación entre las variables
  - Variaciones de la forma general  $f(\mathbf{y}_i) = g(\mathbf{x}_i'\boldsymbol{\beta} + \varepsilon_i)$  permiten crear una gran variedad de formas funcionales que respetan la definición de linealidad
  - Aunque el modelo lineal se puede interpretar como una aproximación para una función subyacente desconocida, es una aproximación muy limitada
  - Aunque el modelo lineal sea muy flexible, no siempre es adecuado para todas las situaciones y no se pueden estimar parámetros ni inferenciar como en este modelo
- La suposición de rango completo o *full rank* establece que no hay relaciones lineales exactas entre variables, también conocida como suposición de no colinearidad perfecta
  - La matriz  $\mathbf{X}$  es una matriz  $n \times K$  con rango  $K$ , de modo que  $\mathbf{X}$  tiene un rango del espacio de columna (y de fila) igual a  $K$  y eso quiere decir que las columnas de  $\mathbf{X}$  son linealmente independientes y que hay como mínimo  $K$  observaciones (dado que  $n$  debe de ser mayor o igual a  $K$ )
  - Esta suposición se conoce como condición de identificación, y es necesaria dado que, de no haber columnas independientes, no se podría estimar los parámetros porque las variables no pueden variar de manera linealmente independiente (variar cuando las otras se quedan constantes) y no se podría inferenciar
  - Hay veces que se pueden utilizar formulaciones que incluyan una misma variable en diferentes formas funcionales (por ejemplo, utilizar cuadrados), pero eso no afecta a la independencia lineal porque aún no siendo funcionalmente independientes, sí que lo son linealmente

- La suposición de exogeneidad de las variables independientes o de independencia de la media expresa que el término de error tiene un valor esperado condicional de cero para cada observación

$$E(\varepsilon_i|\mathbf{X}) = 0 \quad E(\boldsymbol{\varepsilon}|\mathbf{X}) = \begin{bmatrix} E(\varepsilon_1|\mathbf{X}) \\ E(\varepsilon_2|\mathbf{X}) \\ \dots \\ E(\varepsilon_n|\mathbf{X}) \end{bmatrix} = \mathbf{0}$$

- En este caso, la media de cada  $\varepsilon_i$  condicional a todas las observaciones  $\mathbf{x}_i$  es nula, de modo que ninguna observación en  $\mathbf{x}$  proporciona información sobre el valor esperado de las perturbaciones o no explica esta variación (ni  $\mathbf{x}_i$  ni  $\mathbf{x}_j$ )
- Hay situaciones en las que, aunque  $\mathbf{x}_i$  no proporcione ninguna información sobre la esperanza condicional del error,  $\mathbf{x}_j$  puede proporcionar información en otra observación. Este suele ser el caso con series temporales, en donde las observaciones son los periodos y una observación puede de  $\mathbf{x}_j$  puede proporcionar información sobre el siguiente periodo
- Esta suposición implica que la esperanza incondicional de la distribución incondicional del término de error también es nula gracias a las propiedades del valor esperado. No obstante, no ocurre lo converso

$$E(\varepsilon_i) = E_X[E(\varepsilon_i|\mathbf{X})] = E_X(0) = 0$$

$$E(\varepsilon_i|\mathbf{X}) = 0 \Rightarrow E(\varepsilon_i) = 0 \text{ but not } E(\varepsilon_i) = 0 \Rightarrow E(\varepsilon_i|\mathbf{X}) = 0$$

- La suposición no suele ser restrictiva, pero asumir  $E(\varepsilon_i|\mathbf{X}) = 0$  en modelos sin término constante puede alterar el modelo y por eso se asume que los modelos de regresión tienen término constante a no ser que se especifique lo contrario
- Debido a que para cada  $\varepsilon_i$ , la covarianza entre  $E(\varepsilon_i|\mathbf{X})$  y  $\mathbf{X}$  es lo mismo que  $Cov(\varepsilon_i, \mathbf{X})$ , por lo que la suposición también implica que esta es nula para cualquier  $i$

$$\begin{aligned} Cov(\varepsilon_i, \mathbf{X}) &= Cov[E(\varepsilon_i|\mathbf{X}), \mathbf{X}] = \\ &= E \left[ \left( E(\varepsilon_i|\mathbf{X}) - E_X(E(\varepsilon_i|\mathbf{X})) \right) (\mathbf{X} - E(\mathbf{X})) \right] = \\ &= E \left[ \left( E(\varepsilon_i|\mathbf{X}) - E_X(\varepsilon_i) \right) (\mathbf{X} - E(\mathbf{X})) \right] = E[0(\mathbf{X} - E(\mathbf{X}))] = 0 \end{aligned}$$

- Finalmente, la suposición también implica que la función de regresión  $E(\mathbf{y}|\mathbf{X})$  es  $\mathbf{X}\boldsymbol{\beta}$ , de modo que la regresión de  $\mathbf{y}$  en  $\mathbf{X}$  es la esperanza condicional

$$E(\mathbf{y}|\mathbf{X}) = \mathbf{X}\boldsymbol{\beta}$$

- La cuarta suposición expresa que la varianza condicional del error con respecto a  $\mathbf{X}$  es constante y las perturbaciones no están correlacionadas, por lo que la suposición se denomina perturbaciones esféricas

$$Var(\varepsilon_i|\mathbf{X}) = \sigma^2 \text{ for all } i = 1, 2, \dots, n$$

$$Cov(\varepsilon_i, \varepsilon_j|\mathbf{X}) = 0 \text{ for all } i \neq j$$

- La primera parte de la suposición se refiere a que la varianza es constante para todas las observaciones, propiedad conocida como homoscedasticidad. Si la varianza depende de  $\mathbf{X}$  (o, equivalentemente, varía dependiendo de la observación), entonces la propiedad es la heterocedasticidad
- La segunda parte de la suposición se refiere a que no hay correlación entre perturbaciones, propiedad conocida como no autocorrelación. Cuando las observaciones siguen una inercia o patrón (hay dependencia unas de otras) se dice que hay autocorrelación
- No obstante, esta segunda parte no implica que no hay correlación entre  $y_i$  y  $y_j$ , sino que las desviaciones de las observaciones con respecto a sus valores esperados no están correlacionadas (los términos de error)
- La no autocorrelación se puede expresar en términos de esperanzas condicionales como la esperanza condicional del producto de las perturbaciones. A partir de esto, es posible expresar la homoscedasticidad en términos de la covarianza y de la esperanza condicional

$$\begin{aligned} Cov(\varepsilon_i, \varepsilon_j|\mathbf{X}) &= E[(\varepsilon_i - E(\varepsilon_i))(\varepsilon_j - E(\varepsilon_j))|\mathbf{X}] = \\ &= E[(E(\varepsilon_i|\mathbf{X}) - E(\varepsilon_i))(E(\varepsilon_j|\mathbf{X}) - E(\varepsilon_j))|\mathbf{X}] = E(\varepsilon_i \varepsilon_j|\mathbf{X}) = 0 \\ &\text{for all } i \neq j \end{aligned}$$

$$\begin{aligned} Var(\varepsilon_i|\mathbf{X}) &= Cov(\varepsilon_i, \varepsilon_i|\mathbf{X}) = E(\varepsilon_i \varepsilon_i|\mathbf{X}) = E(\varepsilon_i^2|\mathbf{X}) = \sigma^2 \\ &\text{for all } i = 1, 2, \dots, n \end{aligned}$$

$$Var(\varepsilon_i | \mathbf{X}) = E(\varepsilon_i^2 | \mathbf{X}) - [E(\varepsilon_i | \mathbf{X})]^2 = E(\varepsilon_i^2 | \mathbf{X}) = \sigma^2$$

for all  $i = 1, 2, \dots, n$

- Las dos suposiciones juntas implican que  $E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}' | \mathbf{X}) = \sigma^2 \mathbf{I}$ , lo cual se puede expresar de la siguiente forma:

$$\begin{aligned} E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}' | \mathbf{X}] &= \begin{bmatrix} E[\varepsilon_1 \varepsilon_1 | \mathbf{X}] & E[\varepsilon_1 \varepsilon_2 | \mathbf{X}] & \cdots & E[\varepsilon_1 \varepsilon_n | \mathbf{X}] \\ E[\varepsilon_2 \varepsilon_1 | \mathbf{X}] & E[\varepsilon_2 \varepsilon_2 | \mathbf{X}] & \cdots & E[\varepsilon_2 \varepsilon_n | \mathbf{X}] \\ \vdots & \vdots & \ddots & \vdots \\ E[\varepsilon_n \varepsilon_1 | \mathbf{X}] & E[\varepsilon_n \varepsilon_2 | \mathbf{X}] & \cdots & E[\varepsilon_n \varepsilon_n | \mathbf{X}] \end{bmatrix} \\ &= \begin{bmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ & & \ddots & \\ 0 & 0 & \cdots & \sigma^2 \end{bmatrix}, \end{aligned}$$

- Usando la descomposición de la varianza, se puede ver como la varianza incondicional del término de error es la matriz  $\sigma^2 \mathbf{I}$ , lo cual se puede interpretar como la suposición que junta ambas condiciones anteriores (ya que se necesita que no haya otros términos más que la varianza de los errores en la diagonal y eso ocurre cuando no hay correlación entre errores)

$$Var(\boldsymbol{\varepsilon}) = E[Var(\boldsymbol{\varepsilon} | \mathbf{X})] + Var[E(\boldsymbol{\varepsilon} | \mathbf{X})] = \sigma^2 \mathbf{I} + 0 = \sigma^2 \mathbf{I}$$

- La suposición de generación de datos se refiere a que  $\mathbf{X}$  puede ser fija (no estocástica) o aleatoria (estocástica)
  - Es común suponer que  $\mathbf{x}_i$  no es estocástica, como en un experimento en donde el analista escoge los valores de los regresores y observa  $y_i$
  - Esta suposición es una conveniencia matemática porque se trataría el vector  $\mathbf{x}_i$  como una constante en la distribución de probabilidad de  $\mathbf{x}_i$  y usar métodos estadísticos básicos para obtener los resultados y encontrar propiedades de los estimadores. Además, hace que la tercera y la cuarta suposición sean incondicionales (las esperanzas sean incondicionales)
  - En aplicaciones, lo más común es que los datos en  $\mathbf{x}_i$  sean aleatorios como  $y_i$ . Si  $\mathbf{x}_i$  es un vector aleatorio, entonces las suposiciones anteriores se transforman en proposiciones sobre la distribución conjunta de  $y_i$  y  $\mathbf{x}_i$

- La naturaleza del regresor y como se entiende el proceso de muestreo son determinantes para la derivación de las propiedades estadísticas de los estimadores y tests de hipótesis
- Sin embargo,  $X$  puede contener elementos estocásticos y no estocásticos (variables binarias, constantes, tendencias, etc.) a la vez, por lo que se asume que  $X$  puede contener ambos y no solo un tipo
- La suposición de normalidad se refiere que las perturbaciones están normalmente distribuidas con una media nula y con una varianza constante

$$\varepsilon | X \sim N(\mathbf{0}, \sigma^2 I)$$

- En muchos casos, las condiciones del teorema central del límite se darán (aunque sea aproximadamente) y la suposición será razonable
- Esta suposición implica que las observaciones de las perturbaciones  $\varepsilon_i$  son estadísticamente independientes y no están correlacionadas
- Esta suposición no es necesaria para obtener resultados en el modelo de regresión múltiple, pero permite obtener resultados estadísticos más exactos
- Esta suposición implica que hay independencia estadística entre las perturbaciones, que es un resultado más fuerte que el de independencia de media

## Los mínimos cuadrados

- Existen varios métodos para poder estimar los parámetros del modelo, pero el método más popular es el de mínimos cuadrados debido a razones teóricas y prácticas. Este sirve como referencia y otros métodos suelen basarse en modificaciones de este, de modo que es importante estudiar la regresión de mínimos cuadrados
  - Los parámetros desconocidos en la relación estocástica  $y_i = x_i' \beta + \varepsilon_i$  son los objetos de estimación, pero es necesario distinguir entre las cantidades poblacionales, como  $\beta$  y  $\varepsilon_i$ , de las estimaciones muestrales, denotadas por  $b$  y  $e_i$ 
    - De este modo, se puede hacer una distinción entre la regresión poblacional  $E(y_i | x_i) = x_i' \beta$  y la regresión estimada  $\hat{y}_i = x_i' b$

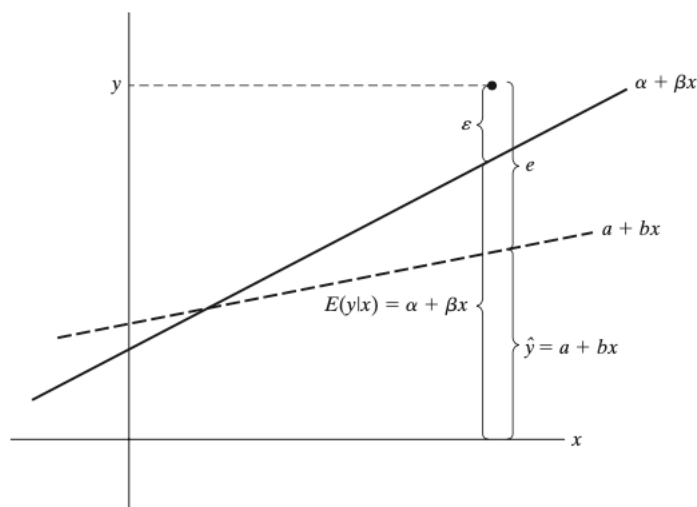
$$E(y_i|x_i) = x_i'\beta \quad \hat{y}_i = x_i'b$$

- En consecuencia, también se puede hacer una distinción entre la perturbación poblacional asociada a la observación  $i$  y la estimación de esta, llamada residuo

$$\varepsilon_i = y_i - x_i'\beta \quad e_i = y_i - x_i'b$$

- A partir de estas definiciones, se puede ver que se cumplen las siguientes identidades para  $y_i$ :

$$y_i = x_i'\beta + \varepsilon_i = x_i'b + e_i$$



- Las cantidades poblacionales  $\beta$  es un vector de parámetros desconocidos de la distribución de probabilidad de  $y_i$  cuyos valores se intentan estimar a través de los datos  $(y_i, x_i)$  para  $i = 1, 2, \dots, n$
- Debido a que se intenta estimar  $\beta$  a través de los datos, este es un problema de inferencia estadística, pero se comienza considerando un problema algebraico en el que se tiene que escoger  $b$  de modo que  $x_i'b$  esté lo más cerca posible a los datos (el criterio de ajuste es la suma de residuos cuadrados), y el candidato más usado es el vector de coeficientes de mínimos cuadrados
- El vector de coeficientes de mínimos cuadrados minimiza la suma de residuos cuadrados, donde  $b_0$  denota la elección del vector de coeficientes

$$\sum_{i=1}^n e_{i0}^2 = \sum_{i=1}^n (y_i - x_i'b_0)^2$$



- En términos matriciales, minimizar la suma de cuadrados requiere escoger  $\mathbf{b}_0$  tal que se minimice la suma de cuadrados

$$\min_{\mathbf{b}_0} S(\mathbf{b}_0) = \mathbf{e}_0' \mathbf{e}_0 = (\mathbf{y} - \mathbf{X}\mathbf{b}_0)'(\mathbf{y} - \mathbf{X}\mathbf{b}_0)$$

$$\begin{aligned} (\mathbf{y} - \mathbf{X}\mathbf{b}_0)'(\mathbf{y} - \mathbf{X}\mathbf{b}_0) &= \mathbf{y}'\mathbf{y} - \mathbf{b}_0'\mathbf{X}'\mathbf{y} - \mathbf{y}'\mathbf{X}\mathbf{b}_0 - \mathbf{b}_0'\mathbf{X}'\mathbf{X}\mathbf{b}_0 = \\ &= \mathbf{y}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}\mathbf{b}_0 - \mathbf{b}_0'\mathbf{X}'\mathbf{X}\mathbf{b}_0 \end{aligned}$$

$$\Rightarrow \min_{\mathbf{b}_0} S(\mathbf{b}_0) = \mathbf{e}_0' \mathbf{e}_0 = \mathbf{y}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}\mathbf{b}_0 - \mathbf{b}_0'\mathbf{X}'\mathbf{X}\mathbf{b}_0$$

- Las condiciones de optimalidad permiten ver como  $\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ , siendo  $\mathbf{b}$  la solución

$$\frac{\partial S(\mathbf{b}_0)}{\partial \mathbf{b}_0} = -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\mathbf{b}_0 = \mathbf{0} \Rightarrow \mathbf{X}'\mathbf{X}\mathbf{b}_0 = \mathbf{X}'\mathbf{y}$$

$$\Rightarrow \mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

$$\frac{\partial S(\mathbf{b}_0)}{\partial \mathbf{b}_0 \partial \mathbf{b}_0'} = 2\mathbf{X}'\mathbf{X} > \mathbf{0} \text{ to be a minimum}$$

$$\Rightarrow \mathbf{c}'\mathbf{X}'\mathbf{X}\mathbf{c} = \sum_{i=1}^n (\mathbf{X}_i \mathbf{c}_i)^2 > 0 \text{ if } \mathbf{c} \neq \mathbf{0} \Rightarrow 2\mathbf{X}'\mathbf{X} > \mathbf{0}$$

- Debido a que  $\mathbf{X}$  tiene rango  $K$  (de modo que  $\mathbf{X}\mathbf{c} \neq \mathbf{0}$  porque no puede ser una combinación lineal de columnas de  $\mathbf{X}$  igual a  $\mathbf{0}$ ), la solución de mínimos cuadrados  $\mathbf{b}$  es única y minimiza la suma de residuos cuadrados
- A partir de la solución, se puede ver que  $\mathbf{X}'\mathbf{X}\mathbf{b} - \mathbf{X}'\mathbf{y} = \mathbf{0}$ , entonces para cada columna  $\mathbf{x}_k$  de  $\mathbf{X}$ ,  $\mathbf{x}_k'\mathbf{e} = 0$ . Esto permite desarrollar tres propiedades algebraicas de la solución  $\mathbf{b}$  siempre que  $\mathbf{X}$  tenga una columna de unos, denotada por  $\mathbf{i}$  (si no, no siempre se mantienen)

$$\mathbf{X}'\mathbf{X}\mathbf{b} - \mathbf{X}'\mathbf{y} = \mathbf{0} \Rightarrow -\mathbf{X}'(\mathbf{y} - \mathbf{X}\mathbf{b}) = \mathbf{0} \Rightarrow -\mathbf{X}'\mathbf{e} = \mathbf{0} \Rightarrow \mathbf{x}_k'\mathbf{e} = 0$$

- La suma de residuos es nula, dado que, si  $\mathbf{X}$  tiene una columna de unos  $\mathbf{i}$ , entonces  $\mathbf{i}'\mathbf{e} = 0$

$$\mathbf{i}'\mathbf{e} = \sum_{i=1}^n e_i = 0$$

- El hiperplano de la regresión pasa por el punto de medias de los datos, dado que  $X'Xb - X'y = 0$  implica que  $\bar{y} = \bar{x}'b$

$$-X'(y - Xb) = 0 \Rightarrow y - Xb = 0 \Rightarrow 1'(y - Xb) = 1'0$$

$$\Rightarrow 1'y - 1'Xb = 0 \Rightarrow \left(\frac{1}{n}\right)[1'y - 1'Xb] = \left(\frac{1}{n}\right)0$$

$$\Rightarrow \left(\frac{1}{n}\right)1'y - \left(\frac{1}{n}\right)1'Xb = 0 \Rightarrow \bar{y} = \bar{x}'b$$

- La media de los valores ajustados de la regresión es equivalente a la media de los valores actuales de  $y_i$ , dado que  $\hat{y} = y + e$

$$\bar{\hat{y}} = \left(\frac{1}{n}\right)1'\hat{y} = \left(\frac{1}{n}\right)1'y + \left(\frac{1}{n}\right)1'e = \left(\frac{1}{n}\right)1'y = \bar{y}$$

- El vector de residuos mínimos cuadrados es  $e = y - Xb$ , así que expandiendo  $b$  se puede ver que  $e$  se puede expresar de la siguiente manera:

$$e = y - X(X'X)^{-1}X'y = (I - X(X'X)^{-1}X')y = My$$

- La matriz  $M \equiv I - X(X'X)^{-1}X'$  es una matriz  $n \times n$  es fundamental en el análisis de regresiones y es una matriz simétrica e idempotente ( $M = M^2$ ). Esta se puede interpretar como la matriz que produce el vector de residuos mínimos cuadrados en la regresión de  $y$  sobre  $X$  cuando se multiplica a la izquierda de  $y$  (llamada fabricante de residuos)
- Cuando se hace una regresión de  $y$  sobre  $X$ , habrá un ajuste perfecto y los residuales serán nulos, de modo que se da la siguiente igualdad:

$$MX = 0$$

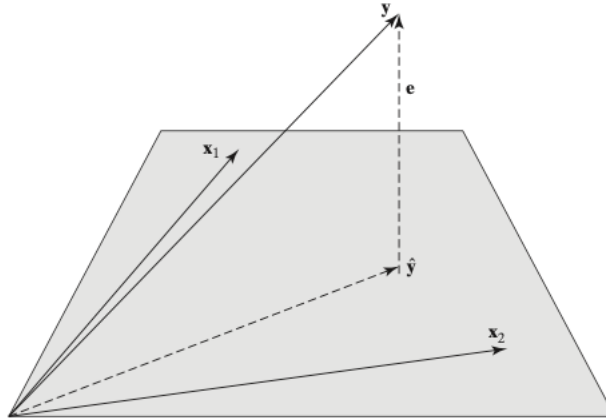
- Como  $e = y - Xb$ , se puede ver que  $y = Xb + e$ , por lo que los resultados de mínimos cuadrados hagan una partición de  $y$  en dos partes: los valores ajustados  $\hat{y} = Xb$  y en los residuos  $e$ , y de este modo se puede derivar la siguiente igualdad:

$$\hat{y} = y - e = (I - M)y = (I - I + X(X'X)^{-1}X')y = (X(X'X)^{-1}X')y = Py$$

$$\text{where } P \equiv X(X'X)^{-1}X'$$

- Debido a que  $MX = 0$ ,  $\hat{y}$  y  $e$  son ortogonales

- Al igual que  $M$ ,  $P$  tiene las propiedades de simetría e idempotencia. Esta matriz  $P$  es una matriz de proyección, la cual es formada a partir de  $X$  de manera que cuando un vector  $P$  se multiplica a la izquierda de  $y$ , el resultado son los valores ajustados en la regresión de mínimos cuadrados, lo cual es una proyección del vector  $y$  en el espacio de columna de  $X$  (denotado por  $W$ )



- Al multiplicar las matrices  $M$  y  $P$ , se puede ver que estas son ortogonales dada a las propiedades de idempotencia y simetría

$$\begin{aligned}
 PM &= MP = (X(X'X)^{-1}X')(I - X(X'X)^{-1}X') = \\
 &= X(X'X)^{-1}X' - X(X'X)^{-1}X' = 0 \\
 PM &= MP = 0
 \end{aligned}$$

- Cuando se hace una regresión de  $X$  sobre  $X$ , la proyección será la misma  $X$  (de manera lógica)

$$PX = X$$

- Las siguientes identidades son útiles al resolver problemas que involucran resultados de mínimos cuadrados (basadas en la propiedad de idempotencia y simetría de las matrices):

$$e'e = y'M'My = y'My = y'e = e'y$$

$$e'e = y'y - b'X'Xb = y'y - b'X'y = y'y - y'Xb$$

- Es común especificar una regresión múltiple cuando solo es interesante un subconjunto del conjunto completo de variables, de modo que es necesario saber que cálculos se involucran para obtener los coeficientes de este subconjunto de variables

- Suponiendo que la regresión involucra dos conjuntos de variables  $X_1$  y  $X_2$ , se puede expresar la ecuación de  $y$  de la siguiente manera:

$$y = X\beta + \varepsilon = X_1\beta_1 + X_2\beta_2 + \varepsilon$$

- Las ecuaciones normales permiten obtener una solución para  $b_1$  y  $b_2$  a través de un sistema de ecuaciones:

$$\begin{bmatrix} X_1'X_1 & X_1'X_2 \\ X_2'X_1 & X_2'X_2 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} X_1'y \\ X_2'y \end{bmatrix} \text{ as } \begin{cases} X_1'X_1b_1 + X_1'X_2b_2 = X_1'y \\ X_2'X_1b_1 + X_2'X_2b_2 = X_2'y \end{cases}$$

$$\Rightarrow b_1 = (X_1'X_1)^{-1}X_1'y - (X_1'X_1)^{-1}X_1'X_2b_2$$

$$\Rightarrow b_2 = (X_2'X_2)^{-1}X_2'y - (X_2'X_2)^{-1}X_2'X_1b_1$$

- Como se puede ver,  $b_1$  y  $b_2$  es el conjunto de coeficientes en la regresión de  $y$  en  $X$  menos un vector de corrección
- Cuando  $X_1'X_2 = 0$  o  $X_2'X_1 = 0$ , entonces  $b_1$  y  $b_2$  corresponden a los vectores de coeficientes de las regresiones por separado, lo cual resulta en el siguiente teorema:

- En la regresión lineal múltiple de mínimos cuadrados de  $y$  en dos conjuntos  $X_1$  y  $X_2$ , si los dos conjuntos son ortogonales, entonces los vectores de coeficientes separados se pueden obtener de las regresiones separadas de  $y$  sobre  $X_1$  e  $y$  sobre  $X_2$
- Si los dos conjuntos son ortogonales, entonces  $X_1'X_2 = X_2'X_1 = 0$ , de modo que  $b_1 = (X_1'X_1)^{-1}X_1'y$  y  $b_2 = (X_2'X_2)^{-1}X_2'y$ , que son los coeficientes que se obtendrían en las regresiones separadas

- Si  $X_1$  y  $X_2$  no son ortogonales, entonces la solución de  $b_1$  y  $b_2$  se pueden obtener a través del teorema de Frisch-Waugh-Lovell

- En la regresión lineal de mínimos cuadrados del vector  $y$  sobre dos conjuntos de variables  $X_1$  y  $X_2$ , el subvector  $b_2$  es el conjunto de coeficientes obtenido cuando se hace una regresión de los residuos de una regresión de  $y$  sobre  $X_1$  sobre el conjunto de los residuos obtenidos de la regresión de cada columna de  $X_2$  sobre  $X_1$
- A partir de la segunda ecuación del sistema de ecuaciones anterior, se puede sustituir  $b_1$  y aislar  $b_2$  para poder obtener una expresión de  $b_2$  en términos de  $M_1X_2$  (matriz de residuos de la regresión de cada columna de  $X_2$  sobre  $X_1$ )

$$X_2'X_1b_1 + X_2'X_2b_2 = X_2'y$$

$$X_2'X_1(X_1'X_1)^{-1}X_1'y - X_2'X_1(X_1'X_1)^{-1}X_1'X_2b_2 + X_2'X_2b_2 = X_2'y$$

$$\Rightarrow b_2 = [X_2'(I - X_1(X_1'X_1)^{-1}X_1')X_2]^{-1}[X_2'(I - X_1(X_1'X_1)^{-1}X_1')y]$$

$$= (X_2'M_1X_2)^{-1}(X_2'M_1y)$$

$$\Rightarrow b_2 = (X_2^{*'}X_2^*)^{-1}X_2^*y^* \text{ where } X_2^* \equiv M_1X_2 \text{ \& } y^* \equiv M_1y$$

- Este proceso normalmente se denomina comúnmente eliminación parcial o compensación del efecto de  $X_1$ , por lo que los coeficientes en una regresión múltiple normalmente se denominan coeficientes de regresión parcial
- El teorema de Frisch-Waugh-Lovell permite derivar diferentes corolarios que permiten estimar coeficientes de regresores individuales, los coeficientes al incluir un término constante y otros
  - El coeficiente de  $z$  en una regresión múltiple de  $y$  sobre  $W = [X|z]$  (sobre dos grupos  $X_1 \equiv X$  y  $X_2 \equiv z$ , de modo que hay  $K$  columnas y una columna adicional  $z$ ) se calcula como  $c = (z'Mz)^{-1}(z'My) = (z^{*'}z^*)^{-1}z^*y^*$  donde  $M \equiv I - X(X'X)^{-1}X'$ ,  $z^* \equiv Mz$  e  $y^* \equiv My$ 
    - De este modo, al añadir un regresor al modelo de regresión lineal múltiple, es posible obtener el coeficiente solo con la matriz  $M$  y las variables  $z$  e  $y$ , dado que se obtendrán los vectores de residuos de las regresiones descritas anteriormente
  - Los coeficientes en una regresión múltiple que contiene un término constante se pueden obtener transformando los datos a desviaciones de sus medias y haciendo una regresión de la desviación de  $y$  sobre las desviaciones de la media muestral de los regresores
    - Considerando el caso en que  $X_1 \equiv i$ , la solución para  $b_2$  en este caso serán los coeficientes de los regresores que no son constantes. El vector de residuos para cualquier variable en  $X_2$  será la desviación de la media muestral de esta

$$x^* = (I - X_1(X_1'X_1)^{-1}X_1')x = x - i(i'i)^{-1}i'x =$$

$$= x - i\left(\frac{1}{n}\right)i'x = x - i\bar{x} = (I - i\bar{x})x = M^0x$$

$$M^0 \equiv I - i\left(\frac{1}{n}\right)i' = I - i\bar{x}$$

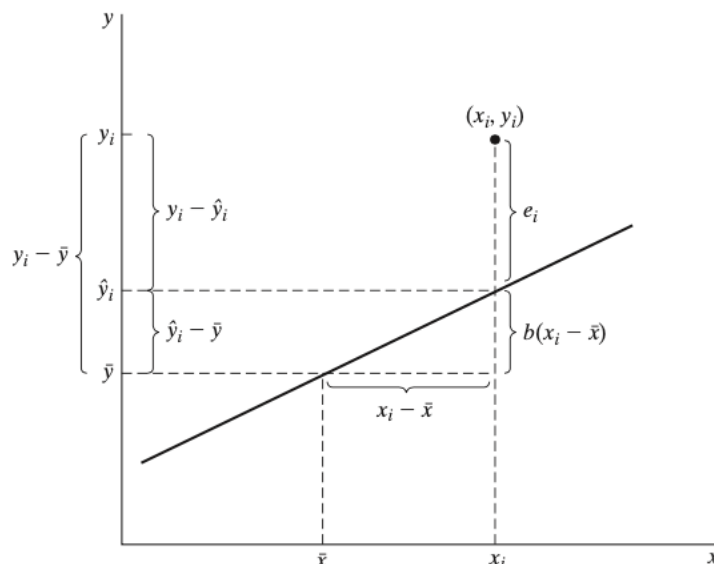
- Por lo tanto, cada columna de  $\mathbf{M}_1\mathbf{X}_2$  es la variable original en forma de desviaciones de sus medias muestrales
- Si la regresión múltiple de  $\mathbf{y}$  sobre  $\mathbf{X}$  contiene un término constante y las variables en la regresión no están correlacionadas, los coeficientes de la regresión múltiple son los mismos que los de las regresiones individuales simples de  $\mathbf{y}$  sobre la constante y cada variable
  - Cuando el regresor contiene un término constante, se pueden calcular los coeficientes con las regresiones de la desviación de  $\mathbf{y}$  de su media muestral sobre las columnas de  $\mathbf{X}$  en forma de desviaciones también
  - De este modo, la ortogonalidad de las columnas significa que las covarianzas muestrales son nulas, debido a que  $\mathbf{x}_k^*'\mathbf{x}_m^* = Cov(\mathbf{x}_k, \mathbf{x}_m) = 0$ , y eso hace que, si  $\mathbf{x}_k^*'\mathbf{x}_m^* = 0$ , entonces los coeficientes de cada regresor no tengan un vector de corrección (como se ha visto anteriormente)
- El uso de la regresión múltiple hace que se tenga que hacer un experimento conceptual que puede no poder hacerse en la práctica, el análisis *Ceteris Paribus*
  - Para poder realizarlo, se utiliza el coeficiente de correlación parcial, el cual usa el hecho de que cuando se hace una regresión sobre una variable, este no tiene ningún poder explicativo para las perturbaciones (por construcción)
- El criterio de ajuste original (la suma de residuos cuadrados) sugiere que esta es una medida para el ajuste de la línea de regresión a los datos, pero como esta se puede escalar arbitrariamente por un escalar, es mejor analizar el ajuste a partir de si la variación de  $\mathbf{x}$  es un buen predictor de la variación de  $\mathbf{y}$ 
  - La variación de la variable dependiente se define en términos de la desviación de esta de su media muestral, de modo que la variación total de  $\mathbf{y}$  es la suma de sus desviaciones cuadradas
    - En términos de la ecuación de regresión y para una observación individual, se pueden formular las siguientes expresiones:

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e} = \hat{\mathbf{y}} + \mathbf{e} \quad y_i = \mathbf{x}_i'\mathbf{b} + e_i = \hat{y}_i + e_i$$

- Si la regresión contiene un término constante, entonces las propiedades algebraicas de los mínimos cuadrados anteriormente vistas aplican y se puede ver que  $y_i - \hat{y}_i$  se puede expresar de la siguiente manera:

$$y_i - \bar{y}_i = \hat{y}_i - \bar{y}_i + e_i = \mathbf{x}_i' \mathbf{b} - \bar{\mathbf{x}}_i' \mathbf{b} + e_i = (\mathbf{x}_i - \bar{\mathbf{x}}_i)' \mathbf{b} + e_i$$

- Intuitivamente, la regresión parece que se ajusta mejor si las desviaciones de  $y$  de su media muestral se deben más a las desviaciones de  $x$  de su media que por los residuos



- Como ambos términos en la descomposición de  $y_i - \bar{y}_i$  suman cero, se utiliza el cuadrado con tal de cuantificar el ajuste. Para el conjunto completo de observaciones, se puede expresar la descomposición en términos de desviaciones de la media muestral

$$\mathbf{M}^0 \mathbf{y} = \mathbf{M}^0 \mathbf{X} \mathbf{b} + \mathbf{M}^0 \mathbf{e} = \mathbf{M}^0 \mathbf{X} \mathbf{b} + \mathbf{e}$$

- Como  $\mathbf{e}' \mathbf{M}^0 \mathbf{X} = \mathbf{e}' \mathbf{X} = \mathbf{0}$ , la suma total de cuadrados se puede expresar de la siguiente manera:

$$\mathbf{y}' \mathbf{M}^0 \mathbf{y} = \mathbf{b}' \mathbf{X}' \mathbf{M}^0 \mathbf{X} \mathbf{b} + \mathbf{e}' \mathbf{e}$$

$$Total\ Sum\ of\ Sq. = Resid.\ Sum\ of\ Sq. + Expl.\ Sum\ of\ Sq.$$

- A partir de esta descomposición, es posible crear una medida de ajuste llamada coeficiente de determinación, denotada por  $R^2$ , y que toma valores del cero al uno (incluidos)

$$\frac{ESS}{TSS} = \frac{\mathbf{b}' \mathbf{X}' \mathbf{M}^0 \mathbf{X} \mathbf{b}}{\mathbf{y}' \mathbf{M}_0 \mathbf{y}} = 1 - \frac{\mathbf{e}' \mathbf{e}}{\mathbf{y}' \mathbf{M}_0 \mathbf{y}}$$

- Este es nulo cuando los coeficientes son nulos (excepto el término constante) y el valor predicho de la regresión es siempre la media de  $y$ , de modo que las desviaciones de  $x$  de su

media no se traducen en diferentes predicciones de  $y$  y  $x$  no tiene poder explicativo

- Este es 1 cuando los valores de  $x$  e  $y$  se encuentran en el mismo hiperplano, de modo que todos los residuos son nulos
- Una manera equivalente de expresar esta medida es con la siguiente ecuación:

$$\hat{y}'M_0\hat{y} = b'X'M^0Xb \Rightarrow \hat{y}'M_0\hat{y} = \hat{y}'M_0y \Rightarrow R^2 = \frac{\hat{y}'M_0\hat{y}}{y'M_0y}$$

$$\Rightarrow R^2 = \frac{\hat{y}'M_0\hat{y}\hat{y}'M_0y}{y'M_0y\hat{y}'M_0\hat{y}} \Rightarrow R^2 = \frac{\sum_{i=1}^n [(y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})]^2}{\sum_{i=1}^n (y_i - \bar{y})^2 \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2}$$

- El cálculo de  $R^2$  se puede resumir en una tabla de análisis de la varianza como la siguiente:

	Source	Degrees of Freedom	Mean Square
Regression	$b'X'y - n\bar{y}^2$	$K - 1$ (assuming a constant term)	
Residual	$e'e$	$n - K$	$s^2$
Total	$y'y - n\bar{y}^2$	$n - 1$	$S_{yy}/(n - 1) = s_y^2$
Coefficient of determination		$R^2 = 1 - e'e/(y'y - n\bar{y}^2)$	

- Sin embargo, existen problemas al analizar el ajuste con  $R^2$ . La primera dificultad reside en que  $R^2$  siempre aumenta cuanto mayor es el número de regresores

- Comparando la regresión de  $y$  en  $X$  con la regresión de  $y$  en  $X$  y una variable adicional  $z$ , se asume que la primera produce una suma de residuos cuadrados  $e'e$  y la segunda una  $u'u$ . Además,  $c = (z^{*'}z^*)^{-1}(z^{*'}y^*)$  es el coeficiente de la variable  $z$ , de modo que si se inserta esto en la ecuación del cambio en la suma de errores cuadrados al añadir una variable  $z$  se obtiene la siguiente igualdad:

$$u'u = e'e - c^2(z^{*'}z^*) = e'e - [(z^{*'}z^*)^{-1}(z^{*'}y^*)]^2(z^{*'}z^*) =$$

$$= e'e - \frac{(z^{*'}y^*)^2}{(z^{*'}z^*)} = e'e(1 - r_{yz}^{*2})$$

- Siendo  $R_{Xz}^2$  el coeficiente de determinación en la regresión de  $y$  sobre  $X$  y la variable adicional  $z$  (que, en el caso multivariante, se puede sustituir por un conjunto de variables),  $R_X^2$  el coeficiente de determinación para la regresión de  $y$  sobre  $X$  y  $r_{yz}^{*2}$  es el coeficiente parcial entre  $y$  y  $z$  (controlando para  $X$ ), entonces se da la siguiente igualdad:



$$R_{xz}^2 = R_X^2 + (1 - R_X^2)r_{yz}^{*2}$$

- Por lo tanto, el  $R^2$  de una regresión no puede ser menor al añadir un regresor
- El  $\bar{R}^2$  es el coeficiente de determinación ajustado (por grados de libertad) se calcula de la siguiente manera:

$$\bar{R}^2 = 1 - \frac{e'e/(n-K)}{y'M^0y/(n-1)}$$

- Este coeficiente también se puede expresar de manera equivalente en términos de  $R^2$ :

$$\bar{R}^2 = 1 - \frac{n-1}{n-K}(1 - R^2)$$

- Este coeficiente puede decrecer cuando se añade un regresor adicional y puede ser hasta negativo (de modo que  $\bar{R}^2 < 0$ ). Si  $\bar{R}^2$  crece o decrece depende de si la contribución de la nueva variable al ajuste de la regresión más que compensa la corrección por la pérdida adicional de grados de libertad
- En una regresión múltiple  $\bar{R}^2$  decrecerá o crecerá cuando la variable  $x$  se elimina de la regresión si el cuadrado de la  $t$ -ratio asociada a la variable es mayor o menor, respectivamente, a 1

$$t_z^2 = \frac{c^2}{\left(\frac{u'u}{n-K-1}\right)(W'W)_{K+1,K+1}^{-1}}$$

- Una segunda dificultad con  $R^2$  está relacionado con el término constante del modelo, dado que si  $X$  no tiene una columna de unos, entonces no se puede calcular bien el valor de  $R^2$ , dado que pueda dar valores mucho más grandes o pequeños que la misma regresión incluyendo un término constante, puede salir del intervalo  $[0,1]$  y hasta podría ser negativo
  - Si no hay una columna de unos en  $X$ , entonces  $M^0e \neq e$  y  $e'M^0X \neq 0$  y eso hace que no se pueda eliminar el término  $2e'M^0Xb$  en  $y'M^0y = (M^0Xb + M^0e)'(M^0Xb + M^0e)$
- Aunque un criterio que se puede utilizar para poder comparar modelos es el  $R^2$ , no hay una referencia para poder comparar modelos porque la bondad del ajuste depende del contexto

- No se puede decir mucho de la calidad relativa de los ajustes de una regresión en diferentes contextos o con diferentes datos, aunque estos provengan del mismo mecanismo generador de datos. Pero, aunque el contexto o los datos sean los mismos, se necesita utilizar una misma referencia a la hora de comparar la bondad de ajuste de los modelos, normalmente relacionado a la variable dependiente
- Además,  $R^2$  es una medida de asociación lineal entre las variables dependientes y la independiente, de modo que, si la relación entre las variables no es lineal, esta será muy baja y puede confundir sobre la bondad del ajuste del modelo
- La interpretación de  $R^2$  como la proporción de varianza del modelo explicado por  $X$  no es razonable cuando no se utilizan mínimos cuadrados y/o no hay un coeficiente constante
- Es posible estudiar un resultado puramente algebraico que es muy útil para entender el cálculo de modelos de regresión lineal
  - En la regresión de  $y$  sobre  $X$  se supone que las columnas de  $X$  son linealmente transformadas, como multiplicando o dividiendo los datos por una constante (cambio de unidades de medida)
    - En este caso se supone que la matriz  $X$  con una transformación lineal es  $Z$ . De este modo, es interesante saber el efecto de una transformación lineal de la regresión  $y$  sobre  $X$  comparada con la regresión  $y$  sobre  $Z$
  - En la regresión lineal de  $y$  sobre  $Z = XP$  en donde  $P$  es una matriz invertible que transforma las columnas de  $X$ , el vector de coeficientes será  $P^{-1}b$  en donde  $b$  es el vector de coeficientes de en la regresión lineal de  $y$  sobre  $X$ 
    - Los coeficientes de la regresión de  $y$  sobre  $Z$  serían los siguientes:

$$\begin{aligned}
 d &= (Z'Z)^{-1}Z'y = ((XP)'(XP))^{-1}(XP)'y = \\
 &= (P'X'XP)^{-1}P'X'y = P^{-1}(X'X)P'^{-1}P'y = P^{-1}b
 \end{aligned}$$

- El vector de residuos  $u$  para la regresión de  $Y$  sobre  $Z$  es idéntico al vector  $e$  de la regresión original de  $Y$  sobre  $X$ , de modo que el numerador de  $1 - R^2$  es el mismo y el denominador no cambia

$$u = y - Z(P^{-1}b) = y - XPP^{-1}b = y - Xb = e$$

- Este resultado es muy útil porque para casos simples en donde se multipliquen o se dividan las unidades por un número  $p$ , el vector de parámetros se multiplica por  $1/p$

## El estimador de mínimos cuadrados

- Una de las razones por las que el estimador de mínimos cuadrados es tan usado es porque es muy fácil de calcular. Además, es un enfoque natural para la estimación porque utiliza la estructura del modelo y sus suposiciones (solo la de rango completo), es un predictor lineal óptimo para la variable dependiente y porque es el estimador más eficiente bajo unas suposiciones concretas
  - Siendo  $\mathbf{x}$  un vector de variables independientes en el modelo de regresión poblacional, se asume que los datos pueden ser estocásticos o fijos, que las perturbaciones de la población son ortogonales a las variables independientes (por lo que  $E(\boldsymbol{\varepsilon}|\mathbf{x}) = 0$  y  $Cov(\boldsymbol{\varepsilon}, \mathbf{x}) = 0$ ), entonces se pueden obtener las siguientes identidades:

$$E(\boldsymbol{\varepsilon}|\mathbf{x}) = 0 \Rightarrow E_{\mathbf{x}}[E_{\boldsymbol{\varepsilon}}(\mathbf{x}\boldsymbol{\varepsilon})] = E_{\mathbf{x}}[E_{\mathbf{y}}(\mathbf{x}(\mathbf{y} - \mathbf{x}'\boldsymbol{\beta}))] = 0$$

$$\Rightarrow E_{\mathbf{x}}[E_{\mathbf{y}}(\mathbf{x}\mathbf{y})] = E_{\mathbf{x}}(\mathbf{x}\mathbf{x}')\boldsymbol{\beta}$$

- Debido a que  $\mathbf{y} = \mathbf{x}'\boldsymbol{\beta}$ ,  $\mathbf{x}\mathbf{x}'\boldsymbol{\beta}$  no es una función de  $\mathbf{y}$  y la esperanza se toma para los valores de  $\mathbf{x}$
- Utilizando las ecuaciones normales de los mínimos cuadrados, se puede ver que  $\mathbf{X}'\mathbf{y} = \mathbf{X}\mathbf{X}'\mathbf{b}$  de modo que dividiendo por  $n$  se obtiene una expresión en términos de sumatoria, la cual sería una expresión análoga para una muestra:

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i y_i = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right) \mathbf{b}$$

- Si se cumplen las suposiciones de la ley de grandes números, la expresión anterior es un estimador de sus contrapartes anteriores (en términos de esperanzas)
- Alternativamente, se puede considerar el problema de encontrar un predictor lineal óptimo para  $\mathbf{y}$ . Ignorando suposiciones anteriores sobre la estocasticidad de  $\mathbf{X}$  y sobre la linealidad del estimador, se puede utilizar la regla del error cuadrático medio (MSE) como función objetivo
  - El valor esperado del error cuadrático medio del predictor se puede expresar de la siguiente manera:

$$\begin{aligned}
MSE &= E_y[E_x[(\mathbf{y} - \mathbf{x}'\boldsymbol{\gamma})^2]] = \\
&= E_y\left[E_x\left[(\mathbf{y} - E(\mathbf{y}|\mathbf{x}))^2\right]\right] + E_y[E_x[(E(\mathbf{y}|\mathbf{x}) - \mathbf{x}'\boldsymbol{\gamma})^2]]
\end{aligned}$$

- Por lo tanto, se quiere encontrar el vector  $\boldsymbol{\gamma}$  que minimice estas expectativas, por lo que solo es necesario minimizar la segunda parte de la expresión con respecto a  $\boldsymbol{\gamma}$  (la otra no depende). Como  $E(\mathbf{y}|\mathbf{x}) = \mathbf{x}'\boldsymbol{\gamma}$ , el segundo término no depende de  $\mathbf{y}$  y se pueden poner el operador de derivada dentro de la expresión

$$\begin{aligned}
\frac{\partial E_y[E_x[(E(\mathbf{y}|\mathbf{x}) - \mathbf{x}'\boldsymbol{\gamma})^2]]}{\partial \boldsymbol{\gamma}} &= E_y\left[E_x\left[\frac{\partial (E(\mathbf{y}|\mathbf{x}) - \mathbf{x}'\boldsymbol{\gamma})^2}{\partial \boldsymbol{\gamma}}\right]\right] = \\
&= -2E_y[E_x[\mathbf{x}(E(\mathbf{y}|\mathbf{x}) - \mathbf{x}'\boldsymbol{\gamma})]] = -2E_y[E_x[\mathbf{x}(\mathbf{x}'\boldsymbol{\gamma} - \mathbf{x}'\boldsymbol{\gamma})]] = 0
\end{aligned}$$

$$\Rightarrow E_y[E_x[\mathbf{x}E(\mathbf{y}|\mathbf{x})]] = E_x[E_y(\mathbf{x}\mathbf{x}')] \boldsymbol{\gamma}$$

- El resultado anterior permite obtener la condición necesaria para obtener el mínimo del MSE, la cual es la misma equivalencia encontrada anteriormente:

$$\begin{aligned}
E_y[E_x[\mathbf{x}(E(\mathbf{y}|\mathbf{x}) - \mathbf{x}'\boldsymbol{\gamma})]] &= Cov(\mathbf{x}, \mathbf{y}) + E(\mathbf{x})E_x[E(\mathbf{y}|\mathbf{x})] = \\
&= Cov(\mathbf{x}, \mathbf{y}) + E(\mathbf{x})E(\mathbf{y}) = E_x[E_y(\mathbf{x}\mathbf{y})] \\
&\Rightarrow E_x[E_y(\mathbf{x}\mathbf{y})] = E_x[E_y(\mathbf{x}\mathbf{x}')] \boldsymbol{\gamma}
\end{aligned}$$

- Si el mecanismo de generación de datos  $(x_i, y_i)_{i=1, \dots, n}$  es tal que la ley de grandes números aplica a los estimadores  $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i y_i = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'\right) \mathbf{b}$  para las matrices  $\mathbf{X}'\mathbf{y}$  y  $\mathbf{X}\mathbf{X}'\mathbf{b}$ , entonces el predictor lineal mínimo del error cuadrático esperado de  $y_i$  es la línea de regresión por mínimos cuadrados

- Asumiendo que los valores esperados existen, se pueden estimar con las sumatorias anteriormente vistas y esto significa que, sin importar la forma de la media condicional (la especificación de la regresión), el estimador de mínimos cuadrados es el predictor lineal que minimiza el MSE
- El estimador de mínimos cuadrados ordinario cumple varias propiedades, entre las cuales hay la de no estar sesgado o su eficiencia bajo las suposiciones del teorema de Gauss-Markov

- El estimador de mínimos cuadrados ordinario no tiene sesgo en ninguna de las muestras. Esto se puede demostrar haciendo una expansión del vector de parámetros y de  $\mathbf{y}$  y tomando el valor esperado condicional a  $\mathbf{X}$

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon}$$

$$\begin{aligned}\Rightarrow E(\mathbf{b}|\mathbf{X}) &= E(\boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon}|\mathbf{X}) = \boldsymbol{\beta} + E((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon}|\mathbf{X}) = \\ &= \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\boldsymbol{\varepsilon}|\mathbf{X}) = \boldsymbol{\beta}\end{aligned}$$

- Consecuentemente, se obtiene que la esperanza de  $\mathbf{b}$  es necesariamente  $\boldsymbol{\beta}$

$$E(\mathbf{b}) = E_X[E(\mathbf{b}|\mathbf{X})] = E_X[\boldsymbol{\beta}] = \boldsymbol{\beta}$$

- Por lo tanto, para cualquier conjunto particular de observaciones  $\mathbf{X}$ , el estimador de mínimos cuadrados ordinario no tiene sesgo (su esperanza es el vector de parámetros poblacionales)

$$E(\mathbf{b}) - \boldsymbol{\beta} = \boldsymbol{\beta} - \boldsymbol{\beta} = 0$$

- Se ha condicionado en  $\mathbf{X}$  y no en  $\mathbf{x}_i$  porque el estimador de mínimos cuadrados se puede entender como la contraparte muestral al vector de parámetros del predictor mínimo del MSE  $\boldsymbol{\gamma}$ , el cual pertenece la función poblacional, y porque, desde un punto de vista bayesiano, se puede entender que  $E(\mathbf{b}|\mathbf{X}) = \boldsymbol{\beta}$  se espera que el estimador se comporte como  $\boldsymbol{\beta}$  dados los datos muestrales finitos  $\mathbf{X}$
- En todo momento se ha utilizado la suposición de que la especificación correcta del modelo de regresión es  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}$ , pero hay muchos tipos de errores de especificación que se pueden hacer al construir el modelo de regresión. Los más comunes son la omisión de variables relevantes y la inclusión de variables superfluas

- Suponiendo que el modelo correctamente especificado es un modelo de la forma  $\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$  (en donde  $\mathbf{X}_1$  y  $\mathbf{X}_2$  son una partición de  $\mathbf{X}$  con  $K_1$  y  $K_2$  columnas respectivamente), si se hace una regresión de  $\mathbf{y}$  sobre  $\mathbf{X}_1$  (sin incluir  $\mathbf{X}_2$ ), el estimador y la esperanza condicionada de este a la muestra entera  $\mathbf{X}$  permite obtener la fórmula de variable omitida

$$\begin{aligned}\mathbf{b}_1 &= (\mathbf{X}_1'\mathbf{X}_1)^{-1}\mathbf{X}_1'\mathbf{y} = (\mathbf{X}_1'\mathbf{X}_1)^{-1}\mathbf{X}_1'(\mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}) = \\ &= (\mathbf{X}_1'\mathbf{X}_1)^{-1}\mathbf{X}_1'\mathbf{X}_1\boldsymbol{\beta}_1 + (\mathbf{X}_1'\mathbf{X}_1)^{-1}\mathbf{X}_1'\mathbf{X}_2\boldsymbol{\beta}_2 + (\mathbf{X}_1'\mathbf{X}_1)^{-1}\mathbf{X}_1'\boldsymbol{\varepsilon} =\end{aligned}$$

$$= \beta_1 + (X_1'X_1)^{-1}X_1'X_2\beta_2 + (X_1'X_1)^{-1}X_1'\varepsilon$$

$$E(b_1|X) = \beta_1 + (X_1'X_1)^{-1}X_1'X_2\beta_2 + (X_1'X_1)^{-1}X_1'E(c|X) =$$

$$= \beta_1 + (X_1'X_1)^{-1}X_1'X_2\beta_2 = \beta_1 + P_{1,2}\beta_2$$

$$\text{where } P_{1,2} \equiv (X_1'X_1)^{-1}X_1'X_2$$

- La matriz  $P_{1,2}$  es una matriz  $K_1 \times K_2$  la cual representa una matriz cuyas columnas son los parámetros de las variables para la regresión de las columnas de  $X_2$  sobre las de  $X_1$
- Con la fórmula de la omisión de variables, es fácil deducir la dirección del sesgo en casos simples como el caso univariante en el que se omite una variable, pero si en  $X_1$  y en  $X_2$  hay más variables, entonces los múltiples coeficientes pueden tener diferentes signos y la dirección no es clara. En el caso simple,  $P_{1,2}$  se puede entender como la *ratio* entre la covarianza de  $X_1$  y  $X_2$  y la varianza  $X_1$ , pero como la esperanza está condicionada a  $X$ , el resultado es constante

$$\frac{Cov(X_1, X_2)}{Var(X_1)} = E[(X_1'X_1)^{-1}|X]E[X_1'X_2|X] = (X_1'X_1)^{-1}X_1'X_2$$

- Si la especificación correcta de la regresión es  $y = X\beta + \varepsilon$  y se estima como si  $y = X_1\beta_1 + X_2\beta_2 + \varepsilon$  fuera correcta, entonces se estarían incluyendo variables innecesarias. No obstante, el modelo no es incorrecto, sino que se estima  $\beta$  correctamente, pero se falla en imponer la restricción  $\beta_2 = 0$  (en el caso anterior, se estimaba  $\beta$  sujeto a una restricción incorrecta  $\beta_2 = 0$ )

$$E(b|X) = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} = \begin{pmatrix} \beta_1 \\ 0 \end{pmatrix}$$

- El problema de esto reside en que incluir variables irrelevantes hace que la precisión de los estimadores, en el sentido que la matriz de varianzas y covarianzas de la regresión de  $y$  sobre  $X_1$  es siempre menor a la de la regresión de  $y$  sobre  $X_1$  y  $X_2$
- Si los regresores se pueden tratar como valores fijos, entonces la varianza muestral del estimador de mínimos cuadrados se puede derivar tratando  $X$  como una matriz de constantes. Alternativamente, se puede permitir que  $X$  sea estocástica y condicionar la esperanza a  $X$

- Considerando el vector de parámetros, si se define  $\mathbf{A} \equiv (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ , entonces bajo las suposiciones hechas,  $\mathbf{b}$  es un estimador lineal y no tiene sesgo independientemente de la distribución de  $\boldsymbol{\varepsilon}$  (solo se necesita la suposición de independencia de media)

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon}$$

$$E(\mathbf{b}|\mathbf{X}) = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\boldsymbol{\varepsilon}|\mathbf{X}) = \boldsymbol{\beta}$$

- A partir de la suposición de que la varianza es constante y de la simetría de  $\mathbf{X}'\mathbf{X}$ , se puede obtener la varianza del estimador de mínimos cuadrados

$$\begin{aligned} Var(\mathbf{b}|\mathbf{X}) &= E[(\mathbf{b} - \boldsymbol{\beta})(\mathbf{b} - \boldsymbol{\beta})'|\mathbf{X}] = \\ &= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon}((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon})'|\mathbf{X}] = \\ &= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}|\mathbf{X}] = E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}|\mathbf{X}] \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'|\mathbf{X}]\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\sigma^2\mathbf{I})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2(\mathbf{X}'\mathbf{X})^{-1} = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})} \end{aligned}$$

- En el modelo de regresión lineal con una matriz de regresores  $\mathbf{X}$ , el estimador de mínimos cuadrados  $\mathbf{b}$  es el estimador lineal sin sesgo de mínima varianza (cc) de  $\boldsymbol{\beta}$ . Para cualquier vector de constantes  $\mathbf{w}$ , el BLUE de  $\mathbf{w}'\boldsymbol{\beta}$  es  $\mathbf{w}'\mathbf{b}$ , donde  $\mathbf{b}$  es el estimador de mínimos cuadrados

- El teorema no utiliza la suposición de que las perturbaciones se distribuyen de manera normal, por lo que solo se necesitan las primeras cuatro (la quinta se puede usar o no)
- Para poder demostrar este teorema, se puede proponer el uso de un estimador lineal  $\mathbf{b}_L = \mathbf{C}\mathbf{y}$  tal que  $E[\mathbf{b}_L|\mathbf{X}] = \boldsymbol{\beta}$  y encontrar un estimador de esta clase que tenga una matriz de varianzas y covarianzas menor a  $\mathbf{b}$ . Esto implica que  $\mathbf{C}\mathbf{X} = \mathbf{I}$  y, por tanto, hay varios candidatos para  $\mathbf{C}$

$$E(\mathbf{b}_L|\mathbf{X}) = E(\mathbf{C}\mathbf{y}|\mathbf{X}) = E(\mathbf{C}\mathbf{X}\boldsymbol{\beta} + \mathbf{C}\boldsymbol{\varepsilon}|\mathbf{X}) = \boldsymbol{\beta} \Rightarrow \mathbf{C}\mathbf{X} = \mathbf{I}$$

- Definiendo una matriz  $\mathbf{D} \equiv \mathbf{C} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ , es posible demostrar que la varianza del estimador  $\mathbf{b}_L$  es mayor a la de  $\mathbf{b}$  porque  $\mathbf{q}'\mathbf{D}\mathbf{D}'\mathbf{q} \geq 0$  (la matriz es semidefinida positiva), obteniendo así el siguiente resultado:

$$Var(\mathbf{b}_L|\mathbf{X}) = Var(\mathbf{b}|\mathbf{X}) + \sigma^2 \mathbf{D}\mathbf{D}' = \sigma^2(\mathbf{X}'\mathbf{X})^{-1} + \sigma^2 \mathbf{D}\mathbf{D}'$$

- La demostración en el caso de  $\mathbf{w}'\mathbf{b}$  seguiría un patrón igual al de esta demostración y se comprobaría la segunda proposición
- Si se permite que los regresores sean estocásticos, entonces se pueden encontrar resultados iguales no condicionados a los regresores  $\mathbf{X}$  a través de propiedades del valor esperado

- Para obtener la varianza incondicional de  $\mathbf{b}$ , se puede utilizar la ley de esperanzas iteradas y el hecho de que  $E(\mathbf{b}|\mathbf{X}) = \boldsymbol{\beta}$

$$\begin{aligned} Var(\mathbf{b}) &= E_X[Var(\mathbf{b}|\mathbf{X})] + Var_X[E(\mathbf{b}|\mathbf{X})] = \\ &= E_X[\sigma^2(\mathbf{X}'\mathbf{X})^{-1}] + Var_X(\boldsymbol{\beta}) = \sigma^2 E_X[(\mathbf{X}'\mathbf{X})^{-1}] \end{aligned}$$

- Este resultado permite ver que la varianza de  $\mathbf{b}$  solo se puede describir en términos del comportamiento medio de  $\mathbf{X}$ . Además, utilizando la demostración anterior, se puede ver como  $Var(\mathbf{b}|\mathbf{X}) \leq Var(\mathbf{b}_L|\mathbf{X})$ , pero debido a que esto se cumple para cualquier  $\mathbf{X}$  particular, también se debe cumplir para  $Var(\mathbf{b}) = E_X[Var(\mathbf{b}|\mathbf{X})]$ , por lo que se da la siguiente desigualdad:

$$Var(\mathbf{b}|\mathbf{X}) \leq Var(\mathbf{b}_L|\mathbf{X}) \Rightarrow Var(\mathbf{b}) \leq Var(\mathbf{b}_L)$$

- En consecuencia, en el modelo de regresión lineal, el estimador de mínimos cuadrados ordinarios es un MVLUE independientemente de si  $\mathbf{X}$  es estocástico o no mientras se cumplan las otras suposiciones
- Para poder hacer contrastes de hipótesis sobre  $\boldsymbol{\beta}$  o para construir intervalos de confianza, se requiere un estimado para la matriz de varianzas y covarianzas. Una manera natural de estimar  $Var(\mathbf{b}|\mathbf{X})$  es a través de los residuos, debido a que las perturbaciones cumplen  $Var(\varepsilon_i|\mathbf{X}) = E(\varepsilon_i^2|\mathbf{X}) - E(\varepsilon_i|\mathbf{X})^2 = E(\varepsilon_i^2|\mathbf{X}) = \sigma^2$  y los residuos son estimadores de los errores

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n e_i^2$$

- No obstante, los residuos son unos estimadores imperfectos de las perturbaciones porque  $e_i$  depende de  $\boldsymbol{\beta}$  y esta no se puede observar directamente

$$e_i = y_i - \mathbf{x}_i'\mathbf{b} = \mathbf{x}_i'\boldsymbol{\beta} + \varepsilon_i - \mathbf{x}_i'\mathbf{b} = \varepsilon_i + \mathbf{x}_i'(\mathbf{b} - \boldsymbol{\beta})$$



- Formalmente, se puede derivar  $\sum_{i=1}^n e_i^2$  para poder ver que este es igual a  $\sigma^2(n - K)$

$$\begin{aligned} E(\mathbf{e}'\mathbf{e}|\mathbf{X}) &= E((\mathbf{M}\boldsymbol{\varepsilon})'\mathbf{M}\boldsymbol{\varepsilon}|\mathbf{X}) = E(\boldsymbol{\varepsilon}'\mathbf{M}'\mathbf{M}\boldsymbol{\varepsilon}|\mathbf{X}) = E(\boldsymbol{\varepsilon}'\mathbf{M}\boldsymbol{\varepsilon}|\mathbf{X}) = \\ &= E(\text{tr}(\boldsymbol{\varepsilon}'\mathbf{M}\boldsymbol{\varepsilon})|\mathbf{X}) = E(\text{tr}(\mathbf{M}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}')|\mathbf{X}) = E(\mathbf{M}\text{tr}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}')|\mathbf{X}) = \\ &= \text{tr}[\mathbf{M}E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'|\mathbf{X})] = \text{tr}(\mathbf{M}\sigma^2\mathbf{I}) = \sigma^2\text{tr}(\mathbf{M}) = \sigma^2(n - K) \end{aligned}$$

$$\begin{aligned} \text{as } \text{tr}(\mathbf{M}) &= \text{tr}(\mathbf{I}_n) - \text{tr}(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') = \\ &= \text{tr}(\mathbf{I}_n) - \text{tr}(\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}) = \text{tr}(\mathbf{I}_n) - \text{tr}(\mathbf{I}_K) = n - K \end{aligned}$$

- Este resultado muestra que el estimador está sesgado hacia cero (aunque se vuelva menor cuando el tamaño de la muestra  $n$  incrementa). Un estimador no sesgado (tanto condicionalmente como de manera incondicional) sería  $s^2 \equiv \mathbf{e}'\mathbf{e}/(n - K)$ , y el error estándar de la regresión sería  $s$

$$s^2 \equiv \frac{\mathbf{e}'\mathbf{e}}{n - K}$$

- Por lo tanto, con  $s^2$  se puede hacer una estimación de la varianza de  $\mathbf{b}$  condicional a  $\mathbf{X}$  con la muestra, y a su vez, obtener el error estándar para cada  $b_k$  a través de la raíz del elemento  $kk$  (de la diagonal) de la matriz

$$\text{Est. Var}(\mathbf{b}|\mathbf{X}) = s^2(\mathbf{X}'\mathbf{X})^{-1}$$

- Hasta ahora no se ha usado la suposición de normalidad de las perturbaciones, la cual es una suposición bastante útil para poder construir intervalos de confianza y tests de hipótesis

- Si se asume que las perturbaciones siguen una distribución normal multivariante, entonces se puede comprobar la distribución del vector de parámetros y de los estimadores siguen una distribución normal multivariante

$$\mathbf{b} | \mathbf{X} \sim N(\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}) \quad b_k | \mathbf{X} \sim N(\beta_k, \sigma^2(\mathbf{X}'\mathbf{X})_{kk}^{-1})$$

- Si no se hace suposiciones sobre la distribución de las perturbaciones, no es posible decir nada de la distribución de parámetros (ya sea condicional o incondicional)

- Hasta ahora se han establecido resultados para una muestra finita, pero también es posible describir las propiedades asintóticas o de muestra grande del estimador de mínimos cuadrados

- Aunque se haya visto que el estimador no es sesgado, basarse solo en un criterio de sesgo no es muy útil

- Esto es debido a que en la mayoría de casos, los estimadores estarán sesgados y solo se puede esperar que estos tiendan a no tener sesgo cuanto más información se use, y, además, que un estimador no esté sesgado no implica que más información sea mejor para la estimación de los parámetros. La propiedad de consistencia mejora el criterio del sesgo en ambas direcciones
- Siendo  $\mathbf{X}$  estocástica o fija, se asume que  $(\mathbf{x}_i, \varepsilon_i)$  es una secuencia de observaciones independientes y se asume que  $(1/n)\mathbf{X}\mathbf{X}' \xrightarrow{P} \mathbf{Q}$  cuando  $n \rightarrow \infty$ , donde  $\mathbf{Q}$  es una matriz semidefinida positiva. Esto permite obtener los siguientes resultados:

$$\mathbf{b} = \boldsymbol{\beta} + \left(\frac{1}{n}\mathbf{X}'\mathbf{X}\right)^{-1} \left(\frac{1}{n}\mathbf{X}'\boldsymbol{\varepsilon}\right)$$

$$\text{plim}_{n \rightarrow \infty} \mathbf{b} = \boldsymbol{\beta} + \mathbf{Q}^{-1} \left( \text{plim}_{n \rightarrow \infty} \frac{1}{n}\mathbf{X}'\boldsymbol{\varepsilon} \right) = \boldsymbol{\beta} + \mathbf{Q}^{-1} \left( \text{plim}_{n \rightarrow \infty} \bar{\mathbf{w}} \right)$$

$$\text{where } \mathbf{w}_i \equiv \mathbf{x}_i \varepsilon_i \text{ \& } \bar{\mathbf{w}} \equiv \frac{1}{n} \sum_{i=1}^n \mathbf{w}_i = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i = \frac{1}{n} \mathbf{X}'\boldsymbol{\varepsilon}$$

- A partir de las suposiciones, es posible obtener la esperanza y la varianza (condicional e incondicional) del estimador para poder comprobar como el estimador es consistente con las suposiciones hechas

$$E(\mathbf{w}_i) = E_X[E(\mathbf{w}_i|\mathbf{x}_i)] = E_X[\mathbf{x}_i E(\varepsilon_i|\mathbf{x}_i)] = 0$$

$$\Rightarrow E(\bar{\mathbf{w}}) = \frac{1}{n} E \left( \frac{1}{n} \sum_{i=1}^n \mathbf{w}_i \right) = 0$$

$$\Rightarrow E(\bar{\mathbf{w}}|\mathbf{X}) = \frac{1}{n} \mathbf{X}' E(\boldsymbol{\varepsilon}|\mathbf{X}) = 0$$

$$\text{Var}(\bar{\mathbf{w}}|\mathbf{X}) = E(\bar{\mathbf{w}}\bar{\mathbf{w}}'|\mathbf{X}) = E \left( \frac{1}{n} \mathbf{X}' \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}' \mathbf{X} \frac{1}{n} \middle| \mathbf{X} \right) =$$

$$\begin{aligned}
&= \frac{1}{n} \mathbf{X}' E(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}' | \mathbf{X}) \mathbf{X} \frac{1}{n} = \frac{1}{n} \mathbf{X}' E(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}' | \mathbf{X}) \mathbf{X} \frac{1}{n} = \frac{1}{n} \mathbf{X}' \sigma^2 \mathbf{I} \mathbf{X} \frac{1}{n} = \\
&= \left( \frac{\sigma^2}{n} \right) \left( \frac{\mathbf{X}' \mathbf{X}}{n} \right)
\end{aligned}$$

$$Var(\bar{\mathbf{w}}) = E_X[Var(\bar{\mathbf{w}} | \mathbf{X})] + Var_X[E(\bar{\mathbf{w}} | \mathbf{X})] = \left( \frac{\sigma^2}{n} \right) E_X \left( \frac{\mathbf{X}' \mathbf{X}}{n} \right)$$

- Debido a la suposición hecha sobre la convergencia de  $(1/n)\mathbf{X}'\mathbf{X}$ , tanto la esperanza como la varianza son nulas cuando  $n \rightarrow \infty$  y eso hace que la media cuadrática tienda a cero y que, por tanto,  $(1/n)\mathbf{X}'\boldsymbol{\varepsilon} \xrightarrow{P} 0$ , haciendo que  $\mathbf{b}$  sea un estimador consistente

$$\lim_{n \rightarrow \infty} \left( \frac{\sigma^2}{n} \right) \left( \frac{\mathbf{X}' \mathbf{X}}{n} \right) = 0 \mathbf{Q} = \mathbf{0}$$

$$\Rightarrow \left( \frac{1}{n} \right) \mathbf{X}' \boldsymbol{\varepsilon} \xrightarrow{P} 0 \Rightarrow \text{plim}_{n \rightarrow \infty} \mathbf{b} = \boldsymbol{\beta} + \mathbf{Q}^{-1} \mathbf{0} = \boldsymbol{\beta}$$

- Las suposiciones hechas pueden ser muy restrictivas dependiendo del contexto, por lo que las condiciones de Grenander suelen ser más débiles y se cumplen en casi todas las ocasiones

**G1.** For each column of  $\mathbf{X}$ ,  $\mathbf{x}_k$ , if  $d_{nk}^2 = \mathbf{x}_k' \mathbf{x}_k$ , then  $\lim_{n \rightarrow \infty} d_{nk}^2 = +\infty$ . Hence,  $\mathbf{x}_k$  does not degenerate to a sequence of zeros. Sums of squares will continue to grow as the sample size increases. No variable will degenerate to a sequence of zeros.

**G2.**  $\lim_{n \rightarrow \infty} x_{ik}^2 / d_{nk}^2 = 0$  for all  $i = 1, \dots, n$ . This condition implies that no single observation will ever dominate  $\mathbf{x}_k' \mathbf{x}_k$ , and as  $n \rightarrow \infty$ , individual observations will become less important.

**G3.** Let  $\mathbf{R}_n$  be the sample correlation matrix of the columns of  $\mathbf{X}$ , excluding the constant term if there is one. Then  $\lim_{n \rightarrow \infty} \mathbf{R}_n = \mathbf{C}$ , a positive definite matrix. This condition implies that the full rank condition will always be met. We have already assumed that  $\mathbf{X}$  has full rank in a finite sample, so this assumption ensures that the condition will never be violated.

- Otra suposición de la que se puede prescindir es la suposición de normalidad de las perturbaciones. Bajo suposiciones generales razonables sobre los procesos generadores de datos muestrales, las distribuciones asintóticas proporcionan una base para la inferencia estadística en el modelo de regresión

- Usando el teorema central del límite y suponiendo que las observaciones son independientes, se puede obtener el tipo de resultado anteriormente visto

$$\left( \frac{\mathbf{X}' \mathbf{X}}{n} \right)^{-1} \xrightarrow{P} \mathbf{Q}^{-1} \quad \& \quad \left( \frac{1}{\sqrt{n}} \right) \mathbf{X}' \boldsymbol{\varepsilon} \xrightarrow{D} N(\mathbf{0}, \sigma^2 \mathbf{Q})$$

$$\Rightarrow \mathbf{Q}^{-1} \left( \frac{1}{\sqrt{n}} \right) \mathbf{X}' \boldsymbol{\varepsilon} \xrightarrow{D} N(\mathbf{Q}^{-1} \mathbf{0}, \mathbf{Q}^{-1} (\sigma^2 \mathbf{Q}) \mathbf{Q}^{-1'})$$

$$\Rightarrow \mathbf{Q}^{-1} \left( \frac{1}{\sqrt{n}} \right) \mathbf{X}' \boldsymbol{\varepsilon} \xrightarrow{D} N(\mathbf{0}, \sigma^2 \mathbf{Q}^{-1})$$

$$\sqrt{n}(\mathbf{b} - \boldsymbol{\beta}) = \left( \frac{\mathbf{X}' \mathbf{X}}{n} \right)^{-1} \left( \frac{1}{\sqrt{n}} \right) \mathbf{X}' \boldsymbol{\varepsilon} \Rightarrow \sqrt{n}(\mathbf{b} - \boldsymbol{\beta}) \xrightarrow{D} N(\mathbf{0}, \sigma^2 \mathbf{Q}^{-1})$$

- Si  $\boldsymbol{\varepsilon}$  se distribuyen independientemente con media nula y varianza finita  $\sigma^2$ , y  $\mathbf{X}$  es tal que las condiciones de Grenander se cumplen, entonces la distribución asintótica de  $\mathbf{b}$  es la siguiente:

$$\mathbf{b} \approx N \left( \boldsymbol{\beta}, \frac{\sigma^2}{n} \mathbf{Q}^{-1} \right)$$

- Para poder completar la derivación de las propiedades asintóticas de  $\mathbf{b}$  es necesario estimar la varianza asintótica  $(\sigma^2/n) \mathbf{Q}^{-1}$  a través de  $s^2$

$$s^2 = \frac{1}{n-K} \boldsymbol{\varepsilon}' \mathbf{M} \boldsymbol{\varepsilon}$$

- Expandiendo el término, se obtiene una expresión con la cual se puede comprobar que el estimador  $s^2$  tiende a la media de los errores cuadráticos

$$\begin{aligned} s^2 &= \frac{1}{n-K} \boldsymbol{\varepsilon}' (\mathbf{I} - \mathbf{X}(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}') \boldsymbol{\varepsilon} = \\ &= \frac{1}{n-K} (\boldsymbol{\varepsilon}' \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}' \mathbf{X}(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \boldsymbol{\varepsilon}) = \\ &= \frac{n}{n-K} \left[ \frac{\boldsymbol{\varepsilon}' \boldsymbol{\varepsilon}}{n} - \left( \frac{\boldsymbol{\varepsilon}' \mathbf{X}}{n} \right) \left( \frac{\mathbf{X}' \mathbf{X}}{n} \right)^{-1} \left( \frac{\mathbf{X}' \boldsymbol{\varepsilon}}{n} \right) \right] \\ &\Rightarrow \frac{n}{n-K} \left[ \frac{\boldsymbol{\varepsilon}' \boldsymbol{\varepsilon}}{n} - \left( \frac{\boldsymbol{\varepsilon}' \mathbf{X}}{n} \right) \left( \frac{\mathbf{X}' \mathbf{X}}{n} \right)^{-1} \left( \frac{\mathbf{X}' \boldsymbol{\varepsilon}}{n} \right) \right] \xrightarrow{P} \frac{\boldsymbol{\varepsilon}' \boldsymbol{\varepsilon}}{n} = \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 \end{aligned}$$

- Por lo tanto, se puede ver como  $E(\varepsilon_i^2) = \sigma^2$  (si son independientes), y asumiendo que se tienen momentos finitos (hasta el segundo como mínimo), se puede ver como  $s^2$  converge en probabilidad a  $\sigma^2$ , y eso hace que el estimador  $s^2 (\mathbf{X}' \mathbf{X}/n)^{-1}$  sea consistente para  $\sigma^2 (\mathbf{X}' \mathbf{X}/n)^{-1}$  y que, por tanto, se pueda estimar la varianza asintótica con este

$$s^2 \xrightarrow{P} \sigma^2 \Rightarrow s^2(\mathbf{X}'\mathbf{X}/n)^{-1} \xrightarrow{P} \sigma^2 \mathbf{Q}^{-1}$$

$$\Rightarrow Est. Asy. Var(\mathbf{b}) = s^2(\mathbf{X}'\mathbf{X})^{-1}$$

- Igual que el teorema de Gauss-Markov asegura la eficiencia del estimador de mínimos cuadrados en muestra finita, se pueden definir las condiciones para las cuales el estimador es asintóticamente eficiente
  - Un estimado es asintóticamente eficiente si es consistente, asintóticamente normal y tiene una matriz asintótica de varianzas y covarianzas que no es mayor a la matriz asintótica de varianzas y covarianzas de cualquier otro estimador
- Debido a que se han establecido propiedades de muestra finita y propiedades asintóticas del estimador sin el uso de la suposición de normalidad, entonces se puede decir que el estimado de métodos cuadrados es robusto a violaciones de la suposición de normalidad. En particular, se puede decir que el estimador es robusto a suposiciones de la distribución de  $\varepsilon$ 
  - El estimador de un modelo es robusto si es insensible a desviaciones de las suposiciones básicas del modelo. En términos econométricos, los estimadores robustos mantienen las propiedades deseadas aún sin que se cumplan las suposiciones que motivan el estimado del modelo
    - Por lo tanto, la robustez de un estimador es relativo a una suposición o conjunto de suposiciones concretas, no es un concepto global (sino relativo)
    - La robustez no es necesariamente una característica bien definida para un estimador, por lo que se considera una caracterización amplia de las propiedades asintóticas de ciertas estimaciones y procedimientos
    - En el contexto de las seis suposiciones básicas del modelo de regresión lineal múltiples, se estableció la consistencia del estimador  $\mathbf{b}$  debido a la suposición de rango completo y a la de independencia de los errores  $\varepsilon$  con las  $\mathbf{X}$ . Por lo tanto, otras suposiciones como la normalidad, la heterocedasticidad y la no autocorrelación no son necesarias para obtener las propiedades del estimador de mínimos cuadrados
  - La derivación de la varianza asintótica se basaba en la suposición de la homoscedasticidad y la no autocorrelación. No obstante, se puede hacer una formulación más general para motivar el uso de un estimador robusto a la heterocedasticidad

- Considerando que  $Var(\varepsilon_i|x_i) = \sigma_i^2$  (en donde se asume que la variación se debe a  $x_i$ ) se puede obtener la matriz de varianzas y covarianzas del vector de parámetros de mínimos cuadrados ordinarios y su matriz de varianzas y covarianzas asintótica

$$\mathbf{b} = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1} \sum_i \mathbf{x}_i \varepsilon_i$$

$$\Rightarrow Var(\mathbf{b}|\mathbf{X}) = (\mathbf{X}'\mathbf{X})^{-1} \left[ \sum_i \sigma_i^2 \mathbf{x}_i \mathbf{x}_i' \right] (\mathbf{X}'\mathbf{X})^{-1}$$

$$\begin{aligned} \Rightarrow Asy. Var(\mathbf{b}|\mathbf{X}) &= \frac{1}{n} \mathbf{Q}^{-1} \left[ \text{plim}_{n \rightarrow \infty} \frac{1}{n} \sum_i \sigma_i^2 \mathbf{x}_i \mathbf{x}_i' \right] \mathbf{Q}^{-1} = \\ &= \frac{1}{n} \mathbf{Q}^{-1} \mathbf{Q}^* \mathbf{Q}^{-1} \end{aligned}$$

- Una estrategia robusta a la heterocedasticidad no especificada es utilizar un estimador factible de  $\mathbf{Q}^*$  tal como el estimador robusto a la heterocedasticidad de White. Por lo tanto, se puede estimar la varianza asintótica de manera robusta a la heterocedasticidad no especificada con él

$$\mathbf{W}_{het} = \frac{1}{n} \sum_i e_i^2 \mathbf{x}_i \mathbf{x}_i' \quad \text{where} \quad e_i = y_i - \mathbf{x}_i' \mathbf{b}$$

$$Est. Asy. Var(\mathbf{b}|\mathbf{X}) = n \mathbf{Q}^{-1} \mathbf{W}_{het} \mathbf{Q}^{-1}$$

- De este modo, no se necesita la suposición de homoscedasticidad para obtener los errores estándar apropiados del estimador de mínimos cuadrados ordinarios
- De manera análoga, se puede derivar un estimador de la varianza asintótica cuando en la base de datos hay clústeres o observaciones en grupos
  - Cuando hay grupos de datos muestrales que están relacionados, se supone que hay  $C$  grupos o clústeres de observaciones (las cuales se indexan por  $c = 1, 2, \dots, C$ ). Las observaciones un clúster se agrupan por las correlaciones entre observaciones de un mismo grupo
  - Asumiendo que hay  $N_c$  observaciones en un clúster, entonces  $n = \sum_{c=1}^C N_c$  y el modelo de regresión es el siguiente:

$$y_{i,c} = \mathbf{x}_{i,c}' \boldsymbol{\beta} + \varepsilon_{i,c}$$

- El estimador de mínimos cuadrados con clústeres permite obtener la matriz de varianzas y covarianzas de  $\mathbf{b}$ , donde  $\mathbf{X}_c$  es la matriz  $N_c \times K$  de variables exógenas para el clúster  $c$ ,  $\varepsilon_c$  es el vector de  $N_c$  perturbaciones y  $\mathbf{\Omega}_c$  es una matriz de varianzas y covarianzas del vector completo  $\varepsilon_c$  que representa sus correlaciones desestructuradas. Asumiendo que los clústeres son independientes, se obtiene la siguiente fórmula:

$$\mathbf{b} = \boldsymbol{\beta} + \left( \sum_{c=1}^C \mathbf{X}_c' \mathbf{X}_c \right)^{-1} \left[ \sum_{c=1}^C \sum_{i=1}^{N_c} x_{i,c} \varepsilon_{i,c} \right] = \boldsymbol{\beta} + (\mathbf{X}' \mathbf{X})^{-1} \left( \sum_{c=1}^C \mathbf{X}_c \varepsilon_c \right)$$

$$\Rightarrow \text{Var}(\mathbf{b}|\mathbf{X}) = (\mathbf{X}' \mathbf{X})^{-1} \left( \sum_{c=1}^C \mathbf{X}_c \mathbf{\Omega}_c \mathbf{X}_c' \right) (\mathbf{X}' \mathbf{X})^{-1}$$

- Como antes, se puede encontrar la varianza asintótica y, a partir de esta, encontrar un estimador de la varianza asintótica parecida a la encontrada con el estimador de White

$$\text{Asy. Var}(\mathbf{b}) = \frac{1}{C} \mathbf{Q}^{-1} \left( \text{plim}_{n \rightarrow \infty} \frac{1}{C} \sum_{c=1}^C \mathbf{X}_c \mathbf{\Omega}_c \mathbf{X}_c' \right) \mathbf{Q}^{-1}$$

$$\mathbf{W}_{cluster} = \frac{1}{C} \sum_{c=1}^C \left( \sum_{i=1}^{N_c} x_{i,c} \mathbf{e}_c \right) \left( \sum_{i=1}^{N_c} x_{i,c} \mathbf{e}_c \right)'$$

$$\Rightarrow \text{Est. Asy. Var}(\mathbf{b}) = C (\mathbf{X}' \mathbf{X})^{-1} \mathbf{W}_{cluster} (\mathbf{X}' \mathbf{X})^{-1}$$

## Los contrastes de hipótesis y la selección del modelo

- Proponiendo un modelo de regresión  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , se puede hacer una pregunta empírica sobre si la especificación del modelo parece ser consistente con los datos a través de los contrastes de hipótesis
  - El enfoque general para contrastar hipótesis es formular un modelo estadístico que contenga las hipótesis como restricciones en sus parámetros
  - Se dice que una teoría económica tiene implicaciones comprobables si implica algunas restricciones contrastables en el modelo. Es decir, si se puede comprobar que uno o más parámetros corresponden a un valor concreto para, así, poder construir una hipótesis nula y contrastarlos

$$\ln I_t = \beta_1 + \beta_2 i_t + \beta_3 \Delta p_t + \beta_4 \ln Y_t + \beta_5 t + \varepsilon_t,$$

$$\ln I_t = \beta_1 + \beta_2(i_t - \Delta p_t) + \beta_4 \ln Y_t + \beta_5 t + \varepsilon_t, \quad \beta_2 + \beta_3 = 0$$

- Por lo tanto, existen modelos que no tienen implicaciones comprobables al no poderse generar una hipótesis nula para los parámetros que se pueda contrastar (ya que una misma variable contiene una combinación de dos o más parámetros, por ejemplo)

$$\ln I_t = \beta_1 + \beta_2 i_t + \beta_3 \Delta p_t + \beta_4 \ln Y_t + \beta_5 t + \varepsilon_t,$$

$$\ln I_t = \beta_1 + \beta_2(i_t - \Delta p_t) + \beta_3 \Delta p_t + \beta_4 \ln Y_t + \beta_5 t + \varepsilon_t.$$

- La descripción de las implicaciones comprobables sugiere que las restricciones contrastables implican que solo algunos de los posibles modelos contenidos en la especificación original serán válidos (consistente con la teoría)
  - El subconjunto de valores que son consistentes con la teoría está contenido dentro de un conjunto de valores no restringido (espacio paramétrico), de modo que los modelos consistentes están anidados y el modelo restringido está contenido en el no restringido. Es decir, a partir del modelo original de la proposición, el otro modelo se puede obtener utilizando restricciones sobre los parámetros de las variables sin eliminar ninguna
  - No obstante, cuando se eliminan variables (restringiendo parámetros a cero), estos modelos no son anidados, aunque se puedan obtener a partir del modelo original (se tienen implicaciones comprobables)
- Para poder contrastar hipótesis, se consideran tres enfoques diferentes: los contrastes de Wald, los contrastes basados en el ajuste y los contrastes de multiplicadores de Lagrange
  - La lógica que sigue cada tipo de procedimiento o contraste es la siguiente:
    - En los contrastes de Wald, la hipótesis expresa que  $\beta$  obedece algunas restricciones, las cuales se pueden expresar como  $c(\beta) = 0$ . Siendo  $b$  el estimador de mínimos cuadrados consistente de  $\beta$ , el contraste de Wald mide que tan cerca  $c(\beta)$  está de cero, basándose en la estimación del modelo no restringido (que tan cerca está el modelo sin restricciones de las restricciones hipotéticas)



- El mejor ajuste posible (mayor  $R^2$ ) se obtiene usando mínimos cuadrados sin poner restricciones, pero añadir restricciones degrada el ajuste del modelo. En los contrastes basados en el ajuste, se mide cuando decrece  $R^2$  cuando se imponen restricciones, comparando el ajuste del modelo restringido con el no restringido
- El contraste de multiplicadores de Lagrange se basa en el modelo restringido, y se basa en un resultado general sobre que, cuando se imponen restricciones, si estas son incorrectas, es posible detectar el fallo con estadísticos medibles. Por lo tanto, se basa en estimadores del modelo restringido (por ejemplo, correlación entre residuos y otros)
- Para desarrollar los procedimientos de contrastes de hipótesis, se asume homocedasticidad y perturbaciones normalmente distribuidas, dado que permiten saber las distribuciones exactas de los estadísticos de contraste
  - No obstante, más adelante se desarrollan resultados alternativos para poder proceder sin las dos suposiciones
  - Es importante distinguir entre la teoría subyacente de los procedimientos de contraste de las mecánicas prácticas de inferencia, basadas en aproximaciones asintóticas y matrices de varianzas y covarianzas robustas
- La hipótesis lineal general de un conjunto de  $J$  restricciones de un modelo de regresión lineal  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  se puede escribir de la siguiente manera:

$$\begin{aligned}
 r_{11}\beta_1 + r_{12}\beta_2 + \cdots + r_{1K}\beta_K &= q_1 \\
 r_{21}\beta_1 + r_{22}\beta_2 + \cdots + r_{2K}\beta_K &= q_2 \\
 &\vdots \\
 r_{J1}\beta_1 + r_{J2}\beta_2 + \cdots + r_{JK}\beta_K &= q_J.
 \end{aligned}$$

- El caso general se puede escribir en notación matricial como  $\mathbf{R}\boldsymbol{\beta} = \mathbf{q}$ , donde cada fila de  $\mathbf{R}$  son los coeficientes de una de las restricciones. Típicamente,  $\mathbf{R}$  solo tendrá unas pocas filas y muchos ceros en cada fila

A set of the coefficients sum to one,  $\beta_2 + \beta_3 + \beta_4 = 1$ ,

$$\mathbf{R} = [0 \quad 1 \quad 1 \quad 1 \quad 0 \quad \cdots]; \mathbf{q} = 1.$$

A subset of the coefficients are all zero,  $\beta_1 = 0$ ,  $\beta_2 = 0$ , and  $\beta_3 = 0$ ,

$$\mathbf{R} = \begin{bmatrix} 1 & 0 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \end{bmatrix} = [\mathbf{I} \mid \mathbf{0}]; \quad \mathbf{q} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

Several linear restrictions,  $\beta_2 + \beta_3 = 1$ ,  $\beta_4 + \beta_6 = 0$ , and  $\beta_5 + \beta_6 = 0$ ,

$$\mathbf{R} = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}; \quad \mathbf{q} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}.$$

All the coefficients in the model except the constant term are zero,

$$\mathbf{R} = [0 \mid \mathbf{I}_{K-1}]; \quad \mathbf{q} = \mathbf{0}.$$

- La matriz  $\mathbf{R}$  tiene  $K$  columnas (para coincidir con las filas de  $\boldsymbol{\beta}$ ),  $J$  filas para un total de  $J$  restricciones, y tiene rango completo de filas, de modo que  $J$  tiene que ser menor o igual a  $K$
- No obstante,  $K \neq J$ , ya que, de no ser así, la matriz sería cuadrada e invertible y  $\boldsymbol{\beta} = \mathbf{R}^{-1}\mathbf{q}$  (no habría problema de estimación o inferencia porque se ha dado un valor para cada parámetro), haciendo que  $J < K$ . La restricción  $\mathbf{R}\boldsymbol{\beta} = \mathbf{q}$  impone  $J$  restricciones en los  $K$  parámetros, haciendo que hayan  $K - J$  parámetros libres
- Para poder extender estos métodos a restricciones no lineales, se plantea una hipótesis no lineal general, que involucra un conjunto  $\mathbf{c}(\boldsymbol{\beta})$  de  $J$  restricciones no lineales posibles de  $\boldsymbol{\beta}$

$$\mathbf{c}(\boldsymbol{\beta}) = \mathbf{q}$$

- La contraparte de los requerimientos vistos para el caso especial en donde  $\mathbf{c}(\boldsymbol{\beta})$  es no lineal es que  $J < K$  y que la matriz jacobiana de estas restricciones (formada por el gradiente de  $\mathbf{c}(\boldsymbol{\beta})$  con respecto a cada uno de los parámetros en  $\boldsymbol{\beta}$ ) tenga rango de fila completo
- En el caso lineal, la jacobiana equivale a  $\mathbf{R}$  y la independencia funcional que se requiere sería equivalente a la independencia lineal
- El contraste de Wald es el procedimiento más utilizado, y normalmente se llama contraste de significación. El principio operador de este procedimiento es ajustar la regresión sin las restricciones y, entonces, evaluar si los resultados parecen, dentro de la variabilidad muestral, ser consistentes con la hipótesis

- El caso más simple es el de contrastar un valor para un solo coeficiente. La distancia de Wald del coeficiente estimado respecto a su valor hipotético es la distancia medida en unidades de desviación estándar, por lo que se puede expresar a través del siguiente estadístico:

$$\begin{cases} H_0: \beta_k = \beta_k^0 \\ H_1: \beta_k \neq \beta_k^0 \end{cases} \Rightarrow W_k = \frac{b_k - \beta_k^0}{\sqrt{\sigma^2 S^{kk}}}$$

where  $S^{kk} = k^{th}$  diagonal elem. of  $(X'X)^{-1}$

- El estadístico  $W_k$  sigue una distribución normal estándar asumiendo que  $E(b_k) = \beta_k^0$ . Si no es igual a  $\beta_k^0$ , el estadístico sigue teniendo una distribución normal, pero la media no será nula. Si, en particular,  $E(b_k) = \beta_k^1$  (diferente de  $\beta_k^0$ ), entonces se obtiene la siguiente igualdad:

$$E[W_k | E(b_k) = \beta_k^1] = \frac{E(b_k) - \beta_k^0}{\sqrt{\sigma^2 S^{kk}}} = \frac{\beta_k^1 - \beta_k^0}{\sqrt{\sigma^2 S^{kk}}}$$

- Con tal de utilizar  $W_k$  para contrastes de hipótesis, la interpretación es que si  $\beta_k$  es igual a  $\beta_k^0$ , entonces  $b_k$  estará cerca de  $\beta_k^0$  (con la distancia medida en unidades de errores estándar). Por lo tanto, la lógica del contraste es rechazar  $H_0$  si  $W_k$  es grande en valor absoluto
- Sin embargo, no se puede usar la medida de Wald propuesta porque  $\sigma^2$  no se sabe, de modo que se utiliza el estimador  $s^2$  para así obtener un estadístico  $t$

$$\begin{cases} H_0: \beta_k = \beta_k^0 \\ H_1: \beta_k \neq \beta_k^0 \end{cases} \Rightarrow t_k = \frac{b_k - \beta_k^0}{\sqrt{s^2 S^{kk}}}$$

where  $S^{kk} = k^{th}$  diagonal elem. of  $(X'X)^{-1}$

- Asumiendo que  $\beta_k$  es igual a  $\beta_k^0$ , el estadístico sigue una distribución  $t$ -Student con  $n - K$  grados de libertad
- A partir de este estadístico, se puede construir un procedimiento de contraste de hipótesis a través de determinar un nivel de confianza (para el intervalo de confianza) y de significación (para el contraste)

$$P\left(t_{n-K, \frac{\alpha}{2}} < t_k < t_{n-K, 1-\frac{\alpha}{2}}\right) = 1 - \alpha$$

$$P(t_k < t_{n-K, 1-\alpha}) = 1 - \alpha \quad \text{or} \quad P(t_k > t_{n-K, \alpha}) = 1 - \alpha$$

- Considerando que se quiere contrastar un conjunto de  $J$  restricciones lineales en la hipótesis nula contra las de la hipótesis alternativa, se puede utilizar el criterio de Wald

$$\begin{cases} H_0: \mathbf{R}\boldsymbol{\beta} - \mathbf{q} = \mathbf{0} \\ H_1: \mathbf{R}\boldsymbol{\beta} - \mathbf{q} \neq \mathbf{0} \end{cases}$$

- Dado el estimador de mínimos cuadrados  $\mathbf{b}$ , el interés se centra en el vector de discrepancia  $\mathbf{Rb} - \mathbf{q} = \mathbf{m}$ . Es improbable que  $\mathbf{m} = \mathbf{0}$ , pero es razonable preguntarse si la desviación se atribuye a la variabilidad muestral o si es significativa
- Como  $\mathbf{b}$  sigue una distribución normal y  $\mathbf{m}$  es una función lineal de  $\mathbf{b}$ ,  $\mathbf{m}$  también se distribuye normalmente. Por lo tanto, si la hipótesis nula es verdad, entonces se obtienen los siguientes resultados:

$$E(\mathbf{m}|\mathbf{X}) = \mathbf{R}E(\mathbf{b}|\mathbf{X}) - \mathbf{q} = \mathbf{R}\boldsymbol{\beta} - \mathbf{q} = \mathbf{0}$$

$$Var(\mathbf{m}|\mathbf{X}) = \mathbf{R}Var(\mathbf{b}|\mathbf{X})\mathbf{R}' = \mathbf{R}[\sigma^2(\mathbf{X}'\mathbf{X})^{-1}]\mathbf{R}'$$

$$\Rightarrow \mathbf{Rb} \sim N(\mathbf{q}, \mathbf{R}[\sigma^2(\mathbf{X}'\mathbf{X})^{-1}]\mathbf{R}')$$

- En consecuencia, es posible basarse en el criterio de Wald para contrastar  $H_0$ . Condicionando en  $\mathbf{X}$ , se puede comprobar que el estadístico  $W$  sigue una distribución  $\chi^2_J$  si la hipótesis nula es cierta, y así ver como valores grandes del estadístico indican que la hipótesis nula es rechazable (la distancia entre  $\mathbf{b}$  y  $\boldsymbol{\beta}$  es muy grande)

$$W = \mathbf{m}'Var(\mathbf{m}|\mathbf{X})^{-1}\mathbf{m} =$$

$$= (\mathbf{Rb} - \mathbf{q})'[\mathbf{R}[\sigma^2(\mathbf{X}'\mathbf{X})^{-1}]\mathbf{R}']^{-1}(\mathbf{Rb} - \mathbf{q}) \sim \chi^2_J$$

- Este estadístico, no obstante, no se puede usar porque depende de  $\sigma^2$ , por lo que se puede obtener un estadístico  $F$  haciendo esa sustitución y dividiendo por los grados de libertad

$$F = \frac{W}{J} \frac{\sigma^2}{s^2} = \frac{(\mathbf{Rb} - \mathbf{q})'[\mathbf{R}[s^2(\mathbf{X}'\mathbf{X})^{-1}]\mathbf{R}']^{-1}(\mathbf{Rb} - \mathbf{q})}{J} \sim F_{J, n-K}$$

- Para comprobar una sola restricción lineal de la forma  $\mathbf{r}'\boldsymbol{\beta} = q$ , el estadístico  $F$  será el siguiente:

$$H_0: r_1\beta_1 + r_2\beta_2 + \dots + r_K\beta_K = \mathbf{r}'\boldsymbol{\beta} = q,$$

$$F[1, n - K] = \frac{(\sum_j r_j b_j - q)^2}{\sum_j \sum_k r_j r_k Est. Var. (b_j, b_k)}$$

- Para contrastar la hipótesis de que el coeficiente  $j$  es igual a un valor particular, entonces  $\mathbf{R}$  tiene una sola fila con un 1 en la posición  $j$ ,  $\mathbf{R}[s^2(\mathbf{X}'\mathbf{X})^{-1}]\mathbf{R}'$  es el elemento diagonal  $j$  de la matriz de varianzas y covarianzas estimada y  $\mathbf{R}\mathbf{b} - \mathbf{q}$  es  $(b_j - q)$ . El estadístico, por tanto, será el siguiente:

$$F[1, n - K] = \frac{(b_j - q)^2}{\text{Est.Var.}(b_j)}$$

- Considerando un enfoque alternativo, se puede plantear una estimación muestral  $\mathbf{r}'\mathbf{b} = \hat{q}$  para  $\mathbf{r}'\boldsymbol{\beta} = q$ . Si  $\hat{q}$  difiere significativamente de  $q$ , se puede concluir que los datos muestrales no son consistentes con la hipótesis nula, de modo que se puede usar el estadístico  $t$

$$t = \frac{\hat{q} - q}{\sqrt{\text{Est.Var}(\hat{q}|\mathbf{X})}}$$

$$\text{where } \text{Est.Var}(\hat{q}|\mathbf{X}) = \mathbf{r}'[s^2(\mathbf{X}'\mathbf{X})^{-1}]\mathbf{r}$$

- Existe una relación útil entre el estadístico  $t$  y el  $F$ , ya que es posible escribir el cuadrado del estadístico  $t$  en términos del estadístico  $F$  cuando hay una sola restricción

$$\begin{aligned} t^2 &= \frac{(\hat{q} - q)^2}{\text{Est.Var}(\hat{q} - q|\mathbf{X})} = \\ &= \frac{(\mathbf{r}'\mathbf{b} - q)'[\mathbf{r}[s^2(\mathbf{X}'\mathbf{X})^{-1}]\mathbf{r}]^{-1}(\mathbf{r}'\mathbf{b} - q)}{1} \sim F_{1, n-K} \end{aligned}$$

- Un enfoque diferente se enfoca en utilizar el ajuste de la regresión, de modo que uno se puede preguntar si escoger valores alternativos para los parámetros permite obtener una pérdida significativa de ajuste
  - El vector de coeficientes  $\mathbf{b}$  es escogido para minimizar la suma de desviaciones al cuadrado  $\mathbf{e}'\mathbf{e}$ 
    - Debido a que  $R^2$  iguala  $1 - \mathbf{e}'\mathbf{e}/\mathbf{y}'\mathbf{M}^0\mathbf{y}$  y  $\mathbf{y}'\mathbf{M}^0\mathbf{y}$  es una constante que no involucra a  $\mathbf{b}$ , si el modelo contiene un término constante,  $\mathbf{b}$  se escoge para maximizar  $R^2$
    - En consecuencia, se puede construir un estadístico de contraste que se base en comparar el  $R^2$  entre dos regresiones

- Suponiendo que se impone explícitamente las restricciones de la hipótesis lineal general en la regresión, es posible obtener un estimador restringido  $\mathbf{b}_*$  a través de utilizar la función lagrangiana

$$\min_{\mathbf{b}_0} \mathbf{e}'\mathbf{e} = (\mathbf{y} - \mathbf{X}\mathbf{b}_0)'(\mathbf{y} - \mathbf{X}\mathbf{b}_0) \quad s. t. \quad \mathbf{R}\mathbf{b}_0 = \mathbf{q}$$

$$\min_{\mathbf{b}_0} L^*(\mathbf{b}_0, \boldsymbol{\lambda}) = (\mathbf{y} - \mathbf{X}\mathbf{b}_0)'(\mathbf{y} - \mathbf{X}\mathbf{b}_0) + 2\boldsymbol{\lambda}'(\mathbf{R}\mathbf{b}_0 - \mathbf{q})$$

$$\frac{\partial L^*}{\partial \mathbf{b}_*} = -2\mathbf{X}'(\mathbf{y} - \mathbf{X}\mathbf{b}_*) + 2\mathbf{R}'\boldsymbol{\lambda}_* = \mathbf{0} \quad \frac{\partial L^*}{\partial \boldsymbol{\lambda}_*} = 2(\mathbf{R}\mathbf{b}_0 - \mathbf{q}) = \mathbf{0}$$

- Dividiendo entre 2 (por eso se ha multiplicado  $\boldsymbol{\lambda}$  por 2) y expandiendo los términos, se puede obtener la siguiente ecuación matricial:

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{R}' \\ \mathbf{R} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{b}_* \\ \boldsymbol{\lambda}_* \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{q} \end{bmatrix}$$

- Asumiendo que la matriz partida es invertible, se tiene la siguiente solución:

$$\begin{bmatrix} \mathbf{b}_* \\ \boldsymbol{\lambda}_* \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{R}' \\ \mathbf{R} & \mathbf{0} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{q} \end{bmatrix} = \mathbf{A}^{-1}\mathbf{d}.$$

- Si además  $\mathbf{X}'\mathbf{X}$  es invertible, entonces se tiene soluciones explícitas:

$$\mathbf{b}_* = \mathbf{b} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'[\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1}(\mathbf{R}\mathbf{b} - \mathbf{q}) = \mathbf{b} - \mathbf{C}\mathbf{m}$$

$$\text{where } \mathbf{C} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'[\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1}$$

$$\boldsymbol{\lambda}_* = [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1}(\mathbf{R}\mathbf{b} - \mathbf{q})$$

- Como  $\boldsymbol{\lambda}_*$  involucra el vector de discrepancia  $\mathbf{m} = \mathbf{R}\mathbf{b} - \mathbf{q}$ , si se satisface la restricción, entonces  $\mathbf{b}_* = \mathbf{b}$  y el estimador no tendría ningún sesgo si las hipótesis son ciertas
- La solución obtenida por mínimos cuadrados restringidos cumple unas propiedades interesantes:

- El estimador obtenido es lineal, por lo que se cumple  $\mathbf{R}\mathbf{b}_* = \mathbf{q}$
- Bajo normalidad, el estimador obtenido es igual al estimador de máxima verosimilitud  $\mathbf{b}_{MLE} = \mathbf{b}_*$ . Por lo tanto, este sigue una distribución normal multivariante y es consistente y asintóticamente eficiente

- La matriz de varianzas y covarianzas para  $\mathbf{b}_*$  es simplemente  $\sigma^2$  veces el bloque superior de la matriz partida  $\mathbf{A}^{-1}$ . Si  $\mathbf{X}'\mathbf{X}$  es invertible, se puede hacer una formulación explícita:

$$\begin{aligned} \text{Var}(\mathbf{b}_*|\mathbf{X}) &= \text{Var}(\mathbf{b}|\mathbf{X}) - \text{nonnegative def. matrix} = \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1} - \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'[\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1}\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \end{aligned}$$

- Una interpretación para esta reducción en la varianza es como el valor de la información contenida en las restricciones. La matriz  $\text{Var}(\mathbf{b}_*|\mathbf{X})$  es más pequeña que  $\text{Var}(\mathbf{b}|\mathbf{X})$  aunque las restricciones sean incorrectas (dado que añadir parámetros hace que aumente la variabilidad de las estimaciones), por lo que el estimador de la varianza es la varianza mínima entre los estimadores lineales y sin sesgo que cumplen el conjunto de restricciones
- Para desarrollar un contraste basado en el estimador de mínimos cuadrados restringido, se puede utilizar el resultado anterior sobre la relación entre el cuadrado del estadístico  $t$  y el estadístico  $F$  para  $J$  restricciones
  - El ajuste de los coeficientes de mínimos cuadrados restringidos no puede ser mejores que la solución no restringida (debido a que se considera un espacio paramétrico más pequeño), lo cual se puede demostrar a través de la suma cuadrada de desviaciones:

$$\begin{aligned} \mathbf{e}_* &= \mathbf{y} - \mathbf{X}\mathbf{b} - \mathbf{X}(\mathbf{b}_* - \mathbf{b}) = \mathbf{e} - \mathbf{X}(\mathbf{b}_* - \mathbf{b}) \\ \Rightarrow \mathbf{e}_*' \mathbf{e}_* &= (\mathbf{e}' - (\mathbf{b}_* - \mathbf{b})' \mathbf{X}')(\mathbf{e} - \mathbf{X}(\mathbf{b}_* - \mathbf{b})) = \\ &= \mathbf{e}' \mathbf{e} - \mathbf{e}' \mathbf{X}(\mathbf{b}_* - \mathbf{b}) - (\mathbf{b}_* - \mathbf{b})' \mathbf{X}' \mathbf{e} + (\mathbf{b}_* - \mathbf{b})' \mathbf{X}' \mathbf{X}(\mathbf{b}_* - \mathbf{b}) = \\ &= \mathbf{e}' \mathbf{e} + (\mathbf{b}_* - \mathbf{b})' \mathbf{X}' \mathbf{X}(\mathbf{b}_* - \mathbf{b}) \geq \mathbf{e}' \mathbf{e} \\ \Rightarrow \mathbf{e}_*' \mathbf{e}_* - \mathbf{e}' \mathbf{e} &= (\mathbf{b}_* - \mathbf{b})' \mathbf{X}' \mathbf{X}(\mathbf{b}_* - \mathbf{b}) = \\ &= (\mathbf{R}\mathbf{b} - \mathbf{q})' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} (\mathbf{R}\mathbf{b} - \mathbf{q}) \end{aligned}$$

- La expresión obtenida para la pérdida de ajuste  $\mathbf{e}_*' \mathbf{e}_* - \mathbf{e}' \mathbf{e}$  aparece en el numerador el estadístico  $F$ , por lo que se puede sustituir y así encontrar una expresión equivalente

$$\frac{(\mathbf{e}_*' \mathbf{e}_* - \mathbf{e}' \mathbf{e})/J}{\mathbf{e}' \mathbf{e}/(n - K)} \sim F_{J, n-K}$$

- Finalmente, dividiendo el numerador y el denominador por la desviación cuadrada de la media muestral  $\sum_i (y_i - \bar{y})^2$  se puede obtener el estadístico en términos de  $R^2$

$$\frac{(R^2 - R_*^2)/J}{(1 - R^2)/(n - K)} \sim F_{J, n-K}$$

- Una de las preguntas más importantes es normalmente ver si la ecuación de la regresión es significativa en su totalidad. Este es un contraste conjunto de las hipótesis de que todos los coeficientes (excepto el término constante) sean nulos
  - Si todos los coeficientes son nulos, entonces el coeficiente de determinación es nulo también, de modo que se puede basar el contraste en la hipótesis sobre el valor de  $R_*^2 = 0$ . Por lo tanto, el estadístico será el siguiente:

$$\frac{R^2/(K - 1)}{(1 - R^2)/(n - K)} \sim F_{J, n-K}$$

- Valores grandes de  $F$  evidencian que la hipótesis nula se puede rechazar, y estos son inducidos por valores grandes de  $R^2$ , por lo que la lógica del contraste es que  $F$  mide la pérdida de ajuste que resulta en imponer la restricción de que  $\beta = 0$

## Las variables binarias o categóricas y los experimentos

- Una de las herramientas más útiles en el análisis de regresión son las variables binarias, las cuales son convenientes para crear desplazamientos discretos de la función en un modelo de regresión
  - Las variables binarias toman el valor de 1 para algunas observaciones para indicar la presencia de un efecto o la membresía a un grupo y toma el valor 0 para el resto de observaciones
    - Las variables binarias normalmente se usan en ecuaciones de regresión que contienen otras variables cuantitativas
    - Una variable binaria que toma el valor de 1 solo para una única observación tiene el efecto de eliminar esa observación del cálculo de los parámetros de mínimos cuadrados y el estimador de la varianza (pero no del  $R^2$ )
  - El grupo de observaciones a las que se les atribuye el 0 se denomina grupo de referencia o *benchmark group*, y la interpretación del coeficiente se hace relativa a este grupo (como una diferencia)



- Por lo tanto, en un modelo con una sola variable binaria, el cambio producido por la variable binaria  $d$  se debe interpretar como el efecto de pertenecer a un grupo o que ocurra un evento (las observaciones para las cuales  $d = 1$ ) frente a no pertenecer o a que no ocurra (grupo de referencia)

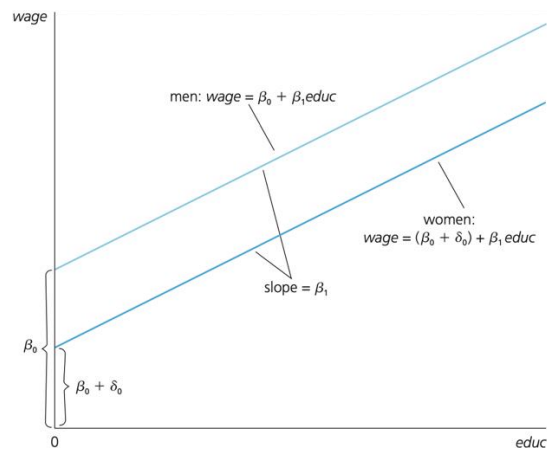
$$y_i = \beta_0 + \beta_1 x_i + \beta_2 d_i + \varepsilon_i$$

$$\Rightarrow \frac{\partial E(y_i | x_i, d_i)}{\partial d_i} = \beta_2 = E(y_i | x_i, d_i = 1) - E(y_i | x_i, d_i = 0)$$

- Como se puede deducir matemáticamente, el efecto de un cambio en  $d$  solo depende de  $d$  y no de cualquier otra variable (al mantenerse constante)
- En un modelo de este tipo, esto se traduce en un desplazamiento del intercepto y, por tanto, de la función de regresión

$$E(y_i | x_i, d_i = 1) = (\beta_0 + \beta_2) + \beta_1 x_i + \varepsilon_i$$

$$E(y_i | x_i, d_i = 0) = \beta_0 + \beta_1 x_i + \varepsilon_i$$



- Debido al teorema de Frisch-Waugh-Lovell, la regresión múltiple sin coeficiente constante es equivalente a hacer una regresión de la variable dependiente y de las variables independientes sobre las variables binarias y hacer una regresión de los residuos de la primera sobre los residuos de la segunda
- Por lo tanto, se puede ver que el vector de coeficientes de las variables cuantitativas seguirá siendo el mismo y se produce el mismo vector de residuos que cuando se incluye el conjunto de variables binarias en la regresión. Esto ocurre porque, si hay  $K - r$  variables continuas y  $r$  variables binarias, entonces la matriz  $X$

tiene un tamaño  $n \times K$  y, debido a que se añaden todas las variables binarias (para que no haya problemas de identificación), entonces hay columnas de ceros y unos que sumadas permiten obtener una columna de unos, y por tanto es como tener un parámetro constante

- Los valores de la variable binaria son arbitrarios, dado que se puede trabajar con cualquier otro par de números al ser combinaciones lineales de 1 y 0, de modo que la codificación no es relevante para la interpretación

$$d = \begin{cases} 1 & \text{if event A} \\ 0 & \text{if event B} \end{cases} \quad g = \begin{cases} a & \text{if event A} \\ b & \text{if event B} \end{cases} \Rightarrow g = ad + b$$

- Cuando cambia la codificación de la variable binaria también cambia el coeficiente (debido a la magnitud diferente de los datos). No obstante, como  $y$  no cambia, entonces se puede plantear la siguiente ecuación para cambiar la codificación

$$E(y_i|x_i, d_i) = \beta_0 + \beta_1 x_i + \beta_2 d_i = \beta_0^* + \beta_1 x_i + \beta_2^* g_i = E(y|x, g)$$

$$E(y_i|x_i, d_i = 1) - E(y_i|x_i, d_i = 0) = E(y|x, d = a) - E(y|x, d = b)$$

- Obviamente, cuando se cambia el grupo base, la codificación es la inversa y los coeficientes serán exactamente los contrarios

$$\beta_2 = E(y_i|x_i, d_i = 1) - E(y_i|x_i, d_i = 0)$$

$$\beta_2^* = E(y_i|x_i, d_i = 0) - E(y_i|x_i, d_i = 1)$$

$$\Rightarrow \beta_2 = \beta_2^* = -\beta_2$$

- Finalmente, como  $y$  no cambia, aunque se cambie la codificación, los modelos de regresión tienen que dar el mismo resultado, por lo que se puede ver que el coeficiente constante cambiará cuando la codificación cambie

$$E(y_i|x_i, d_i) = \beta_0 + \beta_1 x_i + \beta_2 d_i = \beta_0^* + \beta_1 x_i + \beta_2^* g_i = E(y_i|x_i, g_i)$$

$$E(y_i|x_i, d_i = 1) \Rightarrow \beta_0^* = \beta_0 - \beta_2 + \beta_2^* a$$

$$E(y_i|x_i, d_i = 0) \Rightarrow \beta_0^* = \beta_0 + \beta_2^* b$$

$$\beta_2^* = -\frac{\beta_2}{b - a}$$

- En el caso en el que se incluya más de una variable binaria, la interpretación de los coeficientes sigue siendo la misma (la diferencia con el grupo de referencia)

y el efecto gráfico de estas también. No obstante, hay algunas diferencias con el caso anterior

- En la mayoría de aplicaciones, las variables binarias se utilizan para factores puramente cualitativos, pero hay casos en los que se utilizan para representar niveles de un factor subyacente que podría haber sido medido si fuera posible (en donde hay umbrales o categorías)

- En este caso, se pueden dividir los diferentes umbrales o categorías a través de varias variables binarias

$$income = \beta_1 + \beta_2 age + \delta_B B + \delta_M M + \delta_P P + \varepsilon$$

- Sin embargo, la definición del grupo de referencia tendrá un efecto en la función de regresión, dado que se puede hacer que 1 sea pertenecer a una categoría concreta o que sea la categoría más alta (o baja, depende del criterio)

$$\text{High school: } E[income | age, HS] = \beta_1 + \beta_2 age,$$

$$\text{Bachelor's: } E[income | age, B] = \beta_1 + \beta_2 age + \delta_B,$$

$$\text{Master's: } E[income | age, M] = \beta_1 + \beta_2 age + \delta_M,$$

$$\text{Ph.D.: } E[income | age, P] = \beta_1 + \beta_2 age + \delta_P.$$

$$\text{High school: } E[income | age, HS] = \beta_1 + \beta_2 age,$$

$$\text{Bachelor's: } E[income | age, B] = \beta_1 + \beta_2 age + \delta_B,$$

$$\text{Master's: } E[income | age, M] = \beta_1 + \beta_2 age + \delta_B + \delta_M,$$

$$\text{Ph.D.: } E[income | age, P] = \beta_1 + \beta_2 age + \delta_B + \delta_M + \delta_P.$$

- Eso, a su vez, cambia la interpretación de estos coeficientes, dado que, en el primer caso, el coeficiente representa la diferencia entre pertenecer a la categoría y pertenecer al grupo base, y en el segundo caso, representa el efecto marginal de pertenecer a una categoría (dado que también se puede pertenecer a otras)
- Como se puede observar, las observaciones pertenecen al grupo de referencia cuando todas las variables binarias incluidas son nulas
- Cuando hay varias categorías dentro de las observaciones, el uso de variables binarias es necesario
- No obstante, el modelo de regresión no se incluyen variables binarias para todas las categorías debido a que, de ser así, entonces la suma de estas sería igual a 1 y eso haría que se

reprodujera el término constante (las columnas de las variables binarias son combinación lineal de la columna de 1 del coeficiente constante), de modo que no habría rango máximo (habría multicolinealidad perfecta)

- A este hecho se le denomina trampa de las variables binarias, y para poder evitar este problema, es necesario no incluir una de las variables binarias o desestacionalizar los datos a través de quitar el término constante
- El caso en que varios grupos de variables binarias se necesita es casi el mismo que cuando solo se necesita un grupo, pero con una importante excepción
  - Como antes, es necesario quitar una de las variables binarias en el grupo para que este no sea combinación lineal del regresor constante del término constante
  - No obstante, cuando hay  $r$  grupos de  $k$  variables binarias en una regresión, la suma de las variables binarias para cada uno de los  $r$  grupos es 1, por lo que hay que omitir una de las variables en cada uno de los  $r$  grupos (de modo que se incluya  $k - 1$  en cada grupo)
- Como los valores de las variables binarias son arbitrarios, la codificación tampoco es relevante para la interpretación en este caso
  - Cuando cambia la codificación de las variables binarias también cambia el coeficiente (debido a la magnitud diferente de los datos). No obstante, como  $y$  no cambia, entonces se puede plantear la ecuación anteriormente vista para cambiar la codificación
  - No obstante, cuando en un grupo de variables binarias se cambia el grupo de referencia, entonces no solo se modifica uno de los coeficientes, sino todos los de las variables binarias del mismo grupo. Por tanto, como cada coeficiente era la diferencia con respecto al grupo de referencia anterior, se pueden utilizar esas diferencias para poder ver que los nuevos coeficientes serán las diferencias de los coeficientes originales
  - Obviamente, el coeficiente constante también cambiará y se puede plantear una ecuación como la anteriormente vista para poder obtener el nuevo coeficiente constante

## La no linealidad en las variables

- Las variables binarias son particularmente útiles en el modelo de regresión log-lineal para poder saber cual es el cambio porcentual al pasar del grupo de referencia a otro

$$100\% \left( \frac{\partial E(y_i | x_i, d_i)}{\partial d} \right) = 100\%(\beta_2) = 100\%(e^{\beta_2} - 1)$$

- La segunda expresión equivalente para el efecto marginal del cambio en  $d$  se puede obtener a través de la discretización del cambio y de expresar la ecuación en términos de exponenciales

$$\begin{aligned} 100\% \left( \frac{\Delta E(y_i | x_i, d_i)}{\Delta d_i} \right) &= 100\% \left( \frac{E(y_i | x_i, d_i = 1) - E(y_i | x_i, d_i = 0)}{E(y_i | x_i, d_i = 0)} \right) = \\ &= 100\% \left( \frac{e^{\beta_0 + \beta_1 x_i + \beta_2} E(e^\varepsilon) - e^{\beta_0 + \beta_1 x_i} E(e^\varepsilon)}{e^{\beta_0 + \beta_1 x_i} E(e^\varepsilon)} \right) = 100\%(e^{\beta_2} - 1) \end{aligned}$$

## La evaluación de los efectos de tratamientos

- La principal aplicación de los modelos de selección y de endogeneidad es la evaluación de efectos de tratamientos, que se centra en analizar el efecto de la participación en un tratamiento  $T$  en una variable de resultado  $y$ . El objetivo principal de analizar el tratamiento o la intervención es saber el efecto del tratamiento en los tratados, pero se presentan varios problemas
  - El analista tiene el riesgo de atribuir al tratamiento efectos causales que deberían ser atribuidos a factores que motivan el tratamiento y el resultado a la vez (endogeneidad del tratamiento)
    - Los ejemplos más comunes suelen ser la raza, la edad y otras características de los individuos
  - Normalmente falta un contrafactual, dado que hay efectos que se quieren medir pero que, por la naturaleza de los individuos, no se puede porque solo se puede observar un estado
    - Un ejemplo claro son los experimentos sobre el efecto de la educación, en los que se querría medir el efecto en la vida de un individuo con educación y sin ella, pero solo se puede observar uno de los dos estados

- El modelo causal de Rubin proporciona un marco útil para el análisis. En este, cada individuo de una población tiene un resultado potencial  $y$  y puede estar expuesto a un tratamiento  $C$  ( $C_i$  es el indicador de si se ha recibido el tratamiento o no). Por lo tanto, hay dos resultados potenciales:  $y_i|(C_i = 1) = y_{i1}$  y  $y_i|(C_i = 0) = y_{i0}$
- El efecto medio del tratamiento o *average treatment effect* (ATE) es la diferencia media entre los resultados potenciales para toda la población. No obstante, el individuo existe en uno de los dos estados, por lo que no es posible estimar el ATE y se pone más interés en el ATE de los tratados (ATET)

$$ATE = E(y_{i1} - y_{i0}) \quad ATET = E(y_{i1} - y_{i0}|C_i = 1)$$

- El reto de la investigación actualmente está en encontrar estimadores robustos que no se apoyen mucho en suposiciones frágiles como la identificación por la forma funcional o por exclusión de restricciones. Por ello,  $x$  será la información exógena con la que se tiene que lidiar en el problema de estimación, y se hacen suposiciones sobre este vector de variables como las siguientes:

- Recibir el tratamiento  $C_i$  no depende del resultado una vez que los efectos de  $x$  se tienen en cuenta. Por lo que, si la asignación de individuos al grupo de tratamiento es aleatoria, entonces se tendría que omitir el efecto de  $x$  en la suposición de la independencia condicional
- La suposición de independencia condicional se extiende a enfoques de regresión con la suposición de media condicional o *conditional mean assumption*, de modo que los resultados de los individuos que no reciben el tratamiento no afectan a la participación

$$E(y_{i0}|x_i, C_i = 1) = E(y_{i0}|x_i, C_i = 0) = E(y_{i0}|x_i)$$

- La suposición sobre la distribución de resultados potenciales expresa que el modelo que se usa para los resultados potenciales de los tratados y los no tratados es el mismo. En el contexto de una regresión eso quiere decir que la misma regresión aplica para ambos grupos y que las perturbaciones no tienen correlación con  $T$  ( $T$  es exógena). Esta suposición es fuerte, pero permite evitar la endogeneidad

$$f(y|x, T = 1) = f(y|x, T = 0)$$

- La suposición de superposición expresa que para cualquier valor  $x$ ,  $0 < P(C_i = 1|x) < 1$ , de modo que las desigualdades

estrictas expresan que la población contiene una mezcla de individuos tratados y no tratados para cualquier  $x$ . Esta suposición permite suponer que, de media, para cada individuo tratado habrá otro individuo no tratado con las mismas características (útil para la regresión)

- Para poder analizar el efecto de un tratamiento sobre una respuesta, se puede formular un modelo de regresión con una variable binaria  $C$  con valor 1 si se recibe el tratamiento y 0 si no y otras variables  $x$

$$y_i = x_i' \beta + \delta C_i + \varepsilon_i$$

- En este caso,  $\delta$  es el parámetro de desplazamiento de la función de regresión y mide el impacto del tratamiento en la respuesta de los individuos de la muestra (condicionado a las  $x$ )
- Cuando el tratamiento o el análisis de este ocurre a lo largo del tiempo, entonces la variable binaria será  $C_t = 0$  en el primer periodo y  $C_t = 1$  en el segundo, y el efecto  $\delta$  medirá el cambio en la variable de respuesta antes y después del tratamiento
  - Como la suposición de que  $x$  se mantiene constante entre periodos debilita la comparación, una estrategia es incluir un grupo de observaciones de control que no reciban tratamiento, de modo que el cambio en el grupo de tratamiento se puede comparar al cambio del grupo no tratado bajo la suposición de que la diferencia se debe al tratamiento
  - El impacto en este caso se mide con la diferencia de las diferencias anteriormente descritas, de modo que el estimado del parámetro de  $C$  será el estimador de diferencias en diferencias

$$\delta = E(\Delta y_i | \Delta C_i = 1) - E(\Delta y_i | C_i = 0) =$$

$$= [E(y_{i1} | C_i = 1) - E(y_{i0} | C_i = 1)] - [E(y_{i1} | C_i = 0) - E(y_{i0} | C_i = 0)]$$

- Cuando el tratamiento es el resultado de un evento que ocurre completamente fuera del contexto del estudio el análisis se suele denominar experimento natural
  - Cuando esto ocurre, entonces se puede considerar que la asignación del tratamiento y la selección de individuos es aleatoria y, por tanto, las observaciones son independientes e idénticamente distribuidas
  - Uno de los problemas más importantes es la medición de los efectos de los tratamientos cuando  $C$  resulta de una decisión de

participación individual, de modo que los individuos pueden escoger voluntariamente si recibir el tratamiento y puede haber sesgo en el efecto del tratamiento

- Con el enfoque de experimento natural o cuasi-experimento, en el que se intenta replicar la asignación aleatoria, la aleatoriedad viene introducida por características del individuo en  $\mathbf{x}$ . En este caso, se asume que hay independencia condicional
- El modelo básico de selección anterior (en el modelo de Heckman, por ejemplo) se ha podido extender en muchas direcciones, y una de las aplicaciones más interesantes de este es su uso para medir los efectos de los tratamientos y de programas
  - En un modelo en el que se incluye la variable binaria del tratamiento  $C_i$  y las características  $\mathbf{x}_i$  de los individuos, el coeficiente  $\delta$  puede no estimar correctamente el efecto del tratamiento (aunque el resto de la regresión sea especificada correctamente)

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \delta C_i + \varepsilon_i$$

- Esto ocurre porque puede haber problemas de selección muestral, en donde los individuos deciden o son más propensos a participar en el experimento o tratamiento por unas características concretas
- De este modo, el efecto de los tratamientos se sobreestima o se subestima
- Debido a este problema, se puede modelar como un modelo de dos ecuaciones en donde hay truncación incidental

$$C_i^* = \mathbf{w}_i' \boldsymbol{\gamma} + u_i \quad C_i^* = \begin{cases} 1 & \text{if } C_i^* > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned} E(y_i | C_i = 1, \mathbf{x}_i, \mathbf{w}_i) &= \mathbf{x}_i' \boldsymbol{\beta} + \delta + E(\varepsilon_i | C_i = 1, \mathbf{x}_i, \mathbf{w}_i) = \\ &= \mathbf{x}_i' \boldsymbol{\beta} + \delta + \rho \sigma_\varepsilon \lambda(\mathbf{w}_i' \boldsymbol{\gamma}) \end{aligned}$$

- Considerando un modelo de regresión para analizar los efectos del tratamiento en donde hay dos periodos (antes y después del tratamiento), es posible analizar los cambios a través del estimador de diferencias en diferencias

$$y_{it} = \mathbf{x}_{it}' \boldsymbol{\beta} + \gamma C_i + \theta_t + u_i + \varepsilon_{it}$$

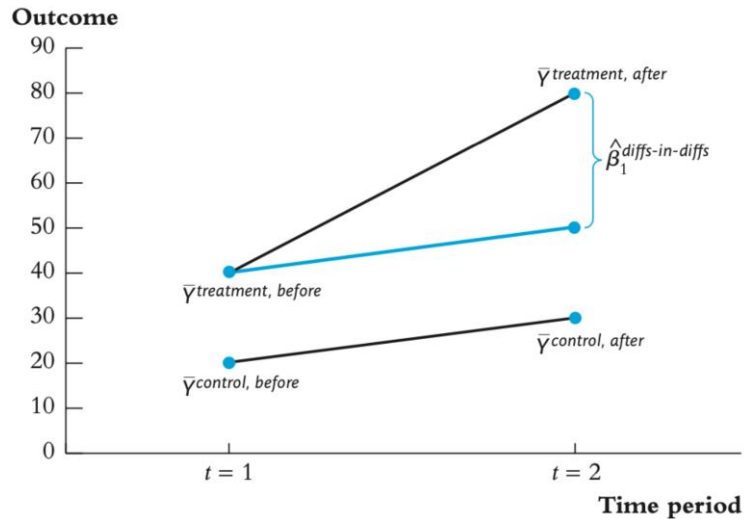


- En este caso,  $x'_{it}$  son las variables de control,  $C_i$  es la variable binaria para indicar el tratamiento,  $u_i$  son las características no observables del individuo y  $\theta_t$  son los efectos temporales
- Considerando la presencia de variables de control en la regresión, la suposición de independencia del término de error y los regresores se cambia por la de independencia de media, en donde solo se pide que la variable del tratamiento sea independiente del error, pero no las variables de control

$$E(\varepsilon_i | C_i, x_{it}) = E(\varepsilon_i | x_{it})$$

- Como hay dos periodos, lo más natural es analizar los cambios cuando  $\Delta C_i = 1$  para los individuos tratados y  $\Delta C_i = 0$  para los no tratados (la diferencia entre ambos periodos). Como se puede ver, esta diferencia elimina los factores individuales no observados (son los mismos en el tiempo)

$$\Delta y_i = y_{i1} - y_{i0} = (\Delta x_{it})' \beta + \gamma \Delta C_i + (\theta_1 - \theta_0) + \Delta \varepsilon_{it}$$



- En la ausencia de variables de control  $x'_{it}$  de las características de los individuos, o asumiendo que no cambian en el tiempo, el estimador del efecto del tratamiento será el estimador de diferencias en diferencias

$$\begin{aligned} \gamma &= E(\Delta y_i | \Delta C_i = 1) - E(\Delta y_i | C_i = 0) = \\ &= [E(y_{i1} | C_i = 1) - E(y_{i0} | C_i = 1)] - [E(y_{i1} | C_i = 0) - E(y_{i0} | C_i = 0)] \\ &\Rightarrow \hat{\gamma} = (\overline{\Delta y} | \Delta C_i = 1) - (\overline{\Delta y} | C_i = 0) \\ &= [(\overline{y_1} | C = 1) - (\overline{y_0} | C = 1)] - [(\overline{y_1} | C = 0) - (\overline{y_0} | C = 0)] \end{aligned}$$

- Aunque el problema se simplifica de manera considerable, el uso de diferencias hace que no se puede discernir lo que indujo al cambio (el tratamiento u otras características en  $x_{it}$ ). Además, las suposiciones son más fuertes de lo que gustaría un analista
- El contrafactual implícito es una observación de lo que la variable dependiente del individuo tratado sería si no fueran tratados (la diferencia para los no tratados), pero los individuos solo pueden estar en un estado, por lo que puede haber sesgo

- El efecto medio del tratamiento en la población sería el ATE, que sería el efecto si se escogiera un individuo aleatoriamente de la población entera

$$ATE = E(y_1 - y_0)$$

- No obstante, puede ser más interesante medir el efecto medio del tratamiento en los tratados, el cual sería el ATET. La dificultad con esto es que el contrafactual  $E(y_0|C = 1)$  no se puede medir

$$ATET = E(y_1 - y_0|C = 1)$$

- Si el tratamiento se asigna de manera totalmente aleatoria, entonces  $E(y_j|C = 1) = E(y_j|C = 0) = E(y_j|C = j)$  para  $j = 0, 1$ , por lo que el resultado potencial no depende del tratamiento y eso quiere decir que se puede medir realmente el efecto del tratamiento para aquellos que lo han recibido comparado con los que no

$$ATE = E(y_1|C = 1) - E(y_0|C = 0)$$

- Sin embargo, si el tratamiento no se asigna de manera exógena o aleatoria, entonces es necesario utilizar otros enfoques como la función de control o variables instrumentales
- Cuando se intenta comparar la media o la esperanza de los resultados de dos poblaciones diferentes, estas se tienen que escalar para poder ser comparables
  - Para poder escalar las métricas, lo que se tiene que hacer es dividir la media por el error estándar de su respectiva población, de modo que la varianza es unitaria y permite comparar ambas métricas (que tendrán la misma variación)

$$\widehat{ATE} = \frac{\bar{y}_1}{\sqrt{S_1^2/n}} - \frac{\bar{y}_0}{\sqrt{S_0^2/n}}$$

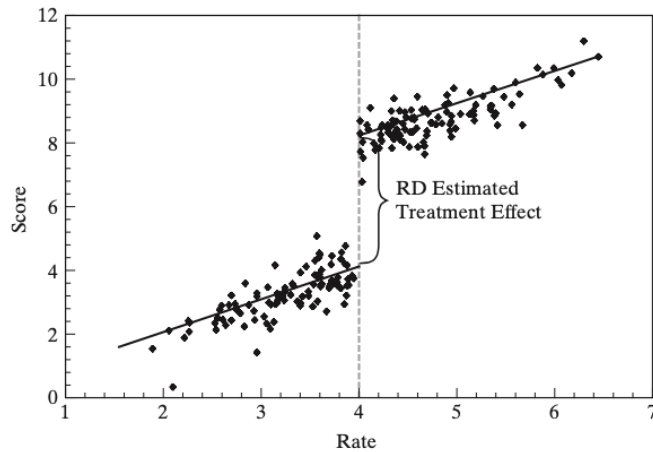
- Del mismo modo, para comparar estimaciones de los coeficientes de dos poblaciones diferentes, es necesario escalar los coeficientes y los errores estándar de los coeficientes estimados dividiendo entre el error estándar de la variables dependiente y

$$group\ effect = \frac{\hat{\beta}_1}{\sqrt{\frac{S_y^2}{n}}}$$

$$group\ effect\ std.\ error = \frac{S_{\hat{\beta}_1}}{\sqrt{S_y^2/n}}$$

- Hay muchas situaciones en las que no hay posibilidad de asignar aleatoriamente individuos a los tratamientos, sobre todo en entornos que estudian las ciencias sociales. Con tal de superar este problema a la hora de la estimación, el diseño de regresión discontinua presenta una alternativa
  - Las condiciones bajo las que este enfoque es efectivo son que  $y$  sea una variable continua, que la variable dependiente varíe suavemente con una variable de asignación  $A$  y que el tratamiento se asigne de forma brusca en función del valor de  $A$ 
    - Se puede considerar que la variable de asignación  $A$  es una característica del individuo que determina su participación en el tratamiento
    - Específicamente, la tercera condición indica que  $C = 1_{A>A^*}$ , donde  $A^*$  es un umbral fijo. Un diseño borroso o *fuzzy design*, en cambio, se basa en que  $P(C = 1|A) = F(A)$ , pero los problemas de identificación con este diseño son más complicados que con el diseño brusco
    - Por lo tanto, se asume que  $y = f(A, C) + \varepsilon$ . Este método, sin embargo, requiere que  $E(\varepsilon|A, C) = E(\varepsilon|A)$ , de modo que la variable de asignación es exógena al experimento
  - Por lo tanto, un modelo de regresión discontinua puede ser el siguiente, en donde  $\alpha$  es el efecto del tratamiento a estimar:

$$y = f(A) + \alpha C + \varepsilon$$



- La especificación de  $f(A)$  puede ser problemática, dado que asumir una función lineal cuando algo es muy general puede sesgar la estimación de  $\alpha$ . Para estos casos se usan métodos no paramétricos, haciendo que el analista pueda utilizar observaciones que están más distantes del umbral
- La identificación del efecto de tratamiento comienza con la suposición de que  $f(A)$  es continua en  $A^*$ , de modo que el efecto del tratamiento puede estimarse por la diferencia de los individuos que están cerca del valor de umbral  $A^*$

$$\lim_{A \rightarrow A^* -} f(A) = \lim_{A \rightarrow A^* +} f(A) = f(A^*)$$

$$\Rightarrow \lim_{A \rightarrow A^* +} E(y|A) - \lim_{A \rightarrow A^* -} E(y|A) =$$

$$= f(A^*) + \alpha + \lim_{A \rightarrow A^* +} E(\varepsilon|A) - f(A^*) - \lim_{A \rightarrow A^* -} E(\varepsilon|A) = \alpha$$

- El efecto causal puede variar de un miembro de la población a otro, y este puede variar por variables observables  $x$  o por variables no observables (cuando se tiene una población heterogénea)
  - Cuando la población es heterogénea, se dan variaciones no observables en el efecto causal de cada individuo, de modo que cada individuo tiene su propio efecto causal

$$y_i = x_i' \beta + \delta_i C_i + \varepsilon_i$$

- Como el coeficiente varía de un individuo a otro, y los individuos se seleccionan aleatoriamente, esta es una variable aleatoria y no una constante como se asume anteriormente. El efecto medio del tratamiento será la esperanza del efecto causal para un individuo seleccionado aleatoriamente de la población a estudiar

$$ATE = E(y_1 - y_0) = E(\delta_i)$$

- Si el tratamiento se asigna aleatoriamente, entonces no está correlacionada con los resultados potenciales, y, por tanto, tampoco con  $\delta_i$  ni con el error
- Cuando hay heterogeneidad en el efecto causal y el tratamiento se asigna aleatoriamente, el estimador de mínimos cuadrados ordinarios es consistente para estimar el efecto causal medio

$$\begin{aligned} \text{plim}_{p \rightarrow \infty} \delta_i &= \frac{\text{Cov}(\mathbf{x}_i' \boldsymbol{\beta} + \delta_i C_i + \varepsilon_i, C_i)}{\text{Var}(C_i | \mathbf{X})} = \frac{\text{Cov}(\delta_i C_i, C_i)}{\text{Var}(C_i | \mathbf{X})} = \frac{E(\delta_i C_i C_i | \mathbf{X})}{\text{Var}(C_i | \mathbf{X})} = \\ &= \frac{E(\delta_i | \mathbf{X}) E(C_i C_i | \mathbf{X})}{\text{Var}(C_i | \mathbf{X})} = \frac{E(\delta_i | \mathbf{X}) \text{Var}(C_i | \mathbf{X})}{\text{Var}(C_i | \mathbf{X})} = E(\delta_i | \mathbf{X}) \end{aligned}$$

- Cuando hay heterogeneidad en el efecto causal y el tratamiento se usa con un instrumento (válido), el estimador de mínimos cuadrados de dos fases o *TSLS* difiere, en general, del efecto medio del tratamiento

$$\begin{aligned} \text{plim}_{p \rightarrow \infty} \delta_i^{TSLS} &= \frac{\text{Cov}(\mathbf{x}_i' \boldsymbol{\beta} + \delta_i (\mathbf{z}_i' \boldsymbol{\pi}_i + v_i) + \varepsilon_i, \mathbf{z}_i)}{\text{Var}(\mathbf{z}_i' \boldsymbol{\pi}_i + v_i | \mathbf{X})} = \frac{E(\delta_i \boldsymbol{\pi}_i | \mathbf{X})}{E(\boldsymbol{\pi}_i | \mathbf{X})} = \\ &= E\left(\frac{\delta_i \boldsymbol{\pi}_i}{E(\boldsymbol{\pi}_i)} \mid \mathbf{X}\right) \end{aligned}$$

- En este caso hay heterogeneidad en el efecto del instrumento sobre el regresor, de modo que los coeficientes  $\boldsymbol{\pi}_i$  también son variables aleatorias e independientes del instrumento y de los términos de error
- El estimador es una media ponderada de los efectos causales individuales, en donde los pesos son  $\boldsymbol{\pi}_i / E(\boldsymbol{\pi}_i)$  y miden el grado relativo en el que el instrumento influencia la recepción del tratamiento para un individuo  $i$ . Este estimador se llama efecto medio del tratamiento local o ATEL, que enfatiza que se pesa más a individuos a los que influye más el instrumento
- Las implicaciones que tiene esto es que los estimadores varíen en función de que variables instrumentales se escojan (dado que se obtienen estimadores ATEL diferentes), y, por tanto, la prueba con el *J-statistic* puede rechazar la hipótesis nula, aunque los instrumentos sean válidos

## Los modelos para datos de panel

- Las bases de datos longitudinales o datos de panel son muy comunes en economía. Este tipo de datos son ricos en información y su análisis permite aprender sobre procesos económicos mientras se tiene en cuenta la heterogeneidad entre individuos y los efectos dinámicos que no aparecen en las secciones transversales. No obstante, su modelaje requiere especificaciones estocásticas complejas
  - Las bases de datos de panel o de datos longitudinales contienen secciones transversales grandes que se siguen en el tiempo, normalmente en pocos periodos
    - Los datos de panel están más orientados a análisis transversales en donde la heterogeneidad entre unidades es una parte integral del análisis, de modo que es posible analizar los cambios en el tiempo de las distribuciones transversales
    - El análisis de datos de panel es muy importante en econometría debido a que estas bases de datos proporcionan un ambiente rico para el desarrollo de técnicas de estimación y resultados teóricos
    - El uso de modelos enfocados a los datos de panel permite aislar los componentes de la varianza y poder controlar heterogeneidad no observada en diseños transversales (tales como la heterogeneidad que proviene del tiempo)
  - La ventaja fundamental de estos modelos de datos de panel sobre los modelos de datos transversales es que permiten mucha flexibilidad para modelar las diferencias de comportamiento entre individuos o grupos (la heterogeneidad individual). El marco básico para la discusión es el modelo de regresión siguiente:

$$y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + \mathbf{z}_i'\boldsymbol{\alpha} + \varepsilon_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + c_i + \varepsilon_{it}$$

- En esta regresión hay  $K$  regresores en  $\mathbf{x}_{it}$  sin incluir un término constante, mientras que la heterogeneidad o efectos individuales se introducen mediante  $c_i \equiv \mathbf{z}_i'\boldsymbol{\alpha}$ , donde  $\mathbf{z}_i$  contiene un término constante y un conjunto de variables específicas de los individuos o grupos que pueden ser observadas (raza, sexo, etc.) o no (características familiares, habilidad, etc.) que son constantes a lo largo del tiempo  $t$
- Si  $\mathbf{z}_i$  se observa para todos los individuos o grupos, entonces el modelo entero se puede tratar como un modelo de regresión ordinario y se puede ajustar con mínimos cuadrados ordinarios. Pero si  $c_i$  no se puede observar, entonces hay más complicaciones (lo cual suele ser el caso)

- Este término se introduce para tener en cuenta de manera explícita el sesgo que se produciría por la heterogeneidad no observada a la hora de estimar los parámetros (permite controlarla)
- El objetivo del análisis es poder estimar de manera consistente y eficiente los efectos parciales de cada variable sobre  $y$  mientras se controla (se tiene en cuenta) la heterogeneidad

$$\beta = \frac{\partial E(y_{it} | x_{it})}{\partial x_{it}}$$

- Que eso sea posible depende de las suposiciones sobre los efectos que no se observan. Una de las suposiciones principales es la exogeneidad estricta de  $x_{it}$ , de modo que el error no tiene correlación con las variables independientes en cualquier periodo pasado, presente y futuro)

$$E(\varepsilon_{it} | x_{i1}, x_{i2}, \dots) = 0$$

- La segunda es la independencia media de  $c_i$ , de modo que si las variables no observadas no tienen correlación con las incluidas  $x_{it}$ , entonces se pueden incluir en el error del modelo. Como es una suposición muy estricta, una más general sería que es una función de las variables independientes, pero es más complicada

$$E(c_i | x_{i1}, x_{i2}, \dots) = \alpha$$

$$\Rightarrow Cov(c_i, X_i | X_i) = Cov(E(c_i | X_i), X_i | X_i) = Cov(\alpha, X_i | X_i) = 0$$

- Hay una variedad de modelos diferentes para datos de panel, pero los principales se pueden clasificar en cuatro:
  - En el modelo de regresión agrupada,  $z_i$  solo contiene un término constante, por lo que los estimadores de mínimos cuadrados ordinarios proporcionan estimadores eficientes y consistentes del parámetro  $\alpha$  común y del vector  $\beta$
  - Si  $z_i$  no se observa, pero está correlacionado con  $x_{it}$ , por lo que el estimador de mínimos cuadrados ordinarios de  $\beta$  está sesgado y es inconsistente por la omisión de variables. No obstante, si  $\alpha_i \equiv c_i = z_i' \alpha$ , este incluye todos los efectos observables y especifica una media condicional estimable, y  $c_i$  será un término constante específico para cada grupo en el modelo, creando el modelo de efectos fijos (refiriéndose a la correlación entre  $c_i$  y  $x_{it}$ )

$$y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + \alpha_i + \varepsilon_{it}$$

- Si la heterogeneidad de los individuos o grupos no observados no está correlacionada con las variables incluidas, se puede construir un modelo de regresión lineal con un error compuesto que puede ser estimado de manera consistente (pero no eficiente), dando lugar al modelo de efectos aleatorios. En este caso,  $u_i$  es un elemento aleatorio específico del grupo (similar a  $\varepsilon_{it}$  pero invariante en el tiempo)

$$y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + \mathbf{z}_i'\boldsymbol{\alpha} + \varepsilon_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + E(\mathbf{z}_i'\boldsymbol{\alpha}) + [\mathbf{z}_i'\boldsymbol{\alpha} - E(\mathbf{z}_i'\boldsymbol{\alpha})] + \varepsilon_{it}$$

$$\Rightarrow y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + \alpha + u_i + \varepsilon_{it}$$

- El modelo de efectos aleatorios se puede ver como un modelo con coeficiente constante aleatorio, pero se puede extender esta idea a incluir otros coeficientes que varíen de manera aleatoria entre individuos o grupos. El modelo de parámetros aleatorios representa una extensión natural para aumentar la heterogeneidad, y  $\mathbf{h}_i$  es un vector aleatorio que induce a la variación en los parámetros entre individuos o grupos

$$y_{it} = \mathbf{x}_{it}'(\boldsymbol{\beta} + \mathbf{h}_i) + (\alpha + u_i) + \varepsilon_{it}$$

- Un aspecto importante del análisis de datos de panel es el balance del panel y su fijación

- Si una base de datos de panel tiene  $n$  conjuntos de observaciones y cada individuo se observa por  $T$  periodos de tiempo, entonces el panel está balanceado
- Si los individuos se observan por diferentes periodos de tiempo  $T_i$ , entonces el panel no está balanceado
- Si la base de datos de panel estudia el mismo conjunto de individuos para toda la duración del estudio, el panel es fijo. Existen otros paneles en los que un conjunto de individuos se rota (queda fuera del estudio) en un periodo concreto, por lo que el panel se denomina rotativo
- Los paneles también se pueden clasificar según la proporción de individuos con respecto a los periodos: si  $T$  es mucho mayor a  $n$ , el panel se denomina panel largo o *long panel*, mientras que en el caso contrario (con  $T$  mucho menor a  $n$ ) se denomina panel corto o *short panel*



- Las propiedades asintóticas de los estimadores en el modelo de regresión múltiple se establecieron bajo las suposiciones anteriormente vistas (linealidad, rango máximo, exogeneidad de variables independientes, homoscedasticidad y no autocorrelación y observaciones i.i.d.)

- No obstante, hay veces que al tratar con datos de panel se tienen hacer excepciones, dado que la muestra consiste de múltiples observaciones en cada una de las unidades observacionales, por lo que las  $x$  pueden estar correlacionadas entre observaciones en el tiempo (al menos dentro de cada unidad, aunque podría haber unidades correlacionadas también)
- Asumiendo que los datos constan de un número fijo de observaciones  $T$  en un conjunto de  $n$  individuos, el número total de filas en  $X$  será  $N = nT$ , y la matriz  $\bar{Q}_n \equiv \frac{1}{n} \sum_{i=1}^n Q_i$  se puede expresar de la siguiente manera:

$$\bar{Q}_n = \frac{1}{n} \sum_{i=1}^n Q_i = \frac{1}{n} \sum_i \frac{1}{T} \sum_{t \in i} Q_{it} = \frac{1}{n} \sum_{i=1}^n \bar{Q}_i$$

$$\bar{Q}_i = \text{average } Q_{it} \text{ for individual}$$

- Las condiciones que se necesitan para establecer convergencia aplican con respecto al número total de unidades observacionales que se consideren (en donde cada grupo puede ser una sola observación o un conjunto de observaciones)
- Un modelo más generalizado para el marco de las regresiones de datos de panel, es el siguiente:

$$y_{it} = x'_{it} \beta_{it} + u_{it} \quad \text{where } u_{it} = \alpha_i + \delta_t + \varepsilon_{it}$$

- En este modelo, el error  $u_{it}$  se puede dividir en perturbaciones que provienen de la dimensión individual, de la dimensión temporal, y de ambas a la vez
- Por lo tanto, las suposiciones de homoscedasticidad y de no autocorrelación no son adecuadas en este caso, aunque lo que normalmente se asume para este modelo son las siguientes:

$$E(u_{it} | x_{it}) = 0 \quad \text{Var}(u_{it}) = \sigma_{it}^2 \quad \text{Cov}(u_{it}, u_{js}) = \sigma_{ij,ts}$$

- Los parámetros pueden ser estimados por mínimos cuadrados ordinarios o por mínimos cuadrados generalizados dadas estas suposiciones, aunque son muy simplistas

- Muchas veces se asume que  $\delta_t = 0$ , aunque no siempre, dado que hay modelos que consideran los efectos temporales
- El análisis se comienza asumiendo la versión más simple del modelo general de datos de panel, el modelo agrupado o *pooled model*. En esta forma, si las suposiciones del modelo de regresión múltiple se cumplen, entonces no es necesario otros estimadores más que el de mínimos cuadrados ordinarios

$$y_{it} = \alpha + \mathbf{x}'_{it}\boldsymbol{\beta} + \varepsilon_{it} \text{ for } i = 1, \dots, n \text{ \& } t = 1, \dots, T$$

$$E(\varepsilon_{it} | \mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iT_i}) = 0$$

$$Var(\varepsilon_{it} | \mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iT_i}) = \sigma_\varepsilon^2$$

$$Cov(\varepsilon_{it}, \varepsilon_{is} | \mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iT_i}) = 0 \text{ if } i \neq j \text{ or } t \neq s$$

- Las suposiciones subyacentes de la estimación por mínimos cuadrados ordinarios que también se asumen en este modelo, no obstante, no suelen cumplirse como se asume en el modelo agrupado
  - En este caso, el modelo muestra que la heterogeneidad no observada es la misma para todos los individuos (se trata como el término constante de la regresión)
  - La heterogeneidad no observada induce a la autocorrelación de los errores en el tiempo, haciendo que el estimador asintótico de mínimos cuadrados de la varianza esté sesgado
  - Esto se puede ver de manera clara a través del modelo de efectos aleatorios que se verá más adelante, en donde el término de error para los grupos  $u_i$  está correlacionado con  $\mathbf{x}_{it}$  y por tanto introduce sesgo en los estimadores

$$\begin{aligned} y_{it} &= \alpha + \mathbf{x}'_{it}\boldsymbol{\beta} + \varepsilon_{it} + (c_i - E(c_i | \mathbf{X}_i)) = \alpha + \mathbf{x}'_{it}\boldsymbol{\beta} + \varepsilon_{it} + u_i = \\ &= \alpha + \mathbf{x}'_{it}\boldsymbol{\beta} + w_{it} \end{aligned}$$

$$\begin{aligned} \Rightarrow E(w_{it}w_{is}) &= Cov(w_{it}, w_{is}) = Cov(\varepsilon_{it} + u_i, \varepsilon_{is} + u_i) = \\ &= Cov(u_i, u_i) = \sigma_u^2 \text{ for } t \neq s \end{aligned}$$

- Suponiendo un modelo aún más general en las que se juntan todas las observaciones  $T_i$  para un individuo  $i$  en una sola ecuación y que el vector de perturbaciones  $\mathbf{w}_i$  incluye  $\varepsilon_{it}$  y los componentes omitidos, se dan las siguientes igualdades:

$$y_i = X_i\beta + w_i \quad \text{Var}(w_i|X_i) = \sigma_\varepsilon^2 I_{T_i} + \Sigma_i = \Omega_i$$

- El estimador de mínimos cuadrados ordinarios  $\beta$ , bajo autocorrelación, será el siguiente:

$$\begin{aligned} b &= (X'X)^{-1}X'y = \left[ \sum_{i=1}^n X_i'X_i \right]^{-1} \sum_{i=1}^n X_i'y_i = \\ &= \left[ \sum_{i=1}^n X_i'X_i \right]^{-1} \sum_{i=1}^n X_i'(X_i\beta + w_i) = \beta + \left[ \sum_{i=1}^n X_i'X_i \right]^{-1} \sum_{i=1}^n X_i'w_i \end{aligned}$$

- La consistencia se puede establecer de la misma manera que anteriormente, y la verdadera matriz de covarianzas asintótica tomaría la forma que se vio para el modelo de regresión generalizado

$$\text{Asy. Var}(b) =$$

$$\begin{aligned} &= \frac{1}{n} \text{plim} \left[ \frac{1}{n} \sum_{i=1}^n X_i'X_i \right]^{-1} \text{plim} \left[ \frac{1}{n} \sum_{i=1}^n X_i'w_iw_i'X_i \right] \text{plim} \left[ \frac{1}{n} \sum_{i=1}^n X_i'X_i \right]^{-1} = \\ &= \frac{1}{n} \text{plim} \left[ \frac{1}{n} \sum_{i=1}^n X_i'X_i \right]^{-1} \text{plim} \left[ \frac{1}{n} \sum_{i=1}^n X_i'\Omega_iX_i \right] \text{plim} \left[ \frac{1}{n} \sum_{i=1}^n X_i'X_i \right]^{-1} \end{aligned}$$

- La matriz central tiene que ser estimada, y se puede estimar de la misma manera que con el estimador de White (la lógica de este también aplica para este caso, aunque no sea exactamente lo mismo), en donde  $\hat{w}_i'$  es el vector de  $T_i$  residuos para el individuo  $i$

$$\text{Est. Asy. Var}(b) =$$

$$= \frac{1}{n} \text{plim} \left[ \frac{1}{n} \sum_{i=1}^n X_i'X_i \right]^{-1} \text{plim} \left[ \frac{1}{n} \sum_{i=1}^n X_i'\hat{w}_i\hat{w}_i'X_i \right] \text{plim} \left[ \frac{1}{n} \sum_{i=1}^n X_i'X_i \right]^{-1}$$

- No obstante, el problema más importante no es la heteroscedasticidad sino la autocorrelación entre observaciones, de modo que el estimador de White no es una solución al problema de inferencia
- Muchos estudios se realizan sacando muestras de clústeres, ya que de otro modo se pueden surgir naturalmente efectos similares a los efectos aleatorios en los datos de panel

- Los efectos de cada clúster dentro de la muestra pueden inducir a la correlación entre observaciones similares a los efectos aleatorios o fijos que se han identificado anteriormente, lo cual, si se ignoran, pueden llevar a errores de inferencia graves
- Para un modelo de dos niveles (en donde hay clústeres dentro de la muestra), un enfoque natural que se puede llevar a cabo es utilizar el enfoque robusto de efectos agrupados visto anteriormente
- El estimador de clústeres para un modelo de un solo nivel es muy parecido al del modelo de regresión de efectos agrupados, pero hay una diferencia en el proceso de generación de datos, dado que, en el primer ajuste, los individuos en el grupo generalmente se observan una vez, y su asociación (los efectos comunes) es menos definida
- Para poder tener en cuenta los efectos en muestras pequeñas cuando el número de clústeres es un porcentaje significativo de un total finito, se utiliza una versión refinada del estimador anteriormente visto con una corrección de los grados de libertad

$$\begin{aligned}
 \text{Est. Asy. Var}(\mathbf{b}) &= \\
 &= \text{plim} \left[ \sum_{g=1}^G \mathbf{X}_g' \mathbf{X}_g \right]^{-1} \text{plim} \left[ \frac{G}{G-1} \sum_{g=1}^G \left( \sum_{i=1}^{n_g} \mathbf{x}_{ig} \hat{\mathbf{w}}_{ig} \right) \left( \sum_{i=1}^{n_g} \mathbf{x}_{ig} \hat{\mathbf{w}}_{ig} \right)' \right] \text{plim} \left[ \sum_{g=1}^G \mathbf{X}_g' \mathbf{X}_g \right]^{-1} = \\
 &= \text{plim} \left[ \sum_{g=1}^G \mathbf{X}_g' \mathbf{X}_g \right]^{-1} \text{plim} \left[ \frac{G}{G-1} \sum_{g=1}^G \mathbf{X}_g' \hat{\mathbf{W}}_g \hat{\mathbf{W}}_g' \mathbf{X}_g \right] \text{plim} \left[ \sum_{g=1}^G \mathbf{X}_g' \mathbf{X}_g \right]^{-1}
 \end{aligned}$$

$$G = n^{\circ} \text{ clusters in sample} \quad n_g = n^{\circ} \text{ indiv. in } g$$

- Si el número de clústeres es pequeño, entonces la corrección de clústeres utilizada después de estimar un modelo de regresión agrupada por mínimos cuadrados ordinarios no es muy útil
- El modelo de regresión agrupado se puede estimar utilizando las medias muestrales de los datos. El modelo de regresión implícito se obtiene multiplicando cada grupo por  $(1/T)\mathbf{i}'$ , en donde  $\mathbf{i}$  es un vector columna de unos

$$(1/T)\mathbf{i}'\mathbf{y}_i = (1/T)\mathbf{i}'\mathbf{x}_i\boldsymbol{\beta} + (1/T)\mathbf{i}'\mathbf{w}_i \Rightarrow \bar{y}_i = \bar{\mathbf{x}}_i\boldsymbol{\beta} + \bar{w}_i$$

- En el modelo de regresión lineal transformado, las perturbaciones siguen teniendo una media condicional nula, pero varianzas heteroscedásticas. Con  $\Omega_i$  sin especificar, la regresión es heteroscedástica y se puede utilizar el estimador de White, pero si las suposiciones clásicas no se mantienen, entonces no se resuelve el problema
- Utilizar la media de grupos no resuelve el problema, pero la pérdida de información que ocurre al hacer la media es relativamente pequeña (aunque los datos desagregados sean mejores) y se permite eliminar la dimensión temporal en el modelo
- Considerando un modelo general en el que  $c_i$  es el efecto latente (variable aleatoria), si la suposición de independencia media  $E(c_i|X_i) = \alpha$  no se cumple, entonces el efecto latente también se transmite a las medias grupales

$$E(c_i|X_i) = h(X_i) = \bar{x}_i\gamma$$

$$\begin{aligned} y_{it} &= x'_{it}\beta + c_i + \varepsilon_{it} = x'_{it}\beta + \bar{x}_i\gamma + [\varepsilon_{it} + c_i - \bar{x}_i\gamma] = \\ &= x'_{it}\beta + \bar{x}_i\gamma + u_{it} \end{aligned}$$

$$\Rightarrow \bar{y}_i = \bar{x}_i\beta + \bar{x}_i\gamma + \bar{u}_i = \bar{x}_i(\beta + \gamma) + \bar{u}_i$$

- En ese caso, el estimador de medias grupales estima  $\beta + \gamma$  y no  $\beta$ , por lo que se agrupan los efectos observados y no observados en el mismo estimador. Si, en otro caso, hubiera error de medición en los regresores, este error se disiparía en la media y el estimador seguiría siendo consistente (aunque el de mínimos cuadrados no)
- El modelo de regresión agrupada se puede estimar de manera consistente (no eficiente) con el estimador de mínimos cuadrados ordinarios. Si se consideran las matrices de sumas de cuadrados y productos cruzados para cada formulación equivalente del modelo, se obtienen las siguientes identidades:
  - En el modelo original, los momentos acumularán la variación en las medias totales  $\bar{\bar{y}}$  y  $\bar{\bar{x}}$  y se usarían las siguientes sumas en el estimador:

$$y_{it} = \alpha + x'_{it}\beta + \varepsilon_{it}$$

$$\mathbf{S}_{xx}^{total} = \sum_{i=1}^n \sum_{t=1}^T (\mathbf{x}_{it} - \bar{\mathbf{x}})(\mathbf{x}_{it} - \bar{\mathbf{x}})'$$

$$\mathbf{S}_{xy}^{total} = \sum_{i=1}^n \sum_{t=1}^T (\mathbf{x}_{it} - \bar{\mathbf{x}})(y_{it} - \bar{y})$$

- En el modelo en términos de la media de grupos, la media de la media de los grupos es la media total, por lo que las matrices de momentos son sumas de cuadrados entre grupos (la variación de las medias de los grupos con respecto a las medias totales) y productos cruzados

$$\bar{y}_i = \alpha + \bar{\mathbf{x}}_i \boldsymbol{\beta} + \bar{\varepsilon}_i$$

$$\mathbf{S}_{xx}^{between} = T \sum_{i=1}^n (\bar{\mathbf{x}}_i - \bar{\mathbf{x}})(\bar{\mathbf{x}}_i - \bar{\mathbf{x}})'$$

$$\mathbf{S}_{xy}^{between} = \sum_{i=1}^n T(\bar{\mathbf{x}}_i - \bar{\mathbf{x}})(\bar{y}_i - \bar{y})$$

- En el modelo en términos de desviaciones de la media grupal, como los datos ya están en desviaciones, las medias de estas son nulas y las matrices de momentos son sumas de cuadrados dentro de los grupos (la variación de las observaciones con respecto a las medias grupales) y productos cruzados

$$y_{it} - \bar{y}_i = (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)' \boldsymbol{\beta} + \varepsilon_{it} - \bar{\varepsilon}_i$$

$$\mathbf{S}_{xx}^{within} = \sum_{i=1}^n \sum_{t=1}^T (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)(\mathbf{x}_{it} - \bar{\mathbf{x}}_i)'$$

$$\mathbf{S}_{xy}^{within} = \sum_{i=1}^n \sum_{t=1}^T (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)(y_{it} - \bar{y}_i)'$$

- Se puede verificar que la suma de cuadrados total y el producto cruzado total es la suma de aquellos entre grupos y dentro de los grupos

$$\mathbf{S}_{xx}^{total} = \mathbf{S}_{xx}^{within} + \mathbf{S}_{xx}^{between}$$

$$\mathbf{S}_{xy}^{total} = \mathbf{S}_{xy}^{within} + \mathbf{S}_{xy}^{between}$$

- Por lo tanto, hay tres posibles estimadores de mínimos cuadrados para el vector  $\beta$  del modelo agrupado, que corresponden a la descomposición de  $S_{xx}^{total}$  y  $S_{xy}^{total}$

- Para el modelo original, el estimador de mínimos cuadrados es el original. Este es consistente y no está sesgado, pero no es el más eficiente

$$b^{total} = [S_{xx}^{total}]^{-1} S_{xy}^{total} = [S_{xx}^{between} + S_{xx}^{within}]^{-1} [S_{xy}^{between} + S_{xy}^{within}]$$

- Para el modelo en términos de medias grupales, el estimador de mínimos cuadrados es el estimador entre grupos, el cual es el estimador de medias grupales visto anteriormente. Este es consistente y no está sesgado cuando  $nT \rightarrow \infty$ , pero, aunque tampoco sea eficiente, tiene menor varianza que el original

$$b^{between} = [S_{xx}^{between}]^{-1} [S_{xy}^{between}]$$

- Para el modelo en términos de medias grupales, el estimador de mínimos cuadrados es el estimador dentro de grupos. Este es consistente y no está sesgado cuando  $nT \rightarrow \infty$ , y además tiene menos varianza que los dos estimadores anteriores (aunque no es el más eficiente)

$$b^{within} = [S_{xx}^{within}]^{-1} [S_{xy}^{within}]$$

- A partir de estos estimadores, es posible encontrar una expresión equivalente para los productos cruzados y así, encontrar una expresión equivalente de  $b^{total}$  en términos de matrices

$$S_{xy}^{between} = S_{xx}^{between} b^{between}$$

$$S_{xy}^{within} = S_{xx}^{within} b^{within}$$

$$\Rightarrow b^{total} = [S_{xx}^{between} + S_{xx}^{within}]^{-1} [S_{xx}^{between} b^{between} + S_{xx}^{within} b^{within}]$$

$$\Rightarrow b^{total} = F^{between} b^{between} + F^{within} b^{within}$$

$$\text{where } F^{between} \equiv [S_{xx}^{between} + S_{xx}^{within}]^{-1} S_{xx}^{between} = I - F^{within}$$

- Cuando  $T = 2$ , se puede demostrar que el estimador de mínimos cuadrados ordinarios, el estimador entre grupos y el de primeras diferencias coinciden entre sí

$$\mathbf{b}^{within} = \frac{\sum_{i=1}^n \left( \frac{\mathbf{x}_{i2} + \mathbf{x}_{i1}}{2} \right) \left( \frac{y_{i2} + y_{i1}}{2} \right)}{\sum_{i=1}^n \left( \frac{\mathbf{x}_{i2} + \mathbf{x}_{i1}}{2} \right) \left( \frac{\mathbf{x}_{i2} + \mathbf{x}_{i1}}{2} \right)'} = \frac{\sum_{i=1}^n (y_{i2} + y_{i1})}{\sum_{i=1}^n (\mathbf{x}_{i2} + \mathbf{x}_{i1})}$$

$$\mathbf{b} = \frac{\sum_{i=1}^n \left( \frac{\mathbf{x}_{i2} + \mathbf{x}_{i1}}{2} \right) \left( \frac{y_{i2} + y_{i1}}{2} \right)}{\sum_{i=1}^n \left( \frac{\mathbf{x}_{i2} + \mathbf{x}_{i1}}{2} \right) \left( \frac{\mathbf{x}_{i2} + \mathbf{x}_{i1}}{2} \right)'} = \frac{\sum_{i=1}^n (y_{i2} + y_{i1})}{\sum_{i=1}^n (\mathbf{x}_{i2} + \mathbf{x}_{i1})}$$

$$\mathbf{b}^{FD} = \frac{\sum_{i=1}^n \left( \frac{\mathbf{x}_{i2} + \mathbf{x}_{i1}}{2} \right) \left( \frac{y_{i2} + y_{i1}}{2} \right)}{\sum_{i=1}^n \left( \frac{\mathbf{x}_{i2} + \mathbf{x}_{i1}}{2} \right) \left( \frac{\mathbf{x}_{i2} + \mathbf{x}_{i1}}{2} \right)'} = \frac{\sum_{i=1}^n (y_{i2} + y_{i1})}{\sum_{i=1}^n (\mathbf{x}_{i2} + \mathbf{x}_{i1})}$$

$$T = 2 \Rightarrow \mathbf{b} = \mathbf{b}^{within} = \mathbf{b}^{FD}$$

- El modelo de efectos fijos nace de la suposición sobre los efectos no observados en el modelo general  $c_i$  están correlacionados con las variables incluidas  $\mathbf{x}_{it}$ , de modo que añadiendo este término  $c_i$  a la regresión se puede compensar por el sesgo introducido por esta correlación

$$y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + c_i + \varepsilon_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + h(\mathbf{X}_i) + \varepsilon_{it} + [c_i - h(\mathbf{X}_i)]$$

$$\Rightarrow y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + \alpha_i + \varepsilon_{it} + [c_i - h(\mathbf{X}_i)]$$

- Se supone que  $c_i$  es una variable aleatoria con  $E(c_i|\mathbf{X}_i) = h(\mathbf{X}_i)$ , de modo que  $c_i - h(\mathbf{X}_i)$  no tiene correlación con  $\mathbf{X}_i$  (por los argumentos vistos anteriormente) y se puede absorber en el término de error  $\varepsilon_{it}$

$$\Rightarrow E(y_{it}|\mathbf{x}_{it}) = E(\mathbf{x}_{it}'\boldsymbol{\beta} + h(\mathbf{X}_i) + \varepsilon_{it} + [c_i - h(\mathbf{X}_i)]|\mathbf{x}_{it}) = \mathbf{x}_{it}'\boldsymbol{\beta} + \alpha_i$$

$$\Rightarrow Cov(c_i - h(\mathbf{X}_i), \mathbf{X}_i|\mathbf{X}_i) = Cov(E(c_i|\mathbf{X}_i) - h(\mathbf{X}_i), \mathbf{X}_i|\mathbf{X}_i) =$$

$$= Cov(h(\mathbf{X}_i) - h(\mathbf{X}_i), \mathbf{X}_i|\mathbf{X}_i) = Cov(0, \mathbf{X}_i|\mathbf{X}_i) = 0$$

$$\Rightarrow y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + \alpha_i + \varepsilon_{it}$$

- En verdad, los efectos fijos sí que podrían estar correlacionados con los regresores, pero se requiere que no para la consistencia de los estimadores en este modelo. De estar correlacionados, los estimadores no serán consistentes
- Otra suposición que se hace que es que  $Var(c_i|\mathbf{X}_i)$  es constante, por lo que con ambas suposiciones se obtiene un modelo de regresión lineal múltiple



$$Var(c_i|X_i) = \sigma_c^2$$

- Qué  $E(c_i|X_i) = h(X_i)$  significa que los efectos latentes son fijos en el tiempo (no dependen del tiempo), y la formulación implica que las diferencias entre grupos se pueden capturar en las diferencias en el término constante  $\alpha_i$  (comparando la de un grupo con otro). Cada  $\alpha_i$  es un parámetro desconocido que se tiene que estimar
- Una de las mayores desventajas de la regresión de efectos fijos es que cualquier variable  $x_{it}$  que no varíe en el tiempo será absorbida en  $\alpha_i$  y los coeficientes de estas variables no se podrían estimar con los dos estimadores vistos anteriormente

$$\bar{x}_i = \frac{\sum_{t=1}^T x_{it}}{T} = x_i \quad \text{for } i = 1, 2, \dots, N$$

- Las suposiciones hechas tienen sentido siempre que la relación  $N/T$  no sea muy grande, ya que de otro modo el modelo tendría pocos grados de libertad

$$\text{Degrees of Freedom} = N(T - 1) - K$$

- Siendo  $y_i$  y  $X_i$  las  $T$  observaciones para la unidad observacional  $i$ ,  $i$  un vector columna  $T \times 1$  y  $\varepsilon_i$  un vector columna  $T \times 1$  de perturbaciones, el modelo se puede expresar de la siguiente manera:

$$y_i = X_i\beta + i\alpha_i + \varepsilon_i$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \beta + \begin{bmatrix} i & 0 & \dots & 0 \\ 0 & i & \dots & 0 \\ & & \ddots & \\ 0 & 0 & \dots & i \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

- Este se puede expresar de manera equivalente creando una matriz  $D \equiv [d_1|d_2|\dots|d_n]$  de tamaño  $nT \times n$ , donde  $d_i$  es una variable binaria que indica la unidad  $i$ . Juntando las  $nT$  filas, se obtiene la siguiente expresión, denominado modelo de mínimos cuadrados de variable binaria (LSDV):

$$y = X\beta + D\alpha + \varepsilon$$

- Como el modelo es uno de regresión múltiple común, se puede utilizar el estimador de mínimos cuadrados ordinarios, que será equivalente al estimador dentro de grupos (en donde se les resta la media temporal a los regresores):

$$\mathbf{b}^{LSDV} = [\mathbf{X}'\mathbf{M}_D\mathbf{X}]^{-1}[\mathbf{X}'\mathbf{M}_D\mathbf{y}] = \mathbf{b}^{within}$$

$$where \mathbf{M}_D = \mathbf{I} - \mathbf{D}(\mathbf{D}'\mathbf{D})^{-1}\mathbf{D}'$$

- Igual que antes, si  $T = 2$ , entonces también se pueden estimar los coeficientes por primeras diferencias y los estimadores de mínimos cuadrados y de dentro de grupos son equivalentes

$$\mathbf{b}^{LSDV} = \mathbf{b}^{within} = \mathbf{b}^{FD}$$

- El estimador entre grupos no se puede utilizar aquí debido a que los efectos fijos se verían absorbidos a utilizar el modelo de medias temporales de las variables

$$\bar{\alpha}_i = \frac{\sum_{t=1}^T \alpha_i}{T} = \alpha_i \quad for \quad i = 1, 2, \dots, N$$

- Los coeficientes de las variables binarias  $\mathbf{d}_i$  se pueden recuperar a partir de la ecuación normal de la regresión parcial, de modo que se puede obtener el coeficiente  $\alpha_i^{LSDV}$  para cada  $i$ , siendo un estimador para  $\alpha_i^{LSDV}$  consistente cuando  $T \rightarrow \infty$

$$\mathbf{D}'\mathbf{D}\mathbf{a} + \mathbf{D}'\mathbf{X}\mathbf{b}^{LSDV} = \mathbf{D}'\mathbf{y} \Rightarrow \mathbf{a} = [\mathbf{D}'\mathbf{D}]^{-1}\mathbf{D}'(\mathbf{y} - \mathbf{X}\mathbf{b}^{LSDV})$$

$$\Rightarrow a_i = \bar{y}_i - \bar{\mathbf{x}}_i'\mathbf{b}^{LSDV}$$

- El estimador asintótico de la matriz de covarianza para  $\mathbf{b}^{LSDV}$  utiliza el segundo momento de la matriz con  $\mathbf{x}$  que se expresan como desviaciones de sus respectivas medias grupales

$$Est. Asy. Var(\mathbf{b}^{LSDV}) = s^2[\mathbf{X}'\mathbf{M}_D\mathbf{X}]^{-1} = s^2[\mathbf{S}_{xx}^{within}]^{-1}$$

$$s^2 = \frac{\sum_{i=1}^n \sum_{t=1}^T (y_{it} - \mathbf{x}_{it}'\mathbf{b}^{LSDV} - a_i)^2}{nT - n - K} =$$

$$= \frac{(\mathbf{M}_D\mathbf{y} - \mathbf{M}_D\mathbf{X}\mathbf{b}^{LSDV})'(\mathbf{M}_D\mathbf{y} - \mathbf{M}_D\mathbf{X}\mathbf{b}^{LSDV})}{nT - n - K}$$

- El numerador en  $s^2$  es exactamente la suma de residuos cuadrados utilizando los coeficientes de mínimos cuadrados ordinarios y los datos en forma de desviaciones de la media grupal

$$e_{it} = y_{it} - \mathbf{x}_{it}'\mathbf{b}^{LSDV} - a_i = y_{it} - \mathbf{x}_{it}'\mathbf{b}^{LSDV} - (\bar{y}_i - \bar{\mathbf{x}}_i'\mathbf{b}^{LSDV})$$

$$= (y_{it} - \bar{y}_i) + (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)'\mathbf{b}^{LSDV}$$

- Para los efectos individuales, el estimador de la varianza asintótica de  $a_i$  es el siguiente:

$$Asy.Var(a_i) = \frac{\sigma_\varepsilon^2}{T} + \bar{x}_i' [Asy.Var(\mathbf{b}^{LSDV})] \bar{x}_i$$

- Es posible utilizar un contraste  $t$  para contrastar  $a_i = 0$ , pero esto es solo para un grupo específico  $i$  y no sirve en este contexto de regresión. Por lo tanto, si uno está interesado en las diferencias entre grupos, se puede contrastar la hipótesis de que todos los términos constantes son iguales con un estadístico  $F$

$$F(n-1, nT-n-K) = \frac{(R_{LSDV}^2 - R_{Pooled}^2)/(n-1)}{(1 - R_{LSDV}^2)/(nT-n-K)}$$

- En este caso se compara el ajuste del modelo estimado con variables binarias y del modelo agrupado, dado que el modelo agrupado solo tiene un término constante (la heterogeneidad no observada) que es igual para todos los grupos
- Si se rechaza la hipótesis nula, entonces el modelo de efectos fijos es mejor que el modelo agrupado para modelar los datos (tiene mejor ajuste) porque la heterogeneidad no observada no es la misma entre individuos. De no rechazarse, la conclusión es la contraria
- El enfoque de regresión de mínimos cuadrados de variable binaria se puede extender para incluir efectos específicos del tiempo, de modo que el modelo se puede extender del siguiente modo:

$$y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + \alpha_i + \delta_t + \varepsilon_{it}$$

- Este modelo se obtiene del anterior incluyendo  $T-1$  variables binarias adicionales (a través del término  $\delta_t$ ) como al construir la matriz  $\mathbf{D}$  anterior (una no se incluye para evitar la trampa)
- No obstante, hay una asimetría en la formulación, dado que cada efecto de grupo tiene un intercepto específico mientras que los efectos temporales son contrastes (comparaciones hechas respecto al periodo base excluido), por lo que se puede formular el siguiente modelo simétrico en donde se incluyen los  $n$  y  $T$  efectos (coeficientes de las variables binarias) de manera restringida:

$$y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + \mu + \alpha_i + \delta_t + \varepsilon_{it} \quad s.t. \quad \sum_i \alpha_i = \sum_t \delta_t = 0$$

- Los estimadores por mínimos cuadrados ordinarios en este modelo se obtienen con una regresión de  $y_{it}^*$  sobre  $x_{it}^*$ , y los coeficientes de las variables binarias y del coeficiente constante se pueden recuperar de las ecuaciones normales

$$y_{it}^* = y_{it} - \bar{y}_i - \bar{y}_t + \bar{\bar{y}} \quad x_{it}^* = x_{it} - \bar{x}_i - \bar{x}_t + \bar{\bar{x}}$$

$$\hat{\mu} = m = \bar{\bar{y}} - \bar{\bar{x}}' \mathbf{b} \quad \hat{a}_i = a_i = (\bar{y}_i - \bar{\bar{y}}) - (\bar{x}_i - \bar{\bar{x}})' \mathbf{b}$$

$$\hat{\delta}_t = d_t = (\bar{y}_t - \bar{\bar{y}}) - (\bar{x}_t - \bar{\bar{x}})' \mathbf{b}$$

$$\bar{y}_t = \frac{1}{n} \sum_{i=1}^n y_{it} \quad \bar{\bar{y}} = \frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T y_{it}$$

$$\bar{x}_t = \frac{1}{n} \sum_{i=1}^n x_{it} \quad \bar{\bar{x}} = \frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T x_{it}$$

- La matriz de covarianzas asintótica para  $\mathbf{b}$  se puede estimar utilizando las sumas cuadradas y productos cruzados de  $x_{it}^*$ , y  $s^2$  se puede calcular con la siguiente fórmula:

$$s^2 = \frac{\sum_{i=1}^n \sum_{t=1}^T (y_{it} - x_{it}' \mathbf{b} - m - a_i + d_t)^2}{nT - (n - 1) - (T - 1) - K - 1}$$

- Este modelo no se suele utilizar en la práctica porque el coste en términos de grados de libertad no suele ser justificado ( $s^2$  es muy pequeño comparado con la formulación asimétrica) y se puede utilizar un modelo más general para modelar la evolución temporal del error
- Además, al incluir los efectos temporales, se puede ver que la estimación por primeras diferencias no permitirá eliminar estas para una estimación de mínimos cuadrados ordinarios
- Si los regresores no están correlacionados en el sentido estricto con los efectos individuales, entonces los efectos de la heterogeneidad provienen del término de error y lo más apropiado es modelar los términos constantes específicos como aleatorios para cada unidad, de modo que se puede utilizar el modelo de efectos aleatorios
  - El modelo de efectos fijos permite que los efectos de los individuos no observados estén correlacionados con las variables incluidas, modelando las diferencias entre unidades como desplazamientos del la

función de regresión (por las variables binarias) con el término  $\alpha_i$  que proviene de  $c_i$  (la heterogeneidad no observada)

- Este modelo puede entenderse como uno que solo aplica para las unidades transversales en el estudio, pero no para unidades observacionales fuera de él (los efectos fijos son para los individuos de la muestra, no para individuos fuera de ella)
  - El enfoque de parámetros aleatorios es más útil si se considera que la muestra transversal se ha sacado de una población grande (dado que el efecto fijo específico de los individuos sería como aleatorio debido al tamaño de la población), dado que los efectos específicos de cada unidad se distribuirán de manera equivalente para todas las unidades (haciendo que los efectos individuales de la heterogeneidad no observable sean globales y no se puedan diferenciar entre individuos como antes)
  - Haciendo esto, es posible reducir el número de parámetros a estimar, aunque el coste es la posibilidad de tener estimador inconsistentes si las suposiciones son inapropiadas
- Si se reconsidera una reformulación del modelo visto anteriormente, pero se incluye la constante en los  $K$  regresores, un término constante  $\alpha$  y un término aleatorio  $u_i$ , entonces la formulación del modelo es la siguiente:

$$y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + E(\mathbf{z}_i'\boldsymbol{\alpha}|\mathbf{x}_{it}) + \mathbf{z}_i'\boldsymbol{\alpha} - E(\mathbf{z}_i'\boldsymbol{\alpha}|\mathbf{x}_{it}) + \varepsilon_{it}$$

$$\Rightarrow y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + E(c_i|\mathbf{x}_{it}) + c_i - E(c_i|\mathbf{x}_{it}) + \varepsilon_{it}$$

$$\Rightarrow y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + \alpha + (u_i + \varepsilon_{it})$$

- En este modelo, se asume que los efectos individuales específicos para cada unidad se distribuyen alrededor de una media común  $\alpha \equiv E(\mathbf{z}_i'\boldsymbol{\alpha}|\mathbf{x}_{it})$  acorde a una distribución de probabilidad
- El término constante  $\alpha$  es la media de la heterogeneidad no observada (constante en el tiempo) y un término aleatorio  $u_i \equiv \mathbf{z}_i'\boldsymbol{\alpha} - E(\mathbf{z}_i'\boldsymbol{\alpha}|\mathbf{x}_{it})$  es la heterogeneidad aleatoria específica o la variación del efecto individual específico para la observación  $i$  con respecto a su media (constante en el tiempo), de modo que se puede interpretar como la incertidumbre introducida por la heterogeneidad no observable y aleatoria de los individuos
- Aunque se incluye la constante en los  $K$  regresores, la constante final del modelo será  $\beta_0 + \alpha$ , dado que no se pueden diferenciar

ambas a la hora de estimar (tal como pasa al estimar el modelo de efectos fijos con variables binarias, ya que el intercepto es  $\beta_0 + \alpha_{-i}$  para aquella binaria que se omite)

- Si se asume exogeneidad estricta, entonces se obtienen las siguientes identidades:

$$E(\varepsilon_{it}|\mathbf{X}) = 0 \quad E(u_{it}|\mathbf{X}) = 0$$

$$E(\varepsilon_{it}^2|\mathbf{X}) = \sigma_\varepsilon^2 \quad E(u_i^2|\mathbf{X}) = \sigma_u^2$$

$$E(\varepsilon_{it}u_j|\mathbf{X}) = 0 \quad \text{for all } i, t \text{ and } j$$

$$E(\varepsilon_{it}\varepsilon_{js}|\mathbf{X}) = 0 \quad \text{if } t \neq s \text{ or } i \neq j$$

$$E(u_u u_j|\mathbf{X}) = 0 \quad \text{if } i \neq j$$

- Igual que antes, es útil reformular el modelo en bloques de  $T$  a través de vectores para el grupo  $i$

$$\mathbf{y}_i = \mathbf{X}_i' \boldsymbol{\beta} + \alpha + u_i \mathbf{i} + \boldsymbol{\varepsilon}_i$$

- Para estas  $T$  observaciones se define un término de error compuesto  $\eta_{it} \equiv \varepsilon_{it} + u_i$ , cuyo vector para las  $T$  observaciones es  $\boldsymbol{\eta}_i' = (\eta_{i1}, \eta_{i2}, \dots, \eta_{iT})$  y hace que haya dos fuentes de perturbaciones

$$\mathbf{y}_i = \mathbf{X}_i' \boldsymbol{\beta} + \alpha + \boldsymbol{\eta}_i$$

- A partir de esta forma de  $\boldsymbol{\eta}_i$ , el modelo se puede reformular, y se llama modelo de componentes del error. En este modelo, se cumplen las siguientes igualdades:

$$E(\eta_{it}\eta_{is}|\mathbf{X}) = \sigma_\varepsilon^2 + \sigma_u^2$$

$$E(\eta_{it}\eta_{is}|\mathbf{X}) = \sigma_u^2 \quad \text{for } t \neq s$$

$$E(\eta_{it}\eta_{js}|\mathbf{X}) = 0 \quad \text{for all } t \text{ and } s \text{ if } i \neq j$$

- Para estas  $T$  observaciones para la unidad observacional  $i$ , siendo  $\boldsymbol{\Sigma} \equiv E(\boldsymbol{\eta}_i \boldsymbol{\eta}_i'|\mathbf{X})$  la matriz de covarianzas,  $\mathbf{i}_T$  un vector  $T \times 1$  de unos e  $\mathbf{I}_T$  una matriz identidad  $T \times T$ , esta se puede expresar de la siguiente manera:

$$\Sigma = \begin{bmatrix} \sigma_\varepsilon^2 + \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & \cdots & \sigma_u^2 \\ \sigma_u^2 & \sigma_\varepsilon^2 + \sigma_u^2 & \sigma_u^2 & \cdots & \sigma_u^2 \\ & \cdots & \cdots & \cdots & \cdots \\ \sigma_u^2 & \sigma_u^2 & \sigma_u^2 & \cdots & \sigma_\varepsilon^2 + \sigma_u^2 \end{bmatrix} = \sigma_\varepsilon^2 \mathbf{I}_T + \sigma_u^2 \mathbf{i}_T \mathbf{i}_T'$$

- Debido a que las observaciones  $i$  y  $j$  son independientes, la matriz de covarianzas de las perturbaciones para las  $nT$  observaciones es la siguiente:

$$\Omega = \begin{bmatrix} \Sigma & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \Sigma & \mathbf{0} & \cdots & \mathbf{0} \\ & & & \ddots & \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \Sigma \end{bmatrix} = \mathbf{I}_n \otimes \Sigma$$

- El modelo anterior es un modelo de regresión generalizado como los vistos anteriormente, por lo que las implicaciones anteriormente vistas del modelo de generalizado también aplican

- Las perturbaciones están autocorrelacionadas (en el tiempo) dentro de un mismo grupo, pero no entre grupos. Esta autocorrelación de  $y_{it}$  entre periodos equivale a la fracción de la varianza de  $y_{it}$  explicada por los efectos aleatorios

$$\text{Corr}(y_{it}, y_{is}) = \frac{E(\eta_{it}\eta_{is}|\mathbf{X})}{\sqrt{E(\eta_{it}\eta_{is}|\mathbf{X})}\sqrt{E(\eta_{it}\eta_{is}|\mathbf{X})}} = \frac{\sigma_u^2}{\sigma_\varepsilon^2 + \sigma_u^2}$$

$$\text{Fraction of variance of rand. effects} = \frac{\sigma_u^2}{\sigma_\varepsilon^2 + \sigma_u^2}$$

- En particular, los parámetros del modelo de efectos aleatorios se pueden estimar consistentemente mediante estimadores de mínimos cuadrados ordinarios, aunque no de manera eficiente. Además, la matriz de covarianzas se puede estimar asintóticamente y de manera robusta a través del estimador visto anteriormente para la del modelo agrupado
- No obstante, existen otros estimadores consistentes disponibles, tales como el estimador LSDV o el de medias grupales
- Asumiendo que no hay regresores invariantes en el tiempo, se pueden considerar las desviaciones de la media con respecto al tiempo y el estimador LSDV sería un estimador consistente de  $\beta$ . No obstante, este estimador es ineficiente

$$y_{it} - \bar{y}_i = (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)' \beta + \alpha + \varepsilon_{it} - \bar{\varepsilon}_i$$

- Considerando la media grupal, se puede estimar consistentemente  $\beta$  a través del estimador de medias grupales, aunque este estima de manera ineficiente

$$\bar{y}_i = \bar{x}_i' \beta + \alpha + u_i + \bar{\varepsilon}_i$$

- Aunque no sean los estimadores más eficientes, como la mayoría de modelos generalizados se estiman en dos partes, estos estimadores pueden servir como un estimador robusto de mínimos cuadrados para los parámetros de la varianza en el modelo. Los siguientes estimadores permiten obtener estimadores consistentes de funciones de varianzas:

$$e'_{pooled} e_{pooled} / nT \xrightarrow{P} \sigma_u^2 + \sigma_\varepsilon^2$$

$$e'_{LSDV} e_{LSDV} / (n(T-1) - K) \xrightarrow{P} \sigma_\varepsilon^2$$

$$e'_{FD} e_{FD} / n(T-1) \xrightarrow{P} 2\sigma_\varepsilon^2$$

$$e'_{between} e_{between} / nT \xrightarrow{P} \sigma_u^2 + \sigma_\varepsilon^2 / T$$

- Diferentes pares de estos estimadores (y otros candidatos) pueden proporcionar un estimador por método de momentos de dos ecuaciones para  $(\sigma_\varepsilon^2, \sigma_u^2)$ , por lo que se puede desarrollar un estimador de mínimos cuadrados generalizados eficiente
- Para obtener el estimador de mínimos cuadrados generalizados de la pendiente en este contexto, es necesario transformar los datos y aplicar mínimos cuadrados ordinarios

$$\hat{\beta} = (X' \Omega^{-1} X)^{-1} X' \Omega^{-1} y = \left( \sum_{i=1}^n X_i' \Sigma^{-1} X_i \right)^{-1} \left( \sum_{i=1}^n X_i' \Sigma^{-1} y_i \right)$$

- Se requiere que  $\Omega^{-1/2} = (I_n \otimes \Sigma)^{-1/2} = I_n \otimes \Sigma^{-1/2}$  (la matriz formada por el producto de cada uno de los componentes de  $I_n$  por la matriz entera  $\Sigma^{-1/2}$ ), de modo que solo es necesario encontrar  $\Sigma^{-1/2}$  y esta se define de la siguiente manera:

$$\Sigma^{-1/2} = \left[ I_T - \frac{\theta_i}{T} i_T i_T' \right] \quad \text{where} \quad \theta_i = 1 - \frac{\sigma_\varepsilon}{\sqrt{\sigma_\varepsilon^2 + T\sigma_u^2}}$$

- La transformación de  $y_i$  y  $X_i$  para el método de mínimos cuadrados generalizados es, por tanto, la siguiente:



$$\Sigma^{-1/2} \mathbf{y}_i = \frac{1}{\sigma_\varepsilon} \begin{bmatrix} y_{i1} - \theta \bar{y}_i \\ y_{i2} - \theta \bar{y}_i \\ \dots \\ y_{iT} - \theta \bar{y}_i \end{bmatrix} \quad \Sigma^{-1/2} \mathbf{X}_i = \frac{1}{\sigma_\varepsilon} \begin{bmatrix} x_{i1} - \theta \bar{x}_i \\ x_{i2} - \theta \bar{x}_i \\ \dots \\ x_{iT} - \theta \bar{x}_i \end{bmatrix}$$

- En el caso en el que  $\theta = 1$ , el procedimiento es equivalente al de estimar a través del estimador LSDV (el estimado para un modelo de efectos fijos), mientras que si  $\theta = 0$ , se obtiene un estimador de mínimos cuadrados ordinarios (como en el modelo agrupado). Esto ocurre porque  $\theta$  se puede entender como la raíz de la fracción de la desviación estándar explicada por el término  $u_i$  (para expresarse en las mismas unidades de  $y_{it}$ )

$$\theta = 1 - \frac{\sigma_\varepsilon}{\sqrt{\sigma_\varepsilon^2 + T\sigma_u^2}} = 1 - \sqrt{\frac{\sigma_\varepsilon^2}{\sigma_\varepsilon^2 + T\sigma_u^2}} = \sqrt{\frac{\sigma_u^2}{\sigma_\varepsilon^2 + T\sigma_u^2}}$$

- Se puede demostrar que el estimador de mínimos cuadrados generalizados es una media ponderada matricial de los estimadores dentro de grupos y entre grupos:

$$\hat{\beta} = \hat{F}^{within} \mathbf{b}^{within} + (\mathbf{I} - \hat{F}^{within}) \mathbf{b}^{between}$$

$$where \quad \hat{F}^{within} = [\mathbf{S}_{xx}^{within} + \lambda \mathbf{S}_{xx}^{between}]^{-1} \mathbf{S}_{xx}^{within} \quad \&$$

$$\lambda = \frac{\sigma_\varepsilon^2}{\sigma_\varepsilon^2 + T\sigma_u^2} = (1 - \theta)^2$$

- Los modelos no balanceados complican el modelo de efectos aleatorios porque la matriz  $\Omega \neq \mathbf{I}_n \otimes \Sigma$  porque los bloques diagonales de  $\Omega$  son de diferentes tamaños. No obstante, la estimación debe verse inafectada dado que la fuente de heteroscedasticidad proviene de las diferencias de tamaños entre grupos

$$\Sigma^{-1/2} = \left[ \mathbf{I}_{T_i} - \frac{\theta_i}{T_i} \mathbf{i}_{T_i} \mathbf{i}_{T_i}' \right] \quad where \quad \theta_i = 1 - \frac{\sigma_\varepsilon}{\sqrt{\sigma_\varepsilon^2 + T_i \sigma_u^2}}$$

- Mientras  $\lambda$  difiera de 1, la ineficiencia de los mínimos cuadrados ordinarios y del estimador dentro de los grupos usado para el modelo de efectos fijos provendrá de una ponderación ineficiente de los dos estimadores
- Comparado con el estimador de mínimos cuadrados generalizados, el estimador MCO pone demasiado peso en las unidades de variación entre grupos (incluye toda la variación en

$X$  y no tiene en cuenta que parte de la variación entre grupos proviene de la variación en  $u_i$  entre unidades)

- En cuanto al modelo de efectos fijos, se puede ver como, de atribuir la variabilidad a más de una fuente, el estimador dentro de los grupos no es eficiente porque no tiene en cuenta  $\sigma_\varepsilon^2$  y atribuye todo el peso a  $\mathbf{b}^{within}$
  - Por lo tanto, cuanto mayor sea  $\sigma_u^2$ , más cerca se estará del estimador usado en el modelo de efectos fijos, y mientras menor sea  $\sigma_u^2$ , más cerca se estará del estimador de mínimos cuadrados ordinarios
- El enfoque de variables binarias es costoso en términos de grados de libertad perdidos, mientras que la suposición de no correlación entre efectos individuales y regresores no está muy justificada. Por lo tanto, para poder saber si usar el modelo de efectos fijos o de efectos aleatorios, se puede utilizar el contraste de especificación de Hausman
- Este contrasta la ortogonalidad de los efectos comunes y de los regresores, basándose en la idea de que, bajo la hipótesis nula de no correlación, tanto el estimador LSDV como el FGLS deben ser consistentes, pero el LSDV debe ser ineficiente (se prefiere un modelo de efectos aleatorios). La hipótesis alternativa es que el estimador LSDV es consistente pero el estimador FGLS no lo es (se prefiere el modelo de efectos fijos)
  - Por lo tanto, bajo la hipótesis nula, las dos estimaciones no deberían diferir sistemáticamente, y se puede hacer un contraste para esa diferencia
  - Otro ingrediente esencial para el contraste es la matriz de varianzas y covarianzas del vector  $(\mathbf{b}_{FE} - \hat{\boldsymbol{\beta}}_{RE})$ . La covarianza entre un estimador eficiente y su diferencia de un estimador ineficiente debe ser nula, lo cual implica que la covarianza entre el estimador eficiente e ineficiente es la varianza del ineficiente
- $$Var(\mathbf{b}_{FE} - \hat{\boldsymbol{\beta}}_{RE}) = Var(\mathbf{b}_{FE}) + Var(\hat{\boldsymbol{\beta}}_{RE}) - 2Cov(\mathbf{b}_{FE}, \hat{\boldsymbol{\beta}}_{RE})$$
- $$Cov[(\mathbf{b}_{FE} - \hat{\boldsymbol{\beta}}_{RE}), \hat{\boldsymbol{\beta}}_{RE}] = Cov(\mathbf{b}_{FE}, \hat{\boldsymbol{\beta}}_{RE}) - Var(\hat{\boldsymbol{\beta}}_{RE}) = 0$$
- Insertando este resultado en la ecuación, se obtiene la matriz de varianzas y covarianzas para el contraste de Wald. Para  $\hat{\Psi}$  se utiliza la matriz de varianzas y covarianzas del estimador de las pendientes en el modelo LSDV y la matriz estimada en el modelo de efectos aleatorios (excluyendo el término constante)

$$Var(\mathbf{b}_{FE} - \hat{\boldsymbol{\beta}}_{RE}) = Var(\mathbf{b}_{FE}) - Var(\hat{\boldsymbol{\beta}}_{RE}) = \boldsymbol{\Psi}$$

$$W = (\mathbf{b}_{FE} - \hat{\boldsymbol{\beta}}_{RE})' \hat{\boldsymbol{\Psi}}^{-1} (\mathbf{b}_{FE} - \hat{\boldsymbol{\beta}}_{RE}) \sim \chi^2_{K-1}$$

- El contraste de Hausman es una herramienta útil para determinar la especificación preferida del modelo de efectos comunes. Un problema de este contraste, sin embargo, es que el estadístico puede ser negativo (no se garantiza que la diferencia sea definida positiva), por lo que se puede considerar la siguiente forma del contraste, la cual asegura una matriz definida no positiva y  $\mathbf{b}_{means}$  es el estimador de medias entre grupos:

$$H' = (\mathbf{b}_{FE} - \mathbf{b}_{means})' [Asy.Var(\mathbf{b})]^{-1} (\mathbf{b}_{FE} - \mathbf{b}_{means})$$

$$where \ Asy.Var(\mathbf{b}) = Asy.Var(\mathbf{b}_{FE}) + Asy.Var(\mathbf{b}_{means})$$

- Escoger entre el modelo de efectos fijos o aleatorios es un dilema, dado que ambas tienen aspectos negativos, tales como la proliferación de parámetros y los regresores invariantes en el tiempo para el modelo de efectos fijos, o la no correlación de la heterogeneidad no observada con los regresores del modelo de efectos aleatorios

- Mundlak, en su investigación de 1978, sugirió una especificación para la esperanza de la heterogeneidad no observada, de modo que el modelo de efectos aleatorios queda de la siguiente manera:

$$E(c_i | \mathbf{X}_i) = \bar{\mathbf{x}}_i' \boldsymbol{\gamma}$$

$$\begin{aligned} y_{it} &= \mathbf{x}_{it}' \boldsymbol{\beta} + c_i + \varepsilon_{it} = \mathbf{x}_{it}' \boldsymbol{\beta} + \bar{\mathbf{x}}_i' \boldsymbol{\gamma} + \varepsilon_{it} + (c_i - E(c_i | \mathbf{X}_i)) = \\ &= \mathbf{x}_{it}' \boldsymbol{\beta} + \bar{\mathbf{x}}_i' \boldsymbol{\gamma} + \varepsilon_{it} + u_i \end{aligned}$$

- Esta preserva la especificación del modelo de efectos aleatorios, pero no se puede lidiar con el problema de la correlación de los efectos y los regresores
- El término  $\bar{\mathbf{x}}_i' \boldsymbol{\gamma}$  solo incluye variables que varían en el tiempo, y la diferencia entre el modelo de efectos fijos y el de efectos aleatorios es el vector diferente de cero  $\boldsymbol{\gamma}$ . De este modo, un test estadístico para  $\boldsymbol{\gamma} = \mathbf{0}$  proporciona un enfoque alternativo para el test de Hausman
- El tratamiento de los datos de panel hasta ahora ha asumido que los parámetros del modelo son constantes, y los interceptos varían aleatoriamente

entre los grupos, de modo que la heterogeneidad se introduce a través de la variación del término constante

- La heterogeneidad de los parámetros entre individuos o grupos se puede modelar como una variación estocástica. Por lo tanto, se puede plantear el siguiente modelo:

$$y_i = X_i \beta_i + \varepsilon_i \quad E(\varepsilon_i | X_i) = 0 \quad E(\varepsilon_i \varepsilon_i' | X_i) = \sigma_\varepsilon^2 I_T$$

$$\beta_i = \beta + u_i \quad E(u_i | X_i) = 0 \quad E(u_i u_i' | X_i) = \Gamma$$

- Asumiendo que no hay autocorrelación o correlación transversal entre observaciones de las perturbaciones y que  $T > K$  (se puede obtener la regresión para cada grupo), entonces  $\beta_i$  es la realización de un proceso aleatorio con media  $\beta$  y matriz de varianzas y covarianzas  $\Gamma$
- Expandiendo los resultados en el modelo de regresión, se puede obtener el siguiente modelo general:

$$y_i = X_i \beta + (\varepsilon_i + X_i u_i)$$

$$\Omega_{ii} = E[(y_i - X_i \beta)(y_i - X_i \beta)' | X_i] = \sigma_\varepsilon^2 I_T + X_i \Gamma X_i'$$

- Para el sistema completo, la matriz de varianzas y covarianzas de las perturbaciones es diagonal en bloque, con un bloque diagonal  $\Omega_{ii}$  de tamaño  $T \times T$
- Se puede escribir el estimador GLS como una media ponderada de matrices de los estimadores mínimos cuadrados ordinarios específicos para cada grupo

$$\hat{\beta} = (X' \Omega^{-1} X)^{-1} X' \Omega^{-1} y = \sum_{i=1}^n w_i b_i$$

$$\text{where } w_i = \left[ \sum_{i=1}^n (\Gamma + \sigma_\varepsilon^2 (X_i' X_i)^{-1})^{-1} \right]^{-1} (\Gamma + \sigma_\varepsilon^2 (X_i' X_i)^{-1})^{-1}$$

- La implementación del modelo requiere un estimador para  $\Gamma$ . Un enfoque es utilizar la varianza empírica del conjunto  $n$  de estimadores de mínimos cuadrados menos el valor medio de  $s_i^2 (X_i' X_i)^{-1}$

$$G = [1/(n-1)] \left( \sum_i b_i b_i' - n \bar{b} \bar{b}' \right) - (1/N) \sum_i v_i$$

$$\bar{\mathbf{b}} = \frac{1}{n} \sum_i \mathbf{b}_i \quad \mathbf{V}_i = s_i^2 (\mathbf{X}_i' \mathbf{X}_i)^{-1}$$

- No obstante, la matriz puede no ser definida positiva, de modo que se tendría que sacar el segundo término
- Un contraste chi cuadrado para contrastar el modelo de coeficientes aleatorios ante la alternativa de la regresión clásica (sin aleatoriedad de los coeficientes) se puede basar en lo siguiente:

$$C = \sum_i (\mathbf{b}_i - \mathbf{b}_*)' \mathbf{V}_i (\mathbf{b}_i - \mathbf{b}_*) \sim \chi_{(n-1)K}^2$$

$$\text{where } \mathbf{b}_* = \left( \sum_i \mathbf{V}_i^{-1} \right)^{-1} \sum_i \mathbf{V}_i^{-1} \mathbf{b}_i$$

- El mejor estimador lineal no sesgado para predictores individuales de los vectores de coeficientes de grupos específicos es una media ponderada de matrices del estimador GLS y del estimador de mínimos cuadrados

$$\hat{\boldsymbol{\beta}}_i = \mathbf{Q}_i \hat{\boldsymbol{\beta}} + (\mathbf{I} - \mathbf{Q}_i) \mathbf{b}_i$$

$$\text{where } \mathbf{Q}_i = [(1/s_i^2) \mathbf{X}_i' \mathbf{X}_i + \mathbf{G}^{-1}]^{-1} \mathbf{G}^{-1}$$

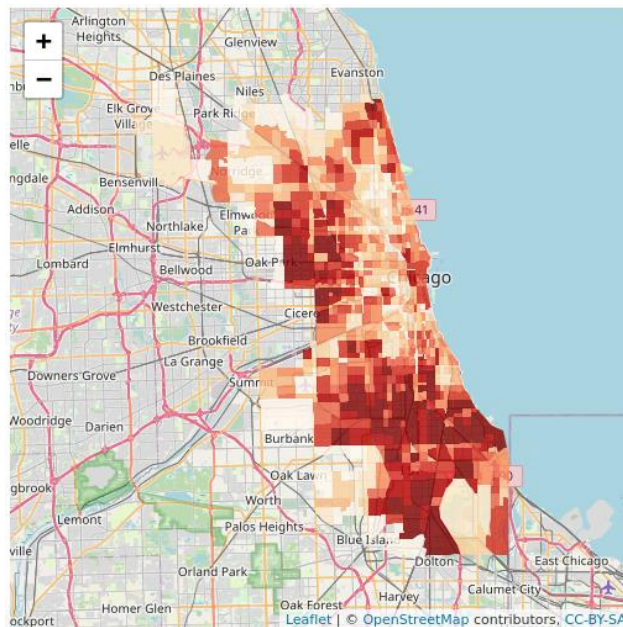
- Muchos investigadores utilizan un enfoque de dos pasos para estimar modelos de dos niveles. En una forma común de la aplicación, se utiliza un conjunto de datos de panel para poder estimar el siguiente modelo:

$$\mathbf{y}_{it} = \mathbf{x}_{it}' \boldsymbol{\beta}_i + \varepsilon_{it} \quad \text{for } i = 1, \dots, n \text{ and } t = 1, \dots, T$$

$$\boldsymbol{\beta}_{i,k} = \mathbf{z}_i' \boldsymbol{\alpha}_k + \mathbf{u}_{i,k} \quad \text{for } i = 1, \dots, n$$

- Asumiendo que el panel es lo suficientemente grande, la primera ecuación se estima  $n$  veces (una para cada individuo) y entonces el coeficiente estimado de  $x_{itk}$  en cada regresión forma una observación para la regresión en el segundo paso
- Los efectos de la clusterización mencionados anteriormente se justifican por la esperanza de que los efectos en localizaciones vecinas se comparten entre estas, creando una especie de correlación entre el espacio (más que en el tiempo) la cual se denomina autocorrelación espacial

- La autocorrelación espacial mide el grado en el que un fenómeno de interés está correlacionado a si mismo en el espacio
  - En otras palabras, valores similares aparecen más juntos uno al otro (o clusterizan) en el espacio (autocorrelación espacial positiva) o los valores vecinos son diferentes (autocorrelación espacial negativa). Una autocorrelación espacial nula indica que el patrón espacial es aleatorio (el espacio no afecta)
  - Para poder tener una idea visual sobre si es necesario considerar la autocorrelación, se puede utilizar un gráfico (mapa) que indique la densidad de los valores de la variable dependiente para cada localización. De este modo, si se puede identificar algún patrón según la cercanía, podría ser necesario considerar la autocorrelación espacial en los datos



- Es posible expresar la existencia de autocorrelación espacial para dos observaciones de una variable aleatoria  $y_i$  y  $y_j$  en las localizaciones  $i$  y  $j$

$$Cov(y_i, y_j) \neq 0 \text{ for } i \neq j$$

- Un modelo de autocorrelación espacial toma la misma forma familiar de las regresiones de datos de panel, pero se define el error de manera diferente para inducir autocorrelación espacial. Este tipo de modelos que introducen la autocorrelación a través del error se denominan modelos de errores espaciales (SEM)

$$y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + \varepsilon_{it} + u_i \text{ for } i = 1, 2, \dots, n \text{ and } t = 1, 2, \dots, T$$

- En esta representación, el subíndice  $i$  representa los individuos o grupos como antes, pero en este caso se pueden interpretar también como localizaciones. El término  $u_i$  sigue representando los efectos individuales o de grupos, pero  $\varepsilon_{it}$  permite que se tenga una estructura de autocorrelación espacial si se define de la siguiente manera:

$$\varepsilon_{it} = \lambda \sum_{j=1}^n W_{ij} \varepsilon_{jt} + v_t \quad \text{for } i = 1, 2, \dots, n$$

- El escalar  $\lambda$  es el coeficiente de autocorrelación espacial, y los elementos  $W_{ij}$  son ponderaciones espaciales o contiguas que se conocen (por suposición). Los elementos que aparecen en la suma de arriba serán una fila del peso espacial o matriz de contigüidad  $W$  de tamaño  $n \times n$ , de modo que para las  $n$  unidades se tiene la siguiente igualdad:

$$\varepsilon_t = \lambda W \varepsilon_t + v_t = \lambda W \varepsilon_t + v_t \mathbf{i}$$

- Asumiendo que  $|\lambda| < 1$  y que los elementos de la matriz  $W$  son tales que  $(I_n - \lambda W)$  sea una matriz invertible, es posible reescribir  $\varepsilon_t$  como  $\varepsilon_t = (I_n - \lambda W)^{-1} v_t$  y expresar el modelo de la siguiente manera:

$$y_t = X_t \beta + (I_n - \lambda W)^{-1} v_t + u$$

- Adicionalmente, se asume que  $u_i$  y  $v_i$  tienen medias nulas y varianzas constantes  $\sigma_u^2$  y  $\sigma_v^2$  y que son independientes entre individuos y entre ellos. Por lo tanto, el modelo de regresión generalizado que aplica a las  $n$  observaciones en un momento  $t$  será el siguiente:

$$E(y_t | X_t) = X_t \beta$$

$$Var(y_t | X_t) = (I_n - \lambda W)^{-1} (\sigma_v^2 \mathbf{i} \mathbf{i}') (I_n - \lambda W)^{-1} + \sigma_u^2 I_n$$

- La estimación y la interpretación se puede hacer como en otros modelos de datos de panel, aunque no hay un estimador natural para  $\lambda$  (a veces se asume normalidad y se usan métodos de máxima verosimilitud)
- Otro tipo de modelo muy utilizado es uno que funciona como un modelo autorregresivo para series temporales, el cual se denomina modelos espaciales autorregresivos (SAR). Este tipo de modelos fueron propuestos por Anselin en 1988 como extensiones del modelo espacial a regresiones dinámicas, y hay varias especificaciones

$$y_{it} = \mathbf{x}'_{it}\boldsymbol{\beta} + \varepsilon_{it} + u_i \quad \text{for } i = 1, 2, \dots, n \text{ and } t = 1, 2, \dots, T$$

- Un modelo recursivo de espacio puro o *pure space-recursive model* especifica que la autocorrelación pertenece a los vecinos en el periodo anterior

$$y_{it} = \gamma[\mathbf{W}\mathbf{y}_{t-1}]_i + \mathbf{x}'_{it}\boldsymbol{\beta} + \varepsilon_{it}$$

$$\text{for } i = 1, 2, \dots, n \text{ and } t = 1, 2, \dots, T$$

- Un modelo recursivo de espacio-tiempo o *time-space recursive model* especifica que la dependencia es puramente autorregresiva respecto a los vecinos en el periodo pasado

$$y_{it} = \rho y_{i,t-1} + \gamma[\mathbf{W}\mathbf{y}_{t-1}]_i + \mathbf{x}'_{it}\boldsymbol{\beta} + \varepsilon_{it}$$

$$\text{for } i = 1, 2, \dots, n \text{ and } t = 1, 2, \dots, T$$

- Un modelo simultáneo de espacio-tiempo o *time-space simultaneous model* especifica que la dependencia espacial es con respecto a los vecinos en el periodo actual

$$y_{it} = \rho y_{i,t-1} + \lambda[\mathbf{W}\mathbf{y}_t]_i + \mathbf{x}'_{it}\boldsymbol{\beta} + \varepsilon_{it}$$

$$\text{for } i = 1, 2, \dots, n \text{ and } t = 1, 2, \dots, T$$

- Finalmente, un modelo dinámico de espacio-tiempo o *time-space dynamic model* especifica que la autorregresión depende de los vecinos en el periodo actual y en el anterior

$$y_{it} = \rho y_{i,t-1} + \lambda[\mathbf{W}\mathbf{y}_t]_i + \gamma[\mathbf{W}\mathbf{y}_{t-1}]_i + \mathbf{x}'_{it}\boldsymbol{\beta} + \varepsilon_{it}$$

$$\text{for } i = 1, 2, \dots, n \text{ and } t = 1, 2, \dots, T$$

- En este caso, serán necesarios métodos de estimación para modelos de datos de panel y la interpretación de los coeficientes en la regresión variarán según el término
- El coeficiente  $\gamma$  es el coeficiente autorregresivo espacial para el periodo anterior, mientras que  $\lambda$  es el mismo coeficiente, pero para el periodo actual (son las correlaciones). El signo del coeficiente indica la dirección de la autocorrelación espacial. El coeficiente  $\rho$  se puede interpretar como el efecto del valor pasado en el presente (como en los modelos autorregresivos)
- Los modelos presentados sirven para modelar datos de panel, pero también tienen una contraparte para datos de sección cruzada:



$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i \text{ where } \varepsilon_i = \lambda \sum_{j=1}^n W_{ij} \varepsilon_j + v_i \text{ for } i = 1, 2, \dots, n \text{ (SEM)}$$

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \lambda \sum_{j=1}^n W_{ij} y_j + \varepsilon_i \text{ for } i = 1, 2, \dots, n \text{ (SAR)}$$

- Con notación matricial, los modelos se pueden expresar equivalentemente:

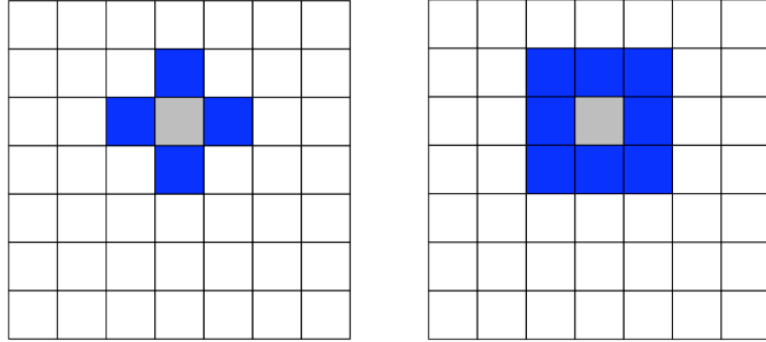
$$\mathbf{y} = \mathbf{X}' \boldsymbol{\beta} + \boldsymbol{\varepsilon} \text{ where } \boldsymbol{\varepsilon} = \lambda \mathbf{W} \boldsymbol{\varepsilon} + \mathbf{v} = \lambda \mathbf{W} \boldsymbol{\varepsilon} + \mathbf{v}_i \mathbf{i} \text{ (SEM)}$$

$$\mathbf{y} = \mathbf{x}_i' \boldsymbol{\beta} + \lambda \mathbf{W} \mathbf{y} + \boldsymbol{\varepsilon} = \mathbf{x}_i' \boldsymbol{\beta} + \lambda \mathbf{W} \mathbf{y} + \varepsilon_i \mathbf{i} \text{ (SAR)}$$

- Para el modelo SAR presentado para datos cruzados, el término  $\mathbf{W} \mathbf{y}$  induce a una correlación no nula con el término de error, ya que los términos siempre están correlacionados independientemente de la estructura de los errores. Esto implica que los estimadores MCO en el modelo no espacial estarán sesgados y serán inconsistentes
  - Si el modelo SEM es real, entonces los estimadores MCO para el modelo no espacial no están sesgados, pero son ineficientes, dado que la correlación proviene de  $\boldsymbol{\varepsilon}$
  - En este caso,  $\lambda$  es el coeficiente de autocorrelación espacial en el modelo SAR y su signo es la dirección de la autocorrelación, mientras que en el modelo SEM la interpretación sigue siendo la misma que para un modelo común de datos cruzados
  - No se consideran las otras especificaciones para los modelos SAR debido a que todas se diferencian al incluir diferentes momentos, y en datos cruzados el tiempo no se tiene en cuenta
- La estructura del modelo proviene de la matriz simétrica de ponderaciones  $\mathbf{W}$ . Esta matriz puede representarse típicamente de varias maneras:
- La primera manera consiste en que  $W_{ij}$  sea 1 cuando un par  $i, j$  son vecinos, mientras que será cero si no lo son para ese par. De este modo, en verdad se impone una restricción sobre cada punto, la cual es la definición de un conjunto de vecindad (en donde, por convención, no se incluyen los pares  $(i, i)$ )

$$W_{ij} = \begin{cases} 1 & \text{if } (i, j) \in N(i, j) \\ 0 & \text{otherwise} \end{cases} \text{ for } \forall (i, j)$$

- La especificación del conjunto de vecindad es un poco arbitraria y hay bastantes maneras de poder hacerlo. Dos criterios muy utilizados son el de la torre (dos unidades son vecinas si comparten un costado de un cuadrado en el espacio) y el de la reina (dos unidades son vecinas si comparten un costado o una esquina)



- Otro criterio usado es el de clasificar dos observaciones o localizaciones como vecinas si están a una cierta distancia, de modo que se cumpla la siguiente condición:

$$(i, j) \in N(i, j) \text{ if } |i - j| < d_{max}$$

- Finalmente, se puede hacer que  $W_{ij}$  sea una distancia en el espacio, de modo que  $W_{ij}$  decrece cuanto más aumente  $|i - j|$  (la ponderación de la autocorrelación baja cuanto más aumenta la lejanía entre localizaciones). En este enfoque, las filas de cada matriz sumarán 1 (en otras palabras,  $\sum_i W_{ij} = 1$ )
- Un paso natural es el de utilizar un contraste hipótesis para los efectos espaciales. El procedimiento estándar para una regresión de datos de sección cruzada es el estadístico  $I$  de Moran, que se calcularía para cada conjunto de residuos  $e_t$

$$I_t = \frac{n \sum_{i=1}^n \sum_{j=1}^n W_{ij} (e_{it} - \bar{e}_t) (e_{jt} - \bar{e}_t)}{(\sum_{i=1}^n \sum_{j=1}^n W_{ij}) \sum_{i=1}^n (e_{it} - \bar{e}_t)^2} = \frac{e_t' W e_t}{e_t' e_t}$$

- Como se puede ver, este estadístico es para conjuntos de cada  $t = 1, 2, \dots, T$ , de modo que se puede obtener el mismo estadístico  $I$  para datos cruzados

$$I = \frac{n \sum_{i=1}^n \sum_{j=1}^n W_{ij} (e_i - \bar{e}) (e_j - \bar{e})}{(\sum_{i=1}^n \sum_{j=1}^n W_{ij}) \sum_{i=1}^n (e_i - \bar{e})^2} = \frac{e' W e}{e' e}$$

- El estadístico  $I$  de Moran es una medida de autocorrelación espacial, y este puede tomar valores entre  $-1$  (dispersión perfecta) y  $1$  (correlación perfecta)
- Para un panel de  $T$  conjuntos de observaciones independientes, el estadístico  $\bar{I} = \frac{1}{T} \sum_{t=1}^T I_t$  usaría todo el conjunto de información. Una aproximación asintótica para la varianza de este estadístico bajo la hipótesis nula:

$$V^2 = \frac{1}{T} \frac{n^2 \sum_{i=1}^n \sum_{j=1}^n W_{ij}^2 + 3 \left( \sum_{i=1}^n \sum_{j=1}^n W_{ij} \right)^2 - n \sum_{i=1}^n \left( \sum_{j=1}^n W_{ij} \right)^2}{(n^2 - 1) \left( \sum_{i=1}^n \sum_{j=1}^n W_{ij} \right)^2}$$

- El estadístico  $\bar{I}/V$  se distribuirá como una distribución normal bajo la hipótesis nula y se puede utilizar un estadístico de contraste. La hipótesis nula que el estadístico  $\bar{I}$  es nulo (lo cual quiere decir que la autocorrelación espacial es nula) mientras que tiene como hipótesis alternativa un valor diferente (habiendo así correlación espacial)
- Baltagi, Song and Koh, en su investigación de 2003, identificaron una variedad de contrastes de multiplicadores de Lagrange bajo la suposición de normalidad
- Para datos cruzados y autocorrelación espacial en el error (para un modelo SEM), se puede utilizar el siguiente contraste, cuya hipótesis nula es que no hay un término  $We$  en el modelo real para  $y$  y la alternativa es que si hay (como en un SEM):

$$LM(1) = \frac{(e'We/s^2)^2}{tr(W'W + W^2)} \quad \text{where } s^2 = e'e/n$$

- Para datos cruzados y para variables dependientes retrasadas (*lagged dependent variables*), se puede utilizar el siguiente contraste, cuya hipótesis nula es que no hay un término  $Wy$  en el modelo real para  $y$  y la alternativa es que si hay (como en un SAR):

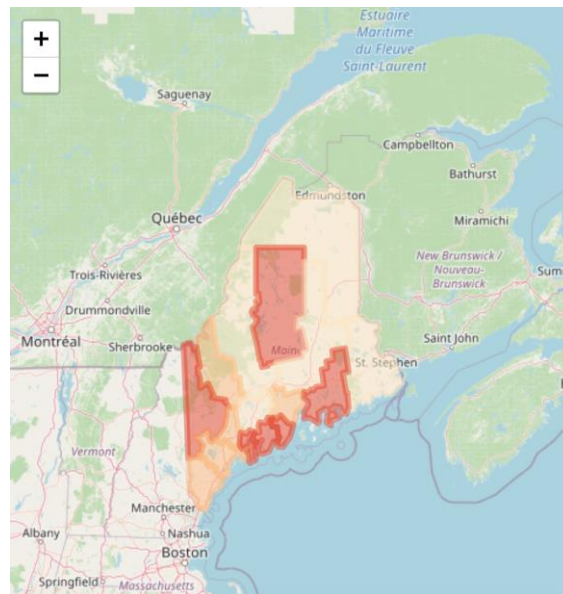
$$LM(2) = \frac{(e'Wy/s^2)^2}{b'X'WMWXb/s^2 + tr(W'W + W^2)}$$

$$\text{where } M = I - X(X'X)^{-1}X'$$

- En este caso, rechazar la hipótesis nula quiere decir que hay evidencia sobre autocorrelación espacial (la hipótesis alternativa), mientras que no rechazar la hipótesis nula quiere

decir que la autocorrelación sería aleatoria (no habría autocorrelación espacial)

- Una manera visual de poder tener una idea del ajuste del modelo a los datos es a través de hacer un gráfico con la densidad de los residuos (como se hizo antes con la variable dependiente)



- Si se detecta que hay algún patrón de clusterización para zonas con residuos significativos, entonces se puede interpretar como una señal de que el ajuste del modelo no es muy bueno porque falta algún factor que explique este patrón (autocorrelación espacial)

## Los resultados binarios y la elección discreta

- El análisis de la elección individual que es foco de la microeconometría se basa en hacer modelar resultados discretos. De este modo, la variable dependiente no es una medida cuantitativa de un resultado económico, sino un indicador de que algo ha ocurrido o no
  - Por lo tanto, los modelos de regresión vistos hasta ahora no son apropiados, de modo que se tiene que pasar al modelaje de probabilidades y de conteo de ocurrencias
    - Los modelos estudiados en esta área son naturalmente no lineales
    - Estos se basan en un modelo de utilidad aleatoria para las elecciones observadas, de modo que el decisor se enfrenta a una situación o conjunto de alternativas y revela algo de sus preferencias subyacentes por la elección que hace. Estas

elecciones hechas se ven afectadas por influencias observables y por características inobservables del decisor

- Los modelos de elección discreta más comunes son el de elección binaria, el de elección multinomial, el de elección ordenada y el de conteo de eventos
  - En la elección binaria, los individuos se enfrentan a dos decisiones y se toma la decisión que proporciona la mayor utilidad. En este contexto, el valor de las variables binarias solo son indicadores numéricos para la elección
  - En la elección multinomial, los individuos se enfrentan a más de dos decisiones y se toma la decisión que proporciona la mayor utilidad. Esta es una variación pequeña de los modelos de elección binaria, pero modelos más elaborados permiten una mayor especificación de las preferencias de los individuos. Además, los valores de la variable dependiente son solo indicadores numéricos para la elección
  - En la elección ordenada, los individuos revelan la fuerza de sus preferencias con respecto a un único resultado. En este contexto, los valores de la variable tienen significado, dado que tienen un orden que revelan la fuerza de las preferencias, y la variable se considera una medida cuantitativa
  - En los modelos de conteo de eventos, el resultado observado es un conteo del número de ocurrencias de un evento. En muchos casos, estos son similar que, en los otros tres contextos, dado que la variable dependiente mide la elección individual, pero en otros casos, este conteo puede pertenecer a un proceso natural. Por lo tanto, hay contextos en los que se hace un modelaje de regresión familia, pero en los más generales se adapta el modelaje para acomodar una variable dependiente discreta y no negativa (a través de probabilidades)
- Con tal de estudiar el comportamiento individual, se construyen modelos que vinculan la decisión o resultado a una serie de factores. El enfoque es analizar cada uno en un marco general de modelos de probabilidad:

$$P(\text{event } j \text{ occurs} | \mathbf{x}_i) = P(Y_i = j | \mathbf{x}_i) = F(\text{relevant effects, parameters, } \mathbf{x}_i)$$

- El estudio de la elección cualitativa se enfoca en una especificación y estimación apropiada, y utiliza modelos para las probabilidades de eventos, en donde los eventos son las elecciones de los individuos entre dos o más alternativas la mayoría de veces

- Por lo tanto, se denota por  $P(Y_i = 1|x_i)$  la probabilidad de que el evento de interés ocurra dado  $x_i$  y  $P(Y_i = 0|x_i)$  es la probabilidad de que el evento no ocurra
- Una interpretación de los datos sobre la elección individual es proporcionada por el modelo de utilidad aleatorio. Si  $U_a$  y  $U_b$  representan la utilidad de un individuo entre dos elecciones diferentes, la elección observada revela cuál de las dos opciones proporciona mayor utilidad, pero no la magnitud de la utilidad subyacente
  - Por lo tanto, el indicador observado es igual a 1 si  $U_a > U_b$  y 0 si  $U_a < U_b$ , entonces  $Y = \mathbf{1}(U > 0)$  si  $U = U_a - U_b$  (lo mismo que el caso de censura visto anteriormente). Una formulación común es el modelo lineal de utilidad aleatoria:

$$U_a = \mathbf{w}'\boldsymbol{\beta}_a + \mathbf{z}'_a\boldsymbol{\gamma}_a + \varepsilon_a \quad U_b = \mathbf{w}'\boldsymbol{\beta}_b + \mathbf{z}'_b\boldsymbol{\gamma}_b + \varepsilon_b$$

- El vector observable de características del individuo se denota  $\mathbf{w}$  (puede incluir ingresos, género, etc.). Los vectores  $\mathbf{z}_a$  y  $\mathbf{z}_b$  denota atributos de las dos elecciones (que pueden ser específicas de la elección), los términos aleatorios  $\varepsilon_a$  y  $\varepsilon_b$  representan los elementos estocásticos que son específicos y solo lo saben los individuos, pero no el observador
- La completitud del modelo para la determinación de resultados observados es la revelación del *ranking* de las preferencias por la elección del individuo. Entonces, si se denota por  $Y = 1$  la elección de consumo de la alternativa  $a$ , se infiere que  $U_a > U_b$  y como el resultado depende de elementos aleatorios, entonces se obtiene el siguiente modelo:

$$\begin{aligned} P(Y = 1|\mathbf{w}, \mathbf{z}_a, \mathbf{z}_b) &= P(U_a > U_b) = \\ &= P[(\mathbf{w}'\boldsymbol{\beta}_a + \mathbf{z}'_a\boldsymbol{\gamma}_a + \varepsilon_a) - (\mathbf{w}'\boldsymbol{\beta}_b + \mathbf{z}'_b\boldsymbol{\gamma}_b + \varepsilon_b) > 0|\mathbf{w}, \mathbf{z}_a, \mathbf{z}_b] = \\ &= P[\mathbf{w}'(\boldsymbol{\beta}_a - \boldsymbol{\beta}_b) - (\mathbf{z}'_a\boldsymbol{\gamma}_a - \mathbf{z}'_b\boldsymbol{\gamma}_b) + (\varepsilon_a + \varepsilon_b) > 0|\mathbf{w}, \mathbf{z}_a, \mathbf{z}_b] = \\ &= P(\mathbf{x}'\boldsymbol{\beta} + \varepsilon > 0|\mathbf{x}) \end{aligned}$$

- En este caso  $\mathbf{x}'\boldsymbol{\beta}$  incluye todos los términos observables de la diferencia de las dos funciones de utilidad y  $\varepsilon$  denota la diferencia entre dos elementos aleatorios
- Los modelos de variable dependiente discreta normalmente se denominan modelos de función de índice. El resultado de una elección discreta es un reflejo de una regresión subyacente, modelando la diferencia entre la utilidad de las opciones con una regresión de la variable no observada  $y^*$

$$y^* = \mathbf{x}'\boldsymbol{\beta} + \varepsilon$$

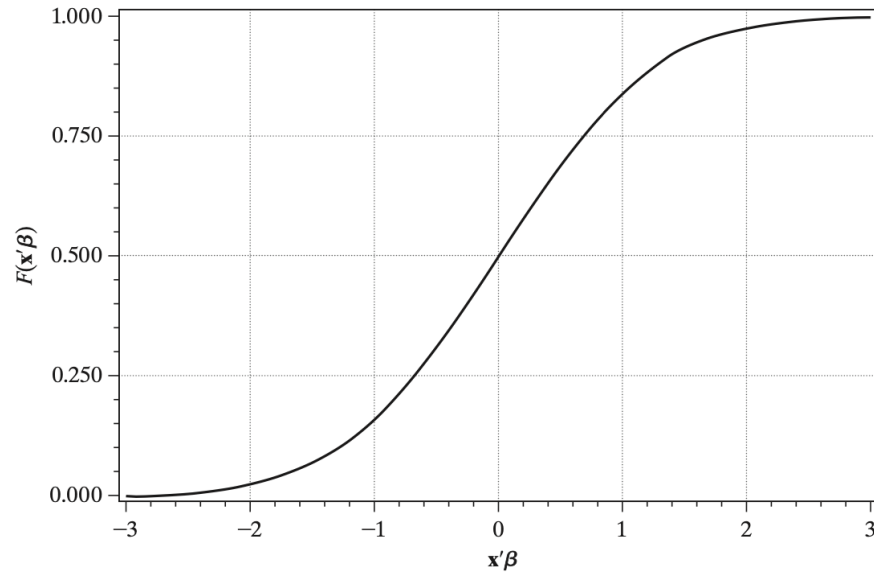
- Este resultado es la utilidad neta mostrada anteriormente, y se asume que  $\varepsilon$  tiene media nula (de modo que en  $\mathbf{x}$  tiene que haber un término constante) y normalmente tiene una distribución logístico u otra distribución conocido
- Como esta variable no se observa, sino que se observa  $y$ , entonces se puede deducir que se observa  $y = 1$  si  $y^* > 0$  y que  $y = 0$  si  $y^* \leq 0$ . Esto se suele denotar por  $\mathbf{1}(y^* > 0)$ . En esta formulación,  $\mathbf{x}'\boldsymbol{\beta}$  se denomina la función índice
- La suposición de varianza conocida de  $\varepsilon$  es una normalización, pero si se considera un parámetro  $\sigma_\varepsilon^2$  no restringido, entonces la regresión latente será  $y^* = \mathbf{x}'\boldsymbol{\beta} + \sigma\varepsilon^*$ , donde ahora  $\varepsilon$  tiene una varianza de 1
  - No obstante,  $(y^*/\sigma) = \mathbf{x}'(\boldsymbol{\beta}/\sigma) + \varepsilon$  es el mismo modelo con los mismos datos. Los datos observados  $y$  no cambian, dependiendo solo del signo de  $y^*$  y no de su escala como se puede ver
  - Esto significa que no hay información sobre  $\sigma$  en los datos muestrales y que, por tanto, no se puede estimar y el vector de parámetros  $\boldsymbol{\beta}$  solo se identifica hasta la escala
- La suposición de que el umbral para la observación es nulo tampoco es realista si el término constante se incluye en la regresión (aunque si lo es si este no se incluye)
  - Siendo  $a$  un umbral diferente de cero y  $\alpha$  un término constante desconocido para la regresión, y pensando que  $\mathbf{x}$  y  $\boldsymbol{\beta}$  incluyen el resto del índice (sin tener un término constante), entonces la probabilidad de  $y$  se puede expresar de la siguiente manera:

$$\begin{aligned} P(y^* > a|\mathbf{x}) &= P(\alpha + \mathbf{x}'\boldsymbol{\beta} + \varepsilon > a|\mathbf{x}) = \\ &= P((\alpha - a) + \mathbf{x}'\boldsymbol{\beta} + \varepsilon > 0|\mathbf{x}) \end{aligned}$$

- Debido a que  $\alpha$  es desconocido, entonces la diferencia  $(\alpha - a)$  es un parámetro desconocido. El resultado final es que, si la regresión contiene un término constante, entonces este no cambia por la elección del umbral, por lo que fijarlo en cero es una normalización sin importancia. Esto permite obtener la siguiente igualdad:

$$P(y^* > 0|\mathbf{x}) = P(\varepsilon > -\mathbf{x}'\boldsymbol{\beta}|\mathbf{x})$$

- Dependiendo del tipo de distribución que se asume para el modelo, entonces se puede obtener una forma funcional para la función de distribución de probabilidad para  $\varepsilon$



- Se supone que un conjunto de factores del vector  $x$  explica la decisión del individuo, mientras que  $\beta$  es un vector de parámetros que refleja el impacto de cambios en  $x$  en la probabilidad. Por lo tanto, solo se requiere una especificación adecuada para la parte derecha de la ecuación

$$P(Y = 1|x) = F(x, \beta) \quad P(Y = 0|x) = 1 - F(x, \beta)$$

- El único requerimiento es que se produzcan predicciones consistentes con la teoría subyacente vista anteriormente. Por lo tanto, para un vector de regresores, se espera que se cumplan las siguientes condiciones:

$$0 \leq P(Y = 1|x) \leq 1$$

$$\lim_{x'\beta \rightarrow -\infty} P(Y = 1|x) = 0 \quad \lim_{x'\beta \rightarrow \infty} P(Y = 1|x) = 1$$

- En principio, cualquier distribución definida sobre la recta real es válida, pero si esta distribución es simétrica (como la logística o la normal), entonces se puede obtener la siguiente igualdad:

$$P(y^* > 0|x) = P(\varepsilon > -x'\beta|x) = P(\varepsilon < x'\beta|x) = F(x'\beta)$$

- Algunas de las distribuciones más utilizadas son la normal, la logística o la distribución de Gumbel, las cuales dan paso a diferentes modelos. No obstante, los marcos más utilizados son los dos primeros



- Si se utiliza la distribución normal, se da paso al modelo probit (con distribución simétrica):

$$P(Y = 1|\mathbf{x}) = \int_{-\infty}^{\mathbf{x}'\boldsymbol{\beta}} \phi(t) dt = \Phi(\mathbf{x}'\boldsymbol{\beta})$$

- Si se utiliza la distribución logística, se da paso al modelo logit (con distribución simétrica):

$$P(Y = 1|\mathbf{x}) = \int_{-\infty}^{\mathbf{x}'\boldsymbol{\beta}} \Lambda(t)[1 - \Lambda(t)] dt = \frac{\exp(\mathbf{x}'\boldsymbol{\beta})}{1 + \exp(\mathbf{x}'\boldsymbol{\beta})} = \Lambda(\mathbf{x}'\boldsymbol{\beta})$$

- Si se utiliza la distribución de Gumbel o el modelo de valores extremos de tipo I, se obtiene el siguiente modelo:

$$P(Y = 1|\mathbf{x}) = \exp[-\exp(-\mathbf{x}'\boldsymbol{\beta})]$$

- El modelo de log-log complementario resta el modelo de Gumbel a la unidad:

$$P(Y = 1|\mathbf{x}) = 1 - \exp[-\exp(-\mathbf{x}'\boldsymbol{\beta})]$$

- El modelo de Burr es parecido al logit, solo que la razón resultante se eleva a un parámetro  $\gamma$

$$P(Y = 1|\mathbf{x}) = \left[ \frac{\exp(\mathbf{x}'\boldsymbol{\beta})}{1 + \exp(\mathbf{x}'\boldsymbol{\beta})} \right]^\gamma = [\Lambda(\mathbf{x}'\boldsymbol{\beta})]^\gamma$$

- Las colas de la distribución logística son más gruesas que la de la distribución normal, de modo que para valores intermedios se pueden obtener probabilidades similares, mientras que para valores en las colas se obtienen probabilidades mucho más grandes para la distribución logística que para la distribución normal
- La mayoría de análisis se basan en examinar las relaciones entre los regresores  $\mathbf{x}$  y la probabilidad del evento  $P(Y = 1|\mathbf{x}) = F(y|\mathbf{x}) = F(\mathbf{x}'\boldsymbol{\beta})$ , las cuales se interpretan a través de los efectos parciales. Como el modelo no es lineal, los parámetros  $\boldsymbol{\beta}$  no son necesariamente los efectos parciales de los modelos

$$\frac{\partial F(y|\mathbf{x})}{\partial \mathbf{x}} = \left[ \frac{\partial F(\mathbf{x}'\boldsymbol{\beta})}{\partial (\mathbf{x}'\boldsymbol{\beta})} \right] \boldsymbol{\beta} = f(\mathbf{x}'\boldsymbol{\beta})\boldsymbol{\beta}$$

- En el caso de la distribución normal y de la distribución logística se obtienen los siguientes resultados:

$$\frac{\partial F(y|x)}{\partial x} = \frac{\partial \Phi(x'\beta)}{\partial (x'\beta)} = \phi(x'\beta_{prob})\beta_{prob}$$

$$\frac{\partial F(y|x)}{\partial x} = \frac{\partial \Lambda(x'\beta)}{\partial (x'\beta)} = \Lambda(x'\beta_{log})[1 - \Lambda(x'\beta_{log})]\beta_{log}$$

- Como se puede ver, los efectos dependerán de los valores de  $x$ . En los modelos de función de índice, generalmente, el vector de efectos parciales es un múltiplo del vector de coeficientes
- En muchas aplicaciones, una regularidad empírica que suele aparecer es que  $1.6\hat{\beta}_{prob} \approx \hat{\beta}_{log}$ , lo cual puede sugerir una gran diferencia entre ambos modelos. Esto ocurre porque el efecto parcial del modelo *probit* para el valor medio sería  $0.4\hat{\beta}_{prob}$ , mientras que para el modelo *logit* el valor sería  $0.5(1 - 0.5)\hat{\beta}_{prob} = 0.25\hat{\beta}_{log}$ , de modo que, para el mismo valor medio, se obtiene que  $\hat{\beta}_{log} = (0.4/0.25)\hat{\beta}_{prob}$
- Existen varias maneras de poder obtener inferencias a partir de estimadores del efecto parcial
  - Para calcular los efectos parciales, normalmente se puede evaluar las expresiones en las medias muestrales de los datos, produciendo efectos parciales en las medias o *partial effects at averages* (PEA)

$$PEA = \hat{\gamma}(\bar{x}) = f(\bar{x}'\hat{\beta})\hat{\beta}$$

- Como las medias no siempre producen un escenario realista, es mejor calcular los efectos parciales en cada punto y obtener la media muestral de los efectos parciales individuales, produciendo los efectos parciales medios o *average partial effects*

$$APE = \hat{\gamma} = \frac{1}{n} \sum_{i=1}^n f(x'_i\hat{\beta})\hat{\beta}$$

- Normalmente el interés reside en el efecto parcial medio o *average partial effect*, que es el valor esperado del efecto parcial

$$APE^0 = \gamma^0 = E_x \left[ \frac{\partial E(y|x)}{\partial x} \right]$$

- Se pueden interpretar los parámetros del modelo utilizando diferentes enfoques

- Debido a que los parámetros no se pueden interpretar directamente como efectos parciales, solo se puede interpretar el signo de la relación con la probabilidad, ya que el signo del coeficiente determinará el signo del efecto parcial

$$\beta_k < 0 \Rightarrow \frac{\partial F(y|\mathbf{x})}{\partial x_k} < 0 \quad \beta_k > 0 \Rightarrow \frac{\partial F(y|\mathbf{x})}{\partial x_k} > 0$$

- Para una variable binaria cualquiera, es común calcular la diferencia de probabilidad con el caso base para analizar el efecto de pasar de una categoría a otra en términos de probabilidad

$$P(y|\mathbf{x}, d = 1) - P(y|\mathbf{x}, d = 0) = F(y|\mathbf{x}, d = 1) - F(y|\mathbf{x}, d = 0)$$

- También se puede comparar la importancia relativa del efecto parcial de una variable con el efecto parcial del de otra

$$\frac{\frac{\partial F(y|\mathbf{x})}{\partial x_k}}{\frac{\partial F(y|\mathbf{x})}{\partial x_l}} = \frac{f(\mathbf{x}'\boldsymbol{\beta})\beta_k}{f(\mathbf{x}'\boldsymbol{\beta})\beta_l} = \frac{\beta_k}{\beta_l}$$

- Una medida que ayuda a interpretar el resultado del modelo es la *odds ratio*, el cual se basa en la medida de *odds in favour*, la cual es una razón  $P(Y = 1|\mathbf{x})/P(Y = 0|\mathbf{x})$  que expresa las probabilidades a favor de que ocurra un evento en el modelo *logit*

$$Odds = \frac{P(Y = 1|\mathbf{x})}{P(Y = 0|\mathbf{x})} = \frac{\exp(\mathbf{x}'\boldsymbol{\beta})/[1 + \exp(\mathbf{x}'\boldsymbol{\beta})]}{1/[1 + \exp(\mathbf{x}'\boldsymbol{\beta})]} = \exp(\mathbf{x}'\boldsymbol{\beta})$$

- Considerando el efecto en las *odds* cuando cambia una variable binaria  $d$ , se obtiene la *odds ratio* (la razón entre las *odds* cuando se incrementa un regresor en una unidad y las *odds* del caso base):

$$\frac{Odds(\mathbf{x}, d = 1)}{Odds(\mathbf{x}, d = 0)} = \frac{\frac{\exp(\mathbf{x}'\boldsymbol{\beta} + \delta)/[1 + \exp(\mathbf{x}'\boldsymbol{\beta} + \delta)]}{1/[1 + \exp(\mathbf{x}'\boldsymbol{\beta} + \delta)]}}{\frac{\exp(\mathbf{x}'\boldsymbol{\beta})/[1 + \exp(\mathbf{x}'\boldsymbol{\beta})]}{1/[1 + \exp(\mathbf{x}'\boldsymbol{\beta})]}} = \exp(\delta)$$

- Por lo tanto, cuando una variable cambia una unidad, la *odds ratio* permite aproximar el efecto parcial de la variable, aunque no es la derivada. Es decir, el exponencial del coeficiente  $\beta_k$  representa el cambio multiplicativo en la *ratio* y  $\exp(\beta_k) - 1$ ,

por tanto, representa el cambio proporcional al aumentar en una unidad la variable  $x_k$

- El resultado binario sugiere un modelo de regresión  $F(\mathbf{x}, \boldsymbol{\beta}) = \mathbf{x}'\boldsymbol{\beta}$  con esperanza condicional  $E(y|\mathbf{x}) = 0[1 - F(\mathbf{x}, \boldsymbol{\beta})] + F(\mathbf{x}, \boldsymbol{\beta}) = F(\mathbf{x}, \boldsymbol{\beta})$ , lo cual implica el siguiente modelo de regresión:

$$y = E(y|\mathbf{x}) + [y - E(y|\mathbf{x})] = \mathbf{x}'\boldsymbol{\beta} + \varepsilon$$

- A este modelo se le llama el modelo de probabilidad lineal, aunque este tiene algunos problemas
- En este modelo,  $\varepsilon$  es heteroscedástico y depende de  $\boldsymbol{\beta}$ , lo cual se debe a la naturaleza binaria de la variable dependiente. Esto se puede ver fácilmente con la varianza, y se puede gestionar utilizando el estimador de mínimos cuadrados generalizado factible (aunque no soluciona el problema teórico, solo el de estimación)

$$\begin{cases} y_i = 1 & \text{if } \varepsilon_i = 1 - \mathbf{x}'\boldsymbol{\beta} \\ y_i = 0 & \text{if } \varepsilon_i = -\mathbf{x}'\boldsymbol{\beta} \end{cases}$$

$$\begin{aligned} \text{Var}(\varepsilon|\mathbf{x}) &= (1 - \mathbf{x}'\boldsymbol{\beta})^2 F(\mathbf{x}, \boldsymbol{\beta}) + (-\mathbf{x}'\boldsymbol{\beta})^2 (1 - F(\mathbf{x}, \boldsymbol{\beta})) = \\ &= (1 - \mathbf{x}'\boldsymbol{\beta})^2 \mathbf{x}'\boldsymbol{\beta} + (-\mathbf{x}'\boldsymbol{\beta})^2 (1 - \mathbf{x}'\boldsymbol{\beta}) = \mathbf{x}'\boldsymbol{\beta}(1 - \mathbf{x}'\boldsymbol{\beta}) \end{aligned}$$

- Un problema más grave es que no se garantiza que el valor de  $y$  esté acotado entre cero y uno (que se comporten como probabilidades). No obstante, a veces se utiliza porque permite obtener efectos parciales parecidos a los obtenidos con modelos *probit* y *logit* y por la robustez del modelo a suposiciones distribucionales
- Exceptuando el modelo de probabilidad lineal, la estimación de los modelos de elección binaria se lleva a cabo a través del método de máxima verosimilitud
  - Cada observación se trata como una realización de una distribución de Bernoulli, en donde la probabilidad de éxito es  $F(\mathbf{x}'\boldsymbol{\beta})$  y las observaciones son independientes. Por lo tanto, la función de verosimilitud es la siguiente:

$$P(Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n | \mathbf{X}) = \prod_{y_i=0} [1 - F(\mathbf{x}'_i \boldsymbol{\beta})] \prod_{y_i=1} F(\mathbf{x}'_i \boldsymbol{\beta})$$

- Esta se puede reexpresar de manera conveniente para una muestra de  $n$  observaciones y para el logaritmo de la función

$$L(\beta|data) = \prod_{i=0}^n [1 - F(x'_i\beta)]^{y_i} [F(x'_i\beta)]^{1-y_i}$$

$$\ln L = \sum_{i=0}^n [y_i \ln F(x'_i\beta) + (1 - y_i) \ln(1 - F(x'_i\beta))]$$

- Por lo tanto, el estimador para  $\beta$  será aquel que satisfazca la siguiente ecuación:

$$\frac{\partial \ln L}{\partial \beta} = \sum_{i=0}^n \left[ y_i \frac{f(x'_i\beta)}{F(x'_i\beta)} + (1 - y_i) \frac{-f(x'_i\beta)}{1 - F(x'_i\beta)} \right] x_i = 0$$

- Mientras que no se esté en el modelo lineal, la ecuación será no lineal y serán necesarios métodos iterativos
- Para el modelo *logit* se obtiene la siguiente condición de primer orden:

$$\frac{\partial \ln L}{\partial \beta} = \sum_{i=0}^n [y_i - \Lambda_i] x_i = 0$$

- Si  $x_i$  contiene un término constante, entonces la media muestral de las probabilidades predichas debe ser igual a la proporción de unos en la muestra. Además, esta implicación tiene similitudes con las ecuaciones normales de mínimos cuadrados si se interpreta  $y_i - \Lambda_i$  como el error
- Para el modelo *probit* se obtiene la siguiente condición de primer orden:

$$\frac{\partial \ln L}{\partial \beta} = \sum_{y_i=0} \frac{-\phi_i}{1 - \Phi_i} x_i + \sum_{y_i=1} \frac{\phi_i}{\Phi_i} x_i = \sum_{y_i=0} \lambda_{0i} x_i + \sum_{y_i=1} \lambda_{1i} x_i = 0$$

- Se han sugerido varias medidas para modelos de respuesta discreta que intentan imitar al  $R^2$ , ya que mide el éxito de un estimador al optimizar el criterio de ajuste y, si se calcula como  $\rho^2$ , mide cómo las predicciones del modelo pueden imitar los datos reales
  - Debido a que muchos modelos se estiman por máxima verosimilitud, de modo que el criterio de ajuste es la verosimilitud
    - Siguiendo con el criterio del  $R^2$ , es interesante ver la hipótesis de que todas las pendientes sean nulas. La log-verosimilitud en el caso de que solo haya un término constante sería:

$$\ln L_0 = n[P_0 \ln P_0 + P_1 \ln P_1] \quad \text{where } P_j \text{ is proportion}$$

- El indicador más usado es la pseudo- $R^2$  es una medida que está acotada entre cero y uno, pero que no puede llegar nunca a uno (la verosimilitud nunca puede ser nula)

$$R_{pseudo}^2 = 1 - \frac{\ln L_{MLE}}{\ln L_0}$$

- Es posible hacer una corrección de grados de libertad utilizando el siguiente estadístico:

$$R_{pseudo}^2 = 1 - \frac{\ln L_{MLE} - K}{\ln L_0}$$

- El pseudo- $R^2$  es muy útil para comparar un modelo con otro, pero si los modelos son anidados, entonces la función de log-verosimilitud es la elección natural. Para casos más generales, los investigadores utilizan criterios de información

- Los dos criterios de información más utilizados son el criterio de información de Akaike y el criterio de información de Bayes

$$AIC = -2 \ln L + 2K \quad \text{or} \quad AIC/n$$

$$BIC = -2 \ln L + K \ln n \quad \text{or} \quad BIC/n$$

- En general, cuanto menor sea el valor del criterio mejor será el ajuste

## Los modelos de elección discreta multinomial desordenada

- Cuando la variable dependiente toma valores discretos, pero no es binaria (no toma solo dos valores, sino que toma tres o más), el modelo lineal sigue siendo un modelo inadecuado para estudiar la elección discreta
  - Los modelos anteriores se pueden generalizar con tal estudiar las probabilidades de que un individuo seleccione una alternativa  $j$  respecto a todas las alternativas posibles  $J$ , de modo que los números de la variable representan diferentes alternativas
    - Se suele utilizar el modelo *logit* multinomial y el modelo *logit* condicional

- Los coeficientes de los regresores  $\mathbf{x}$  siguen siendo no interpretables, de modo que se sigue utilizando el concepto de efecto marginal
- Se sigue usando el modelo de utilidad aleatoria o *random utility model*, pero generalizado al caso en el que hay más de dos alternativas. En este caso, los coeficientes que expresan la diferencia entre un grupo base y otro grupo (el vector  $\beta$ ) variarán dependiendo de la alternativa (a diferencia del caso anterior):

$$U_{ij} = \mathbf{x}'_i \boldsymbol{\gamma}_j + u_{ij} \text{ for } j = 1, 2, \dots, J \Rightarrow U_{ij} - U_{ig} = \mathbf{x}'_i \boldsymbol{\beta}_j + \varepsilon_{ij}$$

$$\text{where } \boldsymbol{\beta}_j = \boldsymbol{\gamma}_j - \boldsymbol{\gamma}_g \text{ for } j = 0, 1, 2, \dots, J$$

- Cuando la variable dependiente es equivalente a una alternativa  $j$ , quiere decir que, según el modelo, la utilidad de la alternativa  $j$  es mayor a la de la alternativa  $g$  (se sigue el modelo de maximización de utilidad del consumidor)

$$U_{ij} = \max(U_{i1}, U_{i2}, \dots)$$

$$Y_i = j \text{ if } U_{ij} > U_{ig} \text{ when } j \neq g \text{ for } j = 0, 1, 2, \dots, J$$

- En este caso la probabilidad de que se escoja la alternativa  $j$  (que  $Y_i = j$ ) se generaliza al caso en el que el grupo base y el grupo de la alternativa no son iguales:

$$P(Y_i = j | \mathbf{x}_i) = P(U_{ij} > U_{ig} | \mathbf{x}_i) = P(U_{ij} - U_{ig} > 0 | \mathbf{x}_i)$$

$$= P(\varepsilon_{ij} > \mathbf{x}'_i \boldsymbol{\beta}_j + \varepsilon_{ig} | \mathbf{x}_i) \text{ for all } j \neq g$$

- Se asume que el término de error para cada alternativa es idéntico e independiente del término de las otras alternativas, y que para los casos multinomiales sigue una distribución de Gumbel o *extreme value type I*

$$\mathbf{u}_i \text{ is i.i.d. } \Rightarrow \boldsymbol{\varepsilon}_i \text{ is i.i.d.}$$

$$F(\varepsilon_{ij}) = \exp\left(-\frac{1}{e^{\varepsilon_{ij}}}\right) \Rightarrow u_{ij}, \varepsilon_{ij} \sim \text{Gumbel}$$

$$\text{where } \mathbf{u}_i = \begin{pmatrix} u_{i1} \\ u_{i2} \\ \dots \end{pmatrix} \text{ and } \boldsymbol{\varepsilon}_i = \begin{pmatrix} \varepsilon_{i1} \\ \varepsilon_{i2} \\ \dots \end{pmatrix}$$

- Si se cree que el error se distribuye de manera normal (con media cero y varianza constante), en vez de seguir una

distribución *Extreme Value Type-I*, se utiliza el modelo *probit* multinomial

- El modelo *logit* multinomial es un modelo generalizado del modelo *logit*. Es un modelo de regresión no lineal que utiliza la función de distribución de probabilidad acumulada de la distribución logística, con la que, a partir del *additive random utility model* se demuestra que la probabilidad de que  $Y = j$  (de entre  $J$  alternativas mutuamente excluyentes) es:

$$P(Y_i = j | \mathbf{x}_i) = \frac{e^{(\mathbf{x}_i' \boldsymbol{\beta}_j)}}{\sum_{m=1}^J e^{(\mathbf{x}_i' \boldsymbol{\beta}_m)}}$$

- En el modelo, igual que en el caso binario, se escoge una alternativa como grupo o categoría base y, por tanto, los coeficientes del modelo pueden ser interpretados respecto a la base como la diferencia entre escoger la alternativa  $j$  o la base
  - Esto implica que los coeficientes de la categoría base serán todos nulos porque si no habría un problema de identificación (demasiados coeficientes)

$$\boldsymbol{\beta}_g = 0$$

- Si los coeficientes de la alternativa base son nulos (por suposición), la probabilidad de escoger la alternativa base serían:

$$P(Y_i = g | \mathbf{x}_i) = \frac{1}{\sum_{m=1}^J e^{(\mathbf{x}_i' \boldsymbol{\beta}_m)}}$$

- Hay dos propiedades que se cumplen en los modelos *logit* si se incluye el coeficiente constante:
  - Como las alternativas tienen que ser mutuamente excluyentes, la suma de las probabilidades de escoger cada alternativa para cada individuo debe ser igual a la unidad (con certeza escogerá una alternativa). Esto implica que el cambio en la probabilidad de una alternativa hace que cambie la probabilidad de las otras

$$\sum_{j=1}^J P(Y_i = j | \mathbf{x}_i) = 1 \quad \text{for } i = 1, 2, \dots, N$$

$$P(Y_i = j | \mathbf{x}_i) \uparrow \Rightarrow P(Y_i = m | \mathbf{x}_i) \downarrow \quad \text{when } j \neq m$$

- Igual que en el caso binario, si se incluye el coeficiente constante, la suma de probabilidades ajustadas del modelo *logit* coincide con el número de individuos con valor  $Y = j$



$$\sum_{i=1}^N P(Y_i = j | \mathbf{x}_i) = N_j \text{ for } j = 1, 2, \dots, J$$

- Si se no se incluye el coeficiente constante, estas propiedades dejan de ser ciertas
- Este modelo se puede modelar a través de variables binarias (respuesta) para cada una de las  $J$  categorías

$$D_{ij} = \begin{cases} 1 & \text{if } Y_i = j \\ 0 & \text{if } Y_i \neq j \end{cases} \text{ for } j = 0, 1, \dots, J-1$$

- Se crean  $J - 1$  variables binarias porque la categoría base no se modelaría
- De este modo, en vez de tener un solo modelo, se tienen  $J$  modelos de regresión binaria *logit*, dado que se modela por separado y no se cuenta el orden

$$p_{ij} = P(D_{ij} = 1 | \mathbf{x}_i) = P(Y_i = j | \mathbf{x}_i)$$

- Igual que en los modelos *logit* y *probit* binomiales, en el modelo *logit* multinomial no se pueden interpretar directamente los coeficientes, de modo que se utiliza el efecto marginal
  - En este caso, los efectos marginales dependen no solo del individuo y del regresor, si no también de la alternativa, debido a que cada alternativa tiene un coeficiente. Además, también depende de las probabilidades de las otras alternativas:

$$ME_{ikj} = \frac{\partial P(Y_i = j | \mathbf{x}_i)}{\partial X_{ik}}$$

$$\frac{\partial P(Y_i = j | \mathbf{x}_i)}{\partial X_{ikj}} = P(Y_i = j | \mathbf{x}_i) \left[ \beta_{jk} - \sum_{m=1}^J \beta_{mk} P(Y_i = m | \mathbf{x}_i) \right]$$

- La sumatoria de las probabilidades de cada alternativa multiplicadas por el coeficiente de un regresor  $k$  (para cada alternativa) se puede interpretar como una media ponderada de los coeficientes del regresor  $k$  a través de las alternativas. De este modo, cada coeficiente de  $X_{ik}$  entra en el efecto marginal de cada alternativa

$$\bar{\beta}_k = \sum_{m=1}^J \beta_{mk} P(Y_i = m | \mathbf{x}_i)$$

- Además, la diferencia entre el coeficiente y la media ponderada no tiene por qué tener el mismo signo que el coeficiente  $\beta_{jk}$ , por lo que no se puede interpretar los signos de los coeficientes (como sí que ocurría en el modelo *logit* binomial)
- Como  $\beta_{jk}$  es uno de los coeficientes que se tienen en cuenta también al hacer la media ponderada  $\bar{\beta}_k$ , solo se sabe con certeza que el coeficiente con el valor máximo es mayor a la media y que aquel con valor mínimo es menor:

$$\beta_{ak} < \sum_{m=1}^J P(Y_i = m | \mathbf{x}_i) \beta_{mk} < \beta_{bk}$$

$a = \text{alternative with minimum } \beta \text{ value}$

$b = \text{alternative with maximum } \beta \text{ value}$

- Para la alternativa con el coeficiente con un valor máximo o mínimo, por tanto, se puede saber el signo del efecto marginal, dado que se sabe el signo de la diferencia

$$\beta_{bk} - \sum_{m=1}^J P(Y_i = m | \mathbf{x}_i) \beta_{mk} > 0 \Rightarrow ME_{ikb} > 0$$

$$\beta_{ak} - \sum_{m=1}^J P(Y_i = m | \mathbf{x}_i) \beta_{mk} < 0 \Rightarrow ME_{ika} < 0$$

- Los coeficientes en el modelo *logit* multinomial también se estiman mediante la función de máxima verosimilitud, la cual se generaliza en este modelo:

$$L = \prod_{i=1}^{N_1} P(Y_i = 1 | \mathbf{x}_i) \dots \prod_{i=N_{J-1}+1}^{N_J} P(Y_i = J | \mathbf{x}_i) = \prod_{j=1}^J \prod_{i=N_{j-1}+1}^{N_j} P(Y_i = j | \mathbf{x}_i)$$

$$\max_{\beta_j} \ln \left[ \prod_{j=1}^J \prod_{i=N_{j-1}+1}^{N_j} P(Y_i = j | \mathbf{x}_i) \right] = \sum_{j=1}^J \sum_{i=N_{j-1}+1}^{N_j} \ln \left[ \frac{e^{(\mathbf{x}'_i \beta_j)}}{\sum_{r=1}^J e_j^{(\mathbf{x}'_i \beta_j)}} \right]$$

- Existen dos casos especiales que simplifican la obtención de los coeficientes maximizadores:

- Cuando solo se tiene que estimar el coeficiente constante es como si se tuviera que estimar la media, y el coeficiente constante que maximiza la función de verosimilitud será aquel que haga que la probabilidad sea la proporción de observaciones en la muestra donde  $Y = j$ :

$$p_j^* = P^*(Y_i = j | x_{ij}) = \frac{e^{\beta_{0j}}}{\sum_{m=1}^J e^{\beta_{0m}}} = \frac{N_j}{N} = \frac{\sum y_i}{N} = \bar{Y}_j$$

$$L = \left(\frac{N_1}{N}\right)^{N_1} \dots \left(\frac{N_J}{N}\right)^{N_J} = \prod_{j=1}^J \left(\frac{N_j}{N}\right)^{N_j}$$

- Cuando se incluye un solo regresor binario, los individuos se separan en dos grupos, lo cual hace que los coeficientes maximizadores sean aquellos que hagan que la probabilidad de que  $Y = 1$  sea igual a la proporción de individuos en la muestra con ese valor para cada grupo:

$$P(Y_i = j | X_i = 1) = \frac{e^{(\beta_{0j} + \beta_{1j})}}{\sum_{m=1}^J e^{(\beta_{0m} + \beta_{1m})}} \Rightarrow P^*(Y_i = j | X_i = 1) = \frac{N_{j,X=1}}{N_{X=1}}$$

$$P(Y = j | X_i = 0) = \frac{e^{\beta_{0j}}}{\sum_{m=1}^J e^{\beta_{0m}}} \Rightarrow P^*(Y = j | X_i = 0) = \frac{N_{j,X=0}}{N_{X=0}}$$

$$L = \left(\frac{N_{1,X=1}}{N_{X=1}}\right)^{N_{1,X=1}} \left(\frac{N_{2,X=1}}{N_{X=1}}\right)^{N_{2,X=1}} \dots \left(\frac{N_{1,X=0}}{N_{X=0}}\right)^{N_{1,X=0}} \left(\frac{N_{2,X=0}}{N_{X=0}}\right)^{N_{2,X=0}} \dots$$

$$L = \prod_{j=1}^J \left(\frac{N_{j,X=1}}{N_{X=1}}\right)^{N_{j,X=1}} \prod_{j=1}^J \left(\frac{N_{j,X=0}}{N_{X=0}}\right)^{N_{j,X=0}}$$

- Este último caso particular puede generalizarse para varias variables binarias
- Cuando se produce un cambio de categoría o grupo base y se tiene el valor de los coeficientes estimados, se pueden estimar los nuevos coeficientes porque se mantienen las diferencias entre coeficientes:

$$\beta_{1j} - \beta_{2j} = \beta'_{1j} - \beta'_{2j}$$

- Esto ocurre porque, cuando se cambia de base, las probabilidades de un individuo  $i$  para las  $J$  alternativas no tiene que variar. En consecuencia, tampoco varía la *odds ratio* (la probabilidad relativa de una alternativa respecto a otra para un individuo  $i$ )

- Tampoco cambian los efectos marginales, debido a que no cambian las probabilidades
- Otro modelo que se utiliza es el *logit* condicional, el cual es una ampliación del modelo *logit* multinomial. En este modelo, se continúan teniendo  $J$  alternativas, pero ahora hay dos tipos de regresores: unos que no dependen de la alternativa (las  $X$ , pero varían según el individuo  $i$ ) y otros que sí lo hacen (por lo que estos varían según la alternativa  $j$  para un mismo individuo  $i$ )

$$U_{ij} = x'_i \gamma_j + z'_{ij} \alpha + u_{ij}$$

$k = \text{number of regressor independent of } j$

- En este modelo, los coeficientes de los regresores que varían según la alternativa son constantes, de modo que en la diferencia de utilidades solo se tiene en cuenta la diferencia del regresor entre alternativas:

$$U_{ij} - U_{ig} = x'_i (\gamma_j - \gamma_g) + (z'_{ij} - z'_{ig}) \alpha + (u_{ij} - u_{ig})$$

$$U_{ij} - U_{ig} = x'_i \beta_j + (z'_{ij} - z'_{ig}) \alpha + \varepsilon_{ij}$$

- Como los coeficientes  $\alpha$  siempre son los mismos, estos no necesitan una alternativa base para ser calculados (si no se tuvieran coeficientes constantes habría un problema de mala especificación). Sin embargo, aún se necesita seleccionar una alternativa base para la estimación de los coeficientes  $\beta$
- Los errores siguen teniendo la misma distribución asumida anteriormente
- De este modo, la probabilidad de escoger la probabilidad  $j$  será:

$$P(Y_i = j | x_{ij}, z_i) = \frac{e^{(x'_i \beta_j + \alpha z'_{ij})}}{\sum_{m=1}^J e^{(x'_i \beta_m + \alpha z'_{im})}}$$

$g = \text{reference category}$

- En este caso, la diferencia entre  $z$  para cada alternativa no se cuenta, dado que cuando se estima el modelo,  $\alpha$  representa el cambio en  $U_{ij} - U_{ig}$  cuando  $z_{ij}$  varía una unidad
- En el modelo *logit* condicional, los coeficientes de los regresores invariantes con la alternativa (las  $\beta$ ) se interpretan del mismo modo que en el modelo *logit* multinomial, pero los coeficientes de los regresores que varían según la alternativa (las  $\alpha$ ) se interpretan como en el modelo *logit* binomial

- Como los coeficientes  $\beta$  dependen de la alternativa  $j$ , solo se puede interpretar aquellos que tengan el valor máximo y el mínimo (a través de sus efectos marginales)

$$\frac{\partial P(Y_i = j | \mathbf{x}_i, \mathbf{z}_{ij})}{\partial X_{ik}} = P(Y_i = j | \mathbf{x}_i, \mathbf{z}_{ij})(\beta_{jk} - \bar{\beta}_k)$$

- Como los coeficientes  $\alpha$  no varían con la alternativa (siempre son los mismos para cualquier alternativa  $j$ ), el efecto marginal de variar una unidad de un regresor  $Z_{ij}$  tendrá el mismo signo que su coeficiente  $\alpha$ , pero la magnitud sigue siendo no interpretable

$$\frac{\partial P(Y_i = j | \mathbf{x}_i, \mathbf{z}_{ij})}{\partial Z_{ikj}} = [P(Y_i = j | \mathbf{x}_i, \mathbf{z}_{ij})[1(j = m) - P(Y_i = m | \mathbf{x}_i, \mathbf{z}_{ij})]] \alpha_k$$

$$\text{where } 1(j = m) = \begin{cases} 1 & \text{if } j = m \\ 0 & \text{if } j \neq m \end{cases}$$

- Cuando incrementa una unidad del regresor  $Z_{ij}$ , el efecto en la probabilidad del individuo  $i$  de escoger la alternativa  $j$  tendrá el mismo signo. Sin embargo, el efecto en la probabilidad de que el individuo  $i$  escoja las otras alternativas tendrá un signo contrario, debido a que el individuo  $i$  tiene que escoger una alternativa de las posibles

$$\sum_{m=1}^J P(Y_i = j | \mathbf{x}_i, \mathbf{z}_{im}) = 1 \Rightarrow \begin{cases} P(Y_i = 1 | \mathbf{x}_i, \mathbf{z}_{ij}) \uparrow \\ P(Y_i = 2 | \mathbf{x}_i, \mathbf{z}_{ij}) \downarrow \\ \dots \downarrow \end{cases}$$

- A diferencia de con variaciones en  $\mathbf{X}_i$ , en donde el efecto del cambio depende de la alternativa  $j$  en la que se esté, el efecto de  $\mathbf{Z}_{ij}$  es el mismo para cualquier alternativa, de modo que el efecto de la variación en los regresores  $\mathbf{Z}_{ij}$  para  $j$  es opuesto para las otras alternativas  $m \neq j$
- Los modelos *logit* multinomiales y *logit* condicionales cumplen una propiedad conveniente para la estimación, pero que no siempre se adecua a la realidad: la independencia de las alternativas irrelevantes o *independence of irrelevant alternatives* (IIA)
  - La *odds ratio* es la probabilidad relativa entre una alternativa  $j$  y una alternativa base  $g$ , la cual se define como la *ratio* entre ambas
    - Para el caso de un modelo *logit* multinomial, se puede ver como la *odds ratio* solo depende de  $j$  y  $r$ , por lo que las otras alternativas son irrelevantes y los coeficientes estimados para

las otras alternativas no tendrían que variar nada o por lo menos no mucho si se extraen alternativas

$$\frac{P(Y_i = j | \mathbf{x}_i)}{P(Y_i = r | \mathbf{x}_i)} = \frac{\frac{e^{(x'_i \beta_j)}}{\sum_{m=1}^J e^{(x'_i \beta_m)}}}{\frac{e^{(x'_i \beta_r)}}{\sum_{m=1}^J e^{(x'_i \beta_m)}}} = e^{x'_i(\beta_j - \beta_r)} = e^{x'_i(\beta_j - \beta_r)}$$

- Para el caso de un modelo *logit* multinomial, se puede ver como la *odds ratio* solo depende de  $j$  y  $r$ , por lo que las otras alternativas son irrelevantes y los coeficientes estimados para las otras alternativas no tendrían que variar nada o por lo menos no mucho si se extraen alternativas

$$\begin{aligned} \frac{P(Y_i = j | \mathbf{x}_i, \mathbf{z}_{ij})}{P(Y_i = r | \mathbf{x}_i, \mathbf{z}_{ir})} &= \frac{\frac{e^{[x'_i \beta_j + z'_{ij} \alpha]}}{\sum_{m=1}^J e^{[x'_i \beta_m + z'_{im} \alpha]}}}{\frac{e^{[x'_i \beta_r + z'_{ir} \alpha]}}{\sum_{m=1}^J e^{[x'_i \beta_m + z'_{im} \alpha]}}} = \\ &= \frac{e^{[x'_i \beta_j + z'_{ij} \alpha]}}{e^{[x'_i \beta_r + z'_{ir} \alpha]}} = \frac{e^{(x'_i \beta_j)} e^{(z'_{ij} \alpha)}}{e^{(x'_i \beta_r)} e^{(z'_{ir} \alpha)}} = e^{x'_i(\beta_j - \beta_r) + (z'_{ij} - z'_{ir}) \alpha} \end{aligned}$$

- En ambos casos, la variación en la *odds ratio* por el incremento de una unidad de  $\mathbf{X}$  es la exponencial de la diferencia de los coeficientes  $\beta$  de cada alternativa. En el caso del *logit* condicional, un aumento en  $\mathbf{Z}$  provoca un incremento de  $e^{(z'_{ij} - z'_{ir}) \alpha}$  en la *odds ratio*

$$\Delta X_{ij} = 1 \Rightarrow \frac{P(Y_i = j | \mathbf{x}_i)}{P(Y_i = r | \mathbf{x}_i)} * e^{x'_i(\beta_j - \beta_r)} = e^{(x'_i + 1)(\beta_j - \beta_r)}$$

$$\begin{aligned} \Delta X_i = 1 &\Rightarrow \frac{P(Y_i = j | \mathbf{x}_i, \mathbf{z}_{ij})}{P(Y_i = r | \mathbf{x}_i, \mathbf{z}_{ir})} * e^{x'_i(\beta_j - \beta_r)} \\ &= e^{(x'_i + 1)(\beta_j - \beta_r) + (z'_{ij} - z'_{ir}) \alpha} \end{aligned}$$

$$\begin{aligned} \Delta Z_{ij} = 1 &\Rightarrow \frac{P(Y_i = j | \mathbf{x}_i, \mathbf{z}_{ij})}{P(Y_i = r | \mathbf{x}_i, \mathbf{z}_{ir})} * e^{(z'_{ij} - z'_{ir}) \alpha} \\ &= e^{x'_i(\beta_j - \beta_r) + 2(z'_{ij} - z'_{ir}) \alpha} \end{aligned}$$

- Para saber los cambios probabilísticos que se dan cuando se incluye o se excluye una alternativa en un modelo *logit* multinomial, se resuelve el siguiente sistema:

$$\begin{cases} \sum_{j=1}^J P(Y_i = j | x_{ij}) = 1 \\ \frac{P(Y_i = j | x_{ij})}{P(Y_i = m | x_{ij})} = a \text{ when } j \neq m \\ P(Y_i = j | x_{ij}) \text{ for } j = 1, 2, \dots, J \end{cases}$$

- Un inconveniente de esta suposición es que a veces da resultados contrarios al sentido común, debido a que la *odds ratio* de las probabilidades no puede variar y hay casos en los que la inclusión o exclusión de una alternativa hace que varíe la probabilidad de las alternativas
- La IIA nace de la suposición inicial de que los errores son independientes y homoscedásticos. Por lo tanto, para comprobar que se cumple la suposición de IIA, se puede usar el test de Hausman o *Hausman test* y así comprobar la consistencia de los estimadores

$$\begin{cases} H_0: \text{the models are equal} \\ H_1: \text{the models are not equal} \end{cases}$$

$$h = (\beta_s - \beta_f)' [V_s - V_f]^{-1} (\beta_s - \beta_f) \sim \chi_r^2$$

$s = \text{estim. on the restr. subset}$      $f = \text{estim. on the full set}$

$V = \text{asympt. cov. matrix}$      $r = n^\circ \text{ of elim. coefficients}$

- El test se basa en estimar dos modelos: uno completo y otro eliminando las alternativas que se sospecha que no cumplen la IIA. De este modo, la diferencia no está en las restricciones que se imponen en coeficientes (como en otros contrastes de hipótesis), si no en el tamaño de la muestra con el que se trabaja (se eliminan las observaciones para las alternativas sospechosas)
- Si se cumple, los modelos tienen que ser muy parecidos o iguales, dado que la eliminación de las alternativas irrelevantes no tiene que afectar a la probabilidad de escoger entre un par de alternativas (la probabilidad relativa de escoger una alternativa  $j$ ) y los estimadores serán consistentes. La hipótesis nula del test expresa que los dos modelos son iguales o muy parecidos (se cumple la IIA), mientras que la hipótesis alternativa expresa que no
- Los coeficientes que se eliminan serán el número de regresores más una unidad ( $k + 1$ ) por cada alternativa que se elimina,

debido a que los coeficientes para cada alternativa son diferentes (pero no varía el número de regresores)

$$r = (k + 1) * n^{\circ} \text{ of eliminated altern.} + d$$

$$d = \text{number of } \mathbf{z} \text{ regressors in remaining altern.}$$

- Las medidas de ajuste y los tests de hipótesis utilizadas en el modelo *logit* binomial pueden utilizarse perfectamente para los modelos *logit* multinomial y para los modelos *logit* condicional, pero las restricciones se ven afectadas
  - En la *likelihood ratio test*, por ejemplo, dado que habrá un número  $k$  de coeficientes de regresores para cada alternativa, los grados de libertad serán el producto entre el número de alternativas y el número de regresores

$$q = k * n^{\circ} \text{ of alternatives} \rightarrow \text{Multinomial Logit}$$

$$q = k * n^{\circ} \text{ of alternatives} + d^{tot} \rightarrow \text{Conditional Logit}$$

$$d^{tot} = \text{number of } \mathbf{z} \text{ regressors}$$

- Si no se cumple la condición de IIA, se necesita utilizar una alternativa al modelo *logit* multinomial o condicional. Una manera de relajar la condición es el modelo *logit* anidado, de modo que se agrupan alternativas similares en subgrupos para que se mantenga la IIA en cada uno
  - Se supone que hay  $J$  alternativas y que se pueden dividir en  $B$  subgrupos o ramas por similitud, de modo que el conjunto de elección se puede dividir de la siguiente manera:

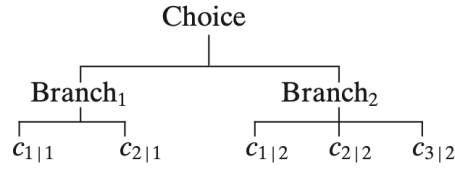
$$\{c_1, c_2, \dots, c_J\} = \left\{ \{c_{1|1}, \dots, c_{J_1|1}\}, \{c_{1|2}, \dots, c_{J_2|2}\}, \dots, \{c_{1|B}, \dots, c_{J_B|B}\} \right\}$$

- Cada subgrupo tiene un número de alternativas similares  $J_b$  que no son las mismas para cada grupo, por lo que el subíndice  $j|b$  indica el número de la alternativa dentro del subgrupo, pero no indica el número de la alternativa dentro del conjunto original de alternativas

$$c_{j|b} \neq c_j$$

- Se puede entender que el proceso de elección consiste en escoger un subconjunto de los  $B$  posibles y entonces seleccionar una alternativa dentro de cada subgrupo





- Suponiendo que existen observaciones para los atributos de las alternativas  $X_{ij|b}$  y para atributos de los subconjuntos de elección  $Z_{ib}$ , el modelo econométrico permite obtener la siguiente probabilidad incondicional para la alternativa  $j_b$  del subgrupo  $b$

$$P(j_b, b | x_{ij|b}, z_{ib}) = \frac{e^{(x'_{ij|b}\beta + z'_{ib}\alpha)}}{\sum_{l=1}^B \sum_{r=1}^{J_l} e^{(x'_{ir|l}\beta + z'_{il}\alpha)}}$$

- En este caso, se asume que los parámetros de los coeficientes son constantes, de modo que no dependen de las alternativas  $j$  ni de los subgrupos  $b$ , por lo que es parecido a un modelo *logit* condicional

$$U_{ij} = x'_{ij|b}\beta + z'_{ib}\alpha + \varepsilon_{ib}$$

- El término de error en este modelo sigue una distribución de valor extremo generalizado o *generalized extreme value* (GEV)

$$\varepsilon_{ib} \sim F(\varepsilon_{i1}, \varepsilon_{i2}, \dots, \varepsilon_{iB}) \quad \text{where}$$

$$F(\varepsilon_{i1}, \varepsilon_{i2}, \dots, \varepsilon_{iB}) = \exp \left[ \sum_{l=1}^B \left( \sum_{r=1}^{J_l} e^{-\frac{\varepsilon_{kr}}{\tau_l}} \right)^{\tau_l} \right] \quad \text{for all } k \in B$$

- Esta probabilidad se puede descomponer en la probabilidad de escoger la alternativa  $j_b$  y la probabilidad de escoger el subgrupo  $b$

$$\begin{aligned} P(j_b, b | x_{ij|b}, z_{ib}) &= \frac{e^{(x'_{ij|b}\beta + z'_{ib}\alpha)}}{\sum_{l=1}^B \sum_{r=1}^{J_l} e^{(x'_{ir|l}\beta + z'_{il}\alpha)}} = \\ &= \frac{e^{(x'_{ij|b}\beta)} e^{(z'_{ib}\alpha)}}{\sum_{l=1}^B \sum_{r=1}^{J_l} e^{(x'_{ir|l}\beta)} e^{(z'_{il}\alpha)}} = \frac{e^{(x'_{ij|b}\beta)}}{\sum_{r=1}^{J_b} e^{(x'_{ir|b}\beta)}} \frac{e^{(z'_{ib}\alpha)}}{\sum_{l=1}^B e^{(z'_{il}\alpha)}} \frac{\left[ \sum_{r=1}^{J_b} e^{(x'_{ir|b}\beta)} \right] \left[ \sum_{l=1}^B e^{(z'_{il}\alpha)} \right]}{\sum_{l=1}^B \sum_{r=1}^{J_b} e^{(x'_{ir|l}\beta)} e^{(z'_{il}\alpha)}} \\ &= P[Y_i = j_b | x_{ij|b}, z_{ib}, b] P[b | x_{ij|b}, z_{ib}] \end{aligned}$$

$$\text{where } \begin{cases} P(Y_i = j_b | x_{ij|b}, z_{ib}, b) = \frac{e^{(x'_{ij|b}\beta)}}{\sum_{r=1}^{J_b} e^{(x'_{ir|b}\beta)}} \\ P(b | x_{ij|b}, z_{ib}) = \frac{e^{(z'_{ib}\alpha)} \left[ \sum_{r=1}^{J_b} e^{(x'_{ir|b}\beta)} \right] \left[ \sum_{l=1}^B e^{(z'_{il}\alpha)} \right]}{\sum_{l=1}^B e^{(z'_{il}\alpha)} \sum_{l=1}^B \sum_{r=1}^{J_b} e^{(x'_{ir|l}\beta)} e^{(z'_{il}\alpha)}} \end{cases}$$

- La variable  $IV_{ib}$  se denomina valor inclusivo o *inclusive value*, y estos dependen del individuo y del subgrupo  $b$ , pero no de la alternativa (por lo que se tiene que calcular para cada alternativa)

$$IV_{ib} = \ln \left( \sum_{r=1}^{J_b} e^{(x'_{ir|b}\beta)} \right)$$

$$\Rightarrow P(j_b, b | x_{ij|b}, z_{ib}) = \frac{e^{(x'_{ij|b}\beta)}}{\sum_{r=1}^{J_b} e^{(x'_{ir|b}\beta)}} \frac{e^{[\tau_b(z'_{ib}\alpha + IV_{ib})]}}{\sum_{l=1}^B e^{[\tau_l(z'_{il}\alpha + IV_{il})]}}$$

$$\Rightarrow \begin{cases} P(Y_i = j_b | x_{ij|b}, z_{ib}, b) = \frac{e^{(x'_{ij|b}\beta)}}{\sum_{r=1}^{J_b} e^{(x'_{ir|b}\beta)}} \\ P(b | x_{ij|b}, z_{ib}) = \frac{e^{[\tau_b(z'_{ib}\alpha + IV_{ib})]}}{\sum_{l=1}^B e^{[\tau_l(z'_{il}\alpha + IV_{il})]}} \end{cases}$$

- Para que el modelo derivado con el valor inclusivo sea igual al modelo para la probabilidad incondicional, los valores  $\tau$  deben ser iguales a 1. Por lo tanto, se utiliza una restricción  $\tau = 1$  para que el modelo sea equivalente al de un *logit* condicional
- No obstante, si los valores  $\tau$  no son unitarios, entonces nace un modelo diferente, que es el modelo *logit* anidado. Estos valores  $\tau$  son el grado de disimilitud, siendo una medida inversa de la correlación entre los errores de diferentes alternativas en un mismo grupo

$$\tau_b = \sqrt{1 - \rho_{j,k}} \quad \text{where } \rho_{j,k} = \text{corr}(\varepsilon_{bj}, \varepsilon_{bk}) \text{ and } j \neq k$$

- La correlación entre los errores de un mismo grupo tiene que tener correlación positiva, dado que las alternativas dentro de cada grupo son similares. Por lo tanto,  $\tau_b$  tiene que estar necesariamente entre 0 y 1

- Para grupos que tengan una sola alternativa, la correlación no estará definida (porque no hay otras alternativas) y por tanto se asume que  $\tau = 1$  para estos grupos
- Cuando hay un subgrupo formado por subgrupos, se utilizan varios valores  $\tau_{hb}$ , donde  $h$  indica el subgrupo superior (formado por subgrupos inferiores) y  $b$  indica el subgrupo inferior
- Debido a que la IIA se debe a la independencia y homoscedasticidad de los errores, al agrupar alternativas en subgrupos se permite que la varianza difiera entre cada subgrupo pero que se mantenga la homoscedasticidad dentro de cada uno
  - Si los valores  $\tau$  difieren de 1 (se relaja la restricción), el modelo es un *logit* anidado. Los parámetros  $\tau$  permiten que exista heterogeneidad entre grupos, pero homoscedasticidad dentro de cada uno, dado que la varianza dentro de cada grupo se define con la siguiente expresión:

$$\sigma_b^2 = \frac{\pi^2}{6\tau_b}$$

- Como los coeficientes no son interpretables, se utilizan los efectos marginales. Comúnmente, se utiliza la semielasticidad para poder entender el efecto de un aumento de un incremento unitario en las variables  $X_{ij|b}$

$$\frac{\partial \ln[P(j_b, b|x_{ij|b}, z_{ib})]}{\partial x_{ij|b}} =$$

$$= \{1(l = b)[1(r = j_b) - P_{j_b|b}] + \tau_b[1(l = b) - P_b]P_{j_b|b}\}\beta$$

$$\text{where } P_{j_b|b} = P(Y_i = j_b|x_{ij|b}, z_{ib}, b) \text{ and } P_b = P(b|x_{ij|b}, z_{ib})$$

- Si el grupo no es el mismo, entonces el efecto tiene el signo opuesto a  $\beta$ , pero
- Para poder estimar los coeficientes de un modelo *logit* anidado se pueden estimar de dos maneras diferentes:
  - Un enfoque de máxima verosimilitud de dos pasos (llamado enfoque de información limitada) que consiste en estimar las  $\beta$  tratando la elección dentro de cada subgrupo como un modelo *logit* condicional, y después calcular los valores intrínsecos para cada subgrupo y estimar las  $\alpha$  y  $\tau$  tratando la elección entre grupos como un modelo *logit* condicional con atributos  $z_{ib}$  e  $IV_{ib}$

- Un enfoque de información completa de máxima verosimilitud, en donde se maximiza la función logarítmica de máxima verosimilitud

$$\ln L = \sum_{i=1}^n \ln[P(Y_i = j_b | \mathbf{x}_{ij|b}, \mathbf{z}_{ib}, b)P(b | \mathbf{x}_{ij|b}, \mathbf{z}_{ib})]$$

- El enfoque de información limitada es menos eficiente debido a que la estimación por dos pasos hace que la matriz de estimaciones se tenga que corregir en el segundo paso
- Para poder contrastar hipótesis sobre si el modelo tiene que ser anidado o no, se tienen que hacer tests de hipótesis para la heteroscedasticidad, la cual se puede comprobar a través de tres estadísticos diferentes
  - El test de hipótesis de Wald es un test de heteroscedasticidad que utiliza
- Otra variante del modelo de elección multinomial es el modelo *logit* de parámetros aleatorios o *random parameters logit model*, también llamado modelo mixto
  - En una posible especificación, se considera que los coeficientes son aleatorios (por lo que tienen una formulación propia) y que siguen una distribución normal, de modo que hay un modelo para la utilidad y otro modelo para los coeficientes

$$U_{ij} = \mathbf{z}'_{ij}\boldsymbol{\beta}_i + \varepsilon_{ij} \text{ where } \beta_i = \beta + u_i$$

$$u_i \sim N(0, \sigma_u^2)$$

- Dado que los coeficientes son aleatorios y dependen de un término de error distribuido normalmente, los coeficientes aleatorios también seguirán una distribución normal con media  $\beta$  y matriz de covarianzas  $\boldsymbol{\Sigma}_\beta$

$$\beta_i \sim N(\beta, \boldsymbol{\Sigma}_\beta)$$

- Otra especificación asume que los coeficientes son constantes pero el término de error está determinado por las mismas variables  $\mathbf{Z}_{ij}$  y un término de error aleatorio

$$U_{ij} = \mathbf{z}'_{ij}\boldsymbol{\beta} + \eta_{ij} \text{ where } \eta_{ij} = \mathbf{z}'_{ij}\boldsymbol{\gamma}_i + \varepsilon_{ij}$$

- Como  $U_{ij}$  y  $\eta_{ij}$  están determinadas por las mismas variables, ambas variables están correlacionadas y su correlación se puede expresar de la siguiente manera:

$$Cov(\eta_{ij}, \eta_{ig}) = \mathbf{z}'_{ij} \boldsymbol{\Sigma}_{\beta} \mathbf{z}_{ig}$$

- Finalmente, una especificación más general es una en la que se tiene un modelo para la utilidad, que depende de las variables  $\mathbf{x}_{ij}$  y de los parámetros  $\boldsymbol{\beta}_{ik}$ , los cuales se determinan por unas variables  $\mathbf{z}_i$  y de los parámetros  $\boldsymbol{\theta}$

$$U_{ij} = \mathbf{x}'_{ij} \boldsymbol{\beta}_i + \varepsilon_{ij} \quad \text{where} \quad \boldsymbol{\beta}_i = \boldsymbol{\beta} + \mathbf{z}'_i \boldsymbol{\theta} + u_{ij}$$

$$u_i \sim N(0, \sigma_u^2)$$

- La distribución de los parámetros  $\boldsymbol{\beta}_i$  será una distribución normal con los siguientes estadísticos:

$$E(\boldsymbol{\beta}_i) = \boldsymbol{\beta} + E(\mathbf{z}'_i) \boldsymbol{\theta} \quad Var(\boldsymbol{\beta}_i) = Var(\mathbf{z}'_i) \boldsymbol{\theta}^2$$

- La probabilidad de que se escoja una alternativa  $j$  se puede calcular a través de la esperanza condicional en las variables explicativas y  $\mathbf{u}_i$

$$E[P(Y_i = j | \mathbf{x}_{ij}, \mathbf{z}_i, \mathbf{u}_i)] = \int \frac{e^{(\mathbf{x}'_{ij} \boldsymbol{\beta}_i)}}{\sum_{l=1}^J e^{(\mathbf{x}'_{il} \boldsymbol{\beta}_i)}} \phi(\boldsymbol{\beta}_i) d\boldsymbol{\beta}_i$$

$$\text{where } \mathbf{u}_i = \begin{pmatrix} u_{i1} \\ \dots \\ u_{ij} \end{pmatrix}$$

- Debido a que los coeficientes dependen del término de error  $u_{ij}$ , se tiene que calcular la esperanza condicional a estos con tal de que las  $\boldsymbol{\beta}_i$  sean constantes para una  $u_{ij}$  determinada
- Para poder estimar esta probabilidad, se puede utilizar métodos numéricos. De este modo, se pueden utilizar las estimaciones para poder estimar la función logarítmica de máxima verosimilitud

$$\hat{E}[P(Y_i = j | \mathbf{x}_{ij}, \mathbf{z}_i, \mathbf{u}_i)] = \frac{1}{S} \sum_{s=1}^S \frac{e^{(\mathbf{x}'_{ij} \boldsymbol{\beta}_i^{(s)})}}{\sum_{l=1}^J e^{(\mathbf{x}'_{il} \boldsymbol{\beta}_i^{(s)})}}$$

$$\ln \hat{L} = \sum_{i=1}^N \sum_{j=1}^J y_{ij} \ln \{ \hat{E}[P(Y_i = j | \mathbf{x}_{ij}, \mathbf{z}_i, \mathbf{u}_i)] \}$$

## Los modelos de elección discreta multinomial ordenada

- Hay veces que las alternativas de la variable dependiente discreta están ordenadas, por lo que se revela la fuerza de las preferencias con respecto a un único resultado y se puede explotar este aspecto con modelos de elección discreta multinomial ordenada
  - Para cualquier individuo, se supone que hay una variable de utilidad continua subyacente que revela la fuerza de la preferencia y que depende de variables independientes para cada individuo

$$U_{ij}^* = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_{ij} \text{ where } \mathbf{x}_i = \begin{pmatrix} 1 \\ X_{i1} \\ \dots \\ X_{iK} \end{pmatrix} \text{ and } \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \dots \\ \beta_K \end{pmatrix}$$

- El modelo sugiere que cada individuo tiene unas características que determinan su utilidad para la alternativa  $j$  pero solo dependen del individuo, y un término de error que depende del individuo y de la alternativa  $j$
- Se asume que la variable subyacente o de utilidad se distribuye por toda la recta real

$$-\infty < U_{ij}^* < \infty$$

- La existencia de un orden entre las alternativas se traduce en la censura de esta variable subyacente, de modo que la elección de  $j$  es una versión censurada de la utilidad

$$Y_i = j \text{ if } \mu_{j-1} < U_{ij}^* \leq \mu_j \text{ for } j = 1, \dots, J$$

$$\text{when } \mu_0 = -\infty \text{ and } \mu_J = \infty$$

- De este modo, cada alternativa  $j$  se escoge cuando la utilidad subyacente está dentro del intervalo de esa alternativa

$$U_{ij} \in (\mu_{j-1}, \mu_j]$$

- Existirán  $J - 1$  límites  $\mu$  para dividir los intervalos de la variable subyacente para los cuales se identifican las observaciones
- La diferencia entre dos alternativas en la jerarquía o estructura de la variable dependiente (una alternativa  $j$  comparada con una alternativa  $g$ ) no es la misma que en la escala de la utilidad subyacente

- Eso quiere decir que se ha hecho una transformación estrictamente no lineal capturada por los límites  $\mu$ , por lo que se pueden estimar los coeficientes con un modelo de elección ordenada
- El modelo *probit* ordenado se construye en base a una regresión latente similar a la del modelo *probit* binomial, solo que en este se tiene en cuenta que las alternativas están ordenadas en la variable dependiente
  - Se asume que existe una variable no observada  $U_{ij}^*$  para un individuo  $i$  y una alternativa  $j$ , la cual depende de características individuales y de un término de error que también depende de la alternativa  $j$

$$U_{ij}^* = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_{ij} \quad \text{where} \quad \mathbf{x}_i = \begin{pmatrix} 1 \\ X_{i1} \\ \dots \\ X_{iK} \end{pmatrix} \quad \text{and} \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \dots \\ \beta_K \end{pmatrix}$$

- Se asume que el término de error se distribuye de manera normal con media nula y varianza unitaria

$$\varepsilon_{ij} \sim N(0,1)$$

- No obstante, se observa una variable  $Y_i$ , la cual está censurada para diferentes intervalos que están determinados por los límites  $\mu$ , los cuales son desconocidos y se tienen que estimar con las  $\boldsymbol{\beta}$

$$Y_i = j \quad \text{if} \quad \mu_{j-1} < U_{ij}^* \leq \mu_j \quad \text{for} \quad j = 1, \dots, J$$

$$\text{when} \quad \mu_0 = -\infty \quad \text{and} \quad \mu_J = \infty$$

- A partir de las suposiciones, se puede ver que la probabilidad de que se escoja una alternativa  $j$  será la probabilidad de que el término de error esté dentro de un intervalo determinado por los límites  $\mu$  y las características individuales  $\mathbf{X}_i$

$$P(Y_i = j | \mathbf{x}_i) = P(\mu_{j-1} < U_{ij}^* \leq \mu_j | \mathbf{x}_i) = P(\mu_{j-1} < \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_{ij} \leq \mu_j | \mathbf{x}_i)$$

$$= P(\mu_{j-1} - \mathbf{x}_i' \boldsymbol{\beta} < \varepsilon_{ij} \leq \mu_j - \mathbf{x}_i' \boldsymbol{\beta} | \mathbf{x}_i) = \Phi(\mu_j - \mathbf{x}_i' \boldsymbol{\beta}) - \Phi(\mu_{j-1} - \mathbf{x}_i' \boldsymbol{\beta})$$

- Como las probabilidades tienen que ser positivas, es necesario que los límites sean estrictamente crecientes

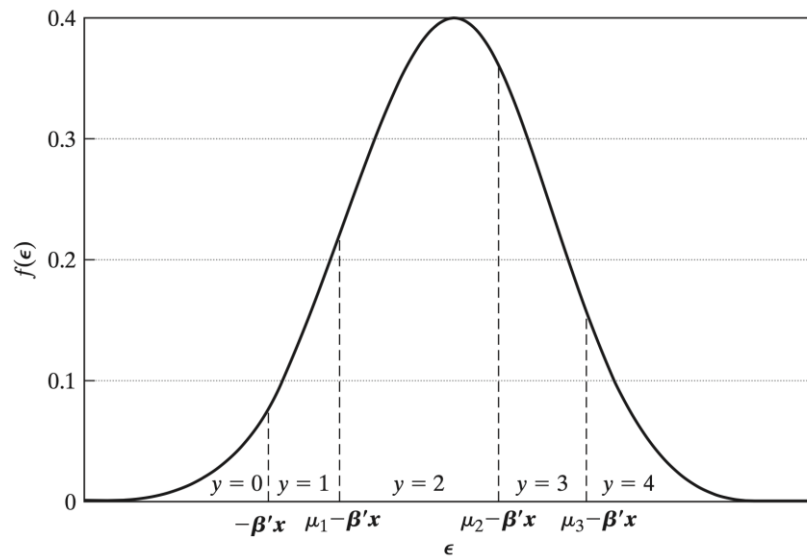
$$\mu_1 < \mu_2 < \dots < \mu_{J-1} < \mu_J$$

- Debido a que se asume que  $\mu_0 = -\infty$  y que  $\mu_J = \infty$ , se puede ver como la probabilidad de la primera opción y de la última se pueden expresar de la siguiente manera:

$$P(Y_i = 1|x_i) = \Phi(\mu_1 - x_i'\beta)$$

$$P(Y_i = J|x_i) = 1 - \Phi(\mu_{J-1} - x_i'\beta)$$

- La estructura del modelo permite ver como la distribución del error permite obtener la probabilidad para una alternativa, en donde el valor  $z$  depende de  $\mu$  y de  $x_i'\beta$

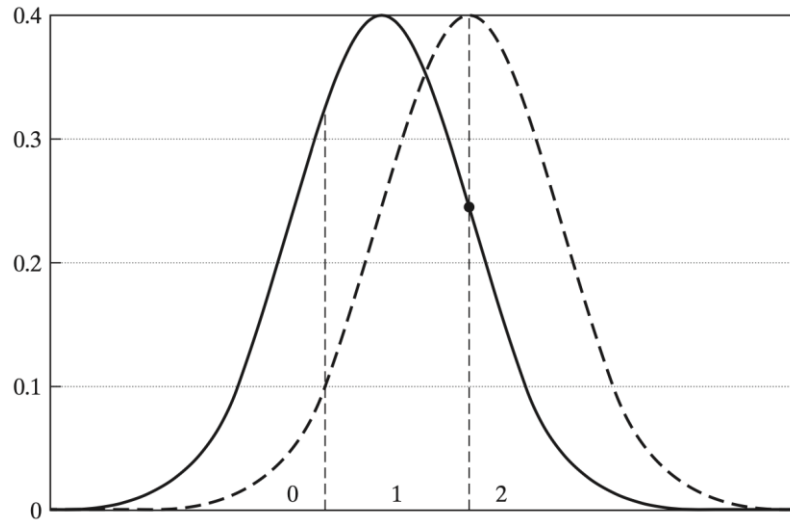


- Los efectos marginales en el modelo dependerán de  $\mu$  y de  $x_i'\beta$ , dado que la probabilidad depende de la función de distribución para  $\mu - x_i'\beta$ , y eso hace que se tenga una interpretación diferente

$$\frac{\partial P(Y_i = j|x_i)}{\partial X_{ik}} = [\phi(\mu_{j-1} - x_i'\beta) - \phi(\mu_j - x_i'\beta)]\beta_{ik}$$

- Siendo  $\mu$  y  $\beta$  constantes, un incremento unitario de una variable  $X_{ik}$  es equivalente a desplazar la distribución, por lo se reduce la masa de probabilidad de la primera o la última categoría (dependiendo de  $\beta$ ) y eso hace que aumenten y disminuyan las probabilidades de diferentes categorías





- Debido a que el efecto marginal de la primera categoría tiene un signo contrario al del coeficiente  $\beta_{ik}$ , mientras que el de la última tiene el mismo signo, el desplazamiento de la masa provocado por un incremento unitario en  $X_{ik}$  hace que se tenga un efecto opuesto en la probabilidad entre la primera y la última categoría

$$\text{If } \beta_{ik} > 0 \Rightarrow \begin{cases} \Delta P(Y_i = 1 | \mathbf{x}_i) < 0 \\ \Delta P(Y_i = J | \mathbf{x}_i) > 0 \end{cases}$$

$$\text{If } \beta_{ik} < 0 \Rightarrow \begin{cases} \Delta P(Y_i = 1 | \mathbf{x}_i) > 0 \\ \Delta P(Y_i = J | \mathbf{x}_i) < 0 \end{cases}$$

$$\frac{\partial P(Y_i = 1 | \mathbf{x}_i)}{\partial X_{ik}} = -\phi(\mu_j - \mathbf{x}_i' \boldsymbol{\beta}) \beta_{ik}$$

$$\frac{\partial P(Y_i = J | \mathbf{x}_i)}{\partial X_{ik}} = \phi(\mu_{j-1} - \mathbf{x}_i' \boldsymbol{\beta}) \beta_{ik}$$

- No obstante, el efecto en las categorías intermedias depende de dos valores diferentes ambas funciones de densidad, de modo que el efecto es ambiguo si no se estima
- El modelo *logit* ordenado se construye en base a una regresión latente similar a la del modelo *logit* binomial, solo que en este se tiene en cuenta que las alternativas están ordenadas en la variable dependiente
  - Se asume que existe una variable no observada  $U_{ij}^*$  para un individuo  $i$  y una alternativa  $j$ , la cual depende de características individuales y de un término de error que también depende de la alternativa  $j$

$$U_{ij}^* = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_{ij} \text{ where } \mathbf{x}_i = \begin{pmatrix} 1 \\ X_{i1} \\ \dots \\ X_{iK} \end{pmatrix} \text{ and } \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \dots \\ \beta_K \end{pmatrix}$$

- Se asume que el término de error se distribuye como en un modelo *logit* multinomial

$$\varepsilon_{ij} \sim \text{Gumbel} \quad F(\varepsilon_{ij}) = \exp\left(-\frac{1}{e^{\varepsilon_{ij}}}\right)$$

- No obstante, se observa una variable  $Y_i$ , la cual está censurada para diferentes intervalos que están determinados por los límites  $\mu$ , los cuales son desconocidos y se tienen que estimar con las  $\boldsymbol{\beta}$

$$Y_i = j \text{ if } \mu_{j-1} < U_{ij}^* \leq \mu_j \text{ for } j = 1, \dots, J$$

$$\text{when } \mu_0 = -\infty \text{ and } \mu_J = \infty$$

- A partir de las suposiciones, se puede ver que la probabilidad de que se escoja una alternativa  $j$  será la probabilidad de que el término de error esté dentro de un intervalo determinado por los límites  $\mu$  y las características individuales  $\mathbf{X}_i$

$$\begin{aligned} P(Y_i = j | \mathbf{X}_i) &= P(\mu_{j-1} < U_{ij}^* \leq \mu_j | \mathbf{x}_i) = P(\mu_{j-1} < \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_{ij} \leq \mu_j | \mathbf{x}_i) \\ &= P(\mu_{j-1} - \mathbf{x}_i' \boldsymbol{\beta} < \varepsilon_{ij} \leq \mu_j - \mathbf{x}_i' \boldsymbol{\beta} | \mathbf{x}_i) = \frac{e^{(\mu_j - \mathbf{x}_i' \boldsymbol{\beta})} - e^{(\mu_{j-1} - \mathbf{x}_i' \boldsymbol{\beta})}}{\sum_{r=1}^J e^{(\mu_r - \mathbf{x}_i' \boldsymbol{\beta})}} \end{aligned}$$

- Como las probabilidades tienen que ser positivas, es necesario que los límites sean estrictamente crecientes

$$\mu_1 < \mu_2 < \dots < \mu_{J-1} < \mu_J$$

- Debido a que se asume que  $\mu_0 = -\infty$  y que  $\mu_J = \infty$ , se puede ver como la probabilidad de la primera opción y de la última se pueden expresar de la siguiente manera:

$$P(Y_i = 1 | \mathbf{X}_i) = \frac{e^{(\mu_1 - \mathbf{x}_i' \boldsymbol{\beta})}}{\sum_{r=1}^J e^{(\mu_r - \mathbf{x}_i' \boldsymbol{\beta})}}$$

$$P(Y_i = J | \mathbf{X}_i) = 1 - \frac{e^{(\mu_{J-1} - \mathbf{x}_i' \boldsymbol{\beta})}}{\sum_{r=1}^J e^{(\mu_r - \mathbf{x}_i' \boldsymbol{\beta})}}$$

- La estructura del modelo permite ver como la distribución del error permite obtener la probabilidad para una alternativa, en donde el valor  $z$  depende de  $\mu$  y de  $\mathbf{x}_i'\boldsymbol{\beta}$
- Los efectos marginales en el modelo dependerán de  $\mu$  y de  $\mathbf{x}_i'\boldsymbol{\beta}$ , dado que la probabilidad depende de la función de distribución para  $\mu - \mathbf{x}_i'\boldsymbol{\beta}$ , y eso hace que se tenga una interpretación diferente

$$\frac{\partial P(Y_i = j | \mathbf{x}_i)}{\partial X_{ik}} = [f(\mu_{j-1} - \mathbf{x}_i'\boldsymbol{\beta}) - f(\mu_j - \mathbf{x}_i'\boldsymbol{\beta})]\beta_{ik}$$

- Debido a que el efecto marginal de la primera categoría tiene un signo contrario al del coeficiente  $\beta_{ik}$ , mientras que el de la última tiene el mismo signo, el desplazamiento de la masa provocado por un incremento unitario en  $X_{ik}$  hace que se tenga un efecto opuesto en la probabilidad entre la primera y la última categoría

$$\text{If } \beta_{ik} > 0 \Rightarrow \begin{cases} \Delta P(Y_i = 1 | \mathbf{x}_i) < 0 \\ \Delta P(Y_i = J | \mathbf{x}_i) > 0 \end{cases}$$

$$\text{If } \beta_{ik} < 0 \Rightarrow \begin{cases} \Delta P(Y_i = 1 | \mathbf{x}_i) > 0 \\ \Delta P(Y_i = J | \mathbf{x}_i) < 0 \end{cases}$$

$$\frac{\partial P(Y_i = 1 | \mathbf{x}_i)}{\partial X_{ik}} = -f(\mu_j - \mathbf{x}_i'\boldsymbol{\beta})\beta_{ik}$$

$$\frac{\partial P(Y_i = J | \mathbf{x}_i)}{\partial X_{ik}} = f(\mu_{j-1} - \mathbf{x}_i'\boldsymbol{\beta})\beta_{ik}$$

- No obstante, el efecto en las categorías intermedias depende de dos valores diferentes ambas funciones de densidad, de modo que el efecto es ambiguo si no se estima
- La formulación del modelo de elección ordenada tiene una suposición implícita debido a como se construye, llamada suposición de regresiones paralelas o *parallel regression assumption*
  - A partir del modelo, es posible construir las siguientes variables binarias:

$$w_{ij} = \begin{cases} 1 & \text{if } Y_i > j \\ 0 & \text{if } Y_i \leq j \end{cases} \quad \text{for } j = 1, 2, \dots, J-1$$

$$\text{with } P(w_{ij} = 1 | \mathbf{x}_i) = F(\mathbf{x}_i'\boldsymbol{\beta} - \mu_j)$$

- En este caso, la función de distribución de probabilidad depende del modelo utilizado

$$F(\mathbf{x}_i'\boldsymbol{\beta} - \mu_j) = \begin{cases} \Phi(\mathbf{x}_i'\boldsymbol{\beta} - \mu_j) & \text{if ord. probit} \\ \frac{e^{(\mathbf{x}_i'\boldsymbol{\beta} - \mu_j)}}{\sum_{r=1}^J e^{(\mathbf{x}_i'\boldsymbol{\beta} - \mu_r)}} & \text{if ord. logit} \end{cases}$$

- Se asume que la probabilidad de que  $w_{ij} = 1$  es la probabilidad de que  $Y_i > j$  para  $j = 1, 2, \dots, J-1$

$$P(Y_i \leq j | \mathbf{x}_i) = P(w_{ij} = 0 | \mathbf{x}_i) = F(\mu_j - \mathbf{x}_i'\boldsymbol{\beta})$$

$$P(Y_i > j | \mathbf{x}_i) = P(w_{ij} = 1 | \mathbf{x}_i) = 1 - F(\mu_j - \mathbf{x}_i'\boldsymbol{\beta}) = F(\mathbf{x}_i'\boldsymbol{\beta} - \mu_j)$$

- Este enfoque muestra como, a partir de  $w_{ij}$  y  $\mathbf{x}_i$ , es posible estimar el mismo vector  $\boldsymbol{\beta}$  para cada regresión binaria que se haga de  $w_{ij}$  sobre  $\mathbf{x}_i$ , aunque el parámetro constante sea diferente para cada una

- Esta igualdad entre vectores de parámetros es la suposición de regresiones paralelas, en donde el nombre viene del hecho de que la pendiente (el vector) es el mismo para cada regresión, pero el término constante difiere (por lo que las regresiones serían paralelas)
- Brant, en su investigación de 1990, mostró el caso para el *proportional odds model* para un modelo *logit* ordenado, propuesto por McCullagh en 1980, en donde se obtiene el mismo resultado de manera más clara

$$\begin{aligned} \frac{P(Y_i \leq j | \mathbf{x}_i)}{1 - P(Y_i \leq j | \mathbf{x}_i)} &= \frac{P(Y_i \leq j | \mathbf{x}_i)}{P(Y_i > j | \mathbf{x}_i)} = e^{(\mathbf{x}_i'\boldsymbol{\beta} - \mu_j)} \\ \Rightarrow \ln \left[ \frac{P(Y_i \leq j | \mathbf{x}_i)}{P(Y_i > j | \mathbf{x}_i)} \right] &= \mathbf{x}_i'\boldsymbol{\beta} - \mu_j \end{aligned}$$

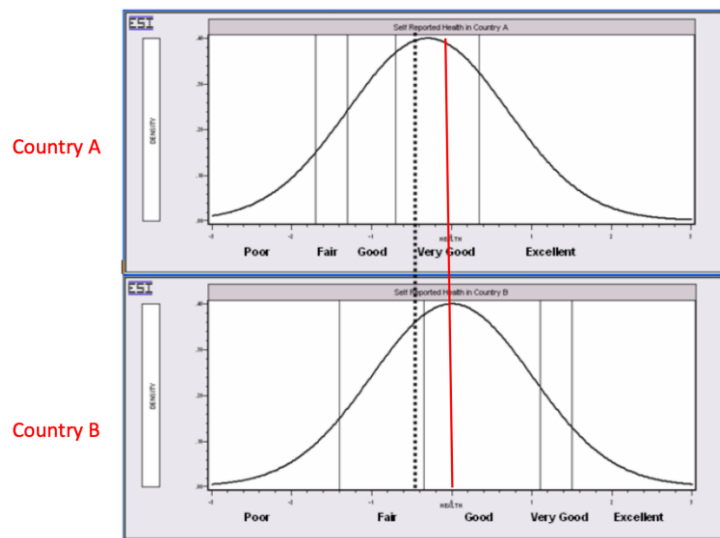
- Debido a este hecho, Brant propuso un test de contraste de hipótesis para contrastar si se cumplía la suposición de regresiones paralelas a través de un test de Wald

$$\begin{cases} H_0: \boldsymbol{\beta}_1 = \boldsymbol{\beta}_2 = \dots = \boldsymbol{\beta}_{J-1} \\ H_1: \text{Otherwise} \end{cases}$$

$$W = (\hat{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}})[\text{Var}(\hat{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}} | \mathbf{x}_i)]^{-1}(\hat{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}})' \sim \chi_{(J-2)K}^2$$

$$\text{where } \hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{\boldsymbol{\beta}}_1 \\ \dots \\ \hat{\boldsymbol{\beta}}_{J-1} \end{pmatrix} \text{ and } \bar{\boldsymbol{\beta}} = \begin{pmatrix} \boldsymbol{\beta}_1 \\ \dots \\ \boldsymbol{\beta}_{J-1} \end{pmatrix}$$

- En este test, rechazar la hipótesis nula quiere decir que el modelo ordenado no es adecuado para poder modelar los datos, dado que la suposición en la que se basa no se cumple
  - No obstante, el uso de otros modelos en donde  $\beta$  dependa de la alternativa  $j$  comporta problemas, dado que el modelo latente dependería de  $j$  (que se realiza con la realización de  $Y_i$ ) pero  $Y_i$  depende del modelo latente. Esto haría que no se preserve el orden de las probabilidades ni se mantenga la coherencia en el modelo
  - Algunas razones propuestas por las cuales se puede haber rechazado la hipótesis nula son una mala especificación de  $X_i'\beta$ , la heteroscedasticidad de  $\varepsilon$  o una mala especificación de la forma distribucional de la variable latente
- El modelo estándar para alternativas ordenadas tiene algunas limitaciones a la hora de comparar y de interpretar los resultados
    - Una de las limitaciones más importantes es que hacer comparaciones entre poblaciones o subgrupos de una población puede comportar problemas debidos a sesgos y heterogeneidad
      - La elección de una alternativa puede estar determinada por las alternativas escogidas en el pasado (por ejemplo, con los estados de salud previos). A esto se le llama *dependent-state bias*
      - Dependiendo de la escala de referencia que se utilice, dos individuos que realmente escogerían una misma alternativa pueden escoger dos diferentes. A esto se le llama sesgo de referencia de escala o *scale reference bias*
      - Puede ser que los individuos entiendan las escalas de manera diferente, de modo que dos individuos que realmente escogerían una misma alternativa pueden escoger dos alternativas diferentes debido a su entendimiento. A esto se le llama heterogeneidad de informe o *reporting heterogeneity*



- Otro problema es que diferentes subgrupos de una población utilizan niveles de límites sistemáticamente diferentes para las categorías, aunque realmente pertenezcan a la misma categoría
  - Esto se debe a que cada persona puede interpretar su estado o la categoría en la que está diferente dependiendo de sus características personales
  - Por ello, una solución puede ser incluir la posibilidad de heterogeneidad en los límites, de modo que se definiría el nivel límite de la alternativa  $j$  como una función  $g$  de las características personales  $Z_i$

$$\mu_{ij} = g(Z_i)$$

## Los modelos para el conteo de eventos

- Que la variable dependiente represente el número de ocurrencias de un evento hace que esta sea una medida cuantitativa como las que se utilizan en los modelos de regresión clásicos
  - No obstante, la preponderancia típica de ceros y valores pequeños y la naturaleza discreta de la variable de respuesta hace que se tenga que utilizar un enfoque que tenga en cuenta estos aspectos
    - Uno de los modelos más importantes es el modelo de Poisson, aunque las suposiciones de este han llevado a que se desarrollen modelos que puedan superar esas limitaciones como el modelo de binomial negativa

- Además, para tener en cuenta la preponderancia de ceros en la variable de respuesta se han desarrollado modelos con datos censurados y de otros tipos
- El modelo de regresión de Poisson especifica que cada  $y_i$  proviene de una distribución de Poisson con parámetro  $\lambda_i$ , la cual está relacionada con los regresores  $x_i$ . La ecuación primordial del modelo es la siguiente:

$$P(Y = y_i | x_i) = \frac{e^{-\lambda_i} \lambda_i^{y_i}}{y_i!} \quad \text{for } y_i = 0, 1, 2, \dots$$

- La formulación más común para  $\lambda_i$  es el modelo de regresión log-lineal

$$\ln \lambda_i = x_i \beta$$

- Se puede demostrar fácilmente que el número de eventos esperados por periodos o por unidad de espacio y, por tanto, su varianza condicional (la cual será heterocedástica) y el efecto marginal son los siguientes:

$$E(y_i | x_i) = Var(y_i | x_i) = \lambda_i = e^{x_i' \beta} \quad \frac{\partial E(y_i | x_i)}{\partial x_i} = \lambda_i \beta$$

- Si se quisiera obtener el efecto medio del tratamiento, es posible realizar una pequeña modificación al modelo. De este modo, se obtienen las siguientes equivalencias:

$$E(y_i | x_i) = e^{x_i' \beta + \gamma T} \Rightarrow ATE = \frac{1}{n} \sum_{i=1}^n [e^{x_i' \beta + \gamma T} - e^{x_i' \beta}]$$

- El ATET se puede calcular a través de hacer esta media con las observaciones en las que  $T = 1$
- Debido a que el modelo de Poisson no es lineal, los coeficientes de los regresores no son directamente interpretables como el efecto marginal:

$$\frac{\partial E(y_i | x_i)}{\partial x_k} = \frac{\partial \lambda_i}{\partial x_k} = \beta_k \lambda_i = \beta_k e^{x_i' \beta}$$

- No se puede interpretar su magnitud (porque no es el coeficiente únicamente), pero si se puede interpretar el signo de los coeficientes
- El efecto marginal medio estimado, al ser la media de los efectos marginales estimados para cada individuo, también se puede

expresar como el producto entre la media de la variable  $y$  y el coeficiente del regresor  $x_k$  para el que se estima el efecto

- Aunque la magnitud no sea interpretable en el efecto marginal, sí tiene una interpretación para el efecto sobre el valor esperado de  $\lambda_i$  (dado que depende de los coeficientes)

$$e^{\beta_0 + \beta_1 x_1 + \dots + \beta_k (x_k + 1)} = e^{\beta_0 + \beta_1 x_{1i} + \dots + \beta_k x_{ki}} e^{\beta_k} = \lambda_i e^{\beta_k}$$

- De este modo, la variación de una unidad en un regresor hace que  $E(y_i | x_i)$  se multiplique por la exponencial del coeficiente del regresor
- Cuando la variable está en logaritmos, sin embargo, el coeficiente  $\beta_k$  se puede interpretar como la elasticidad gracias a las propiedades de los exponentes

$$y = e^{\ln x_k} \Rightarrow \ln y = \ln x_k$$

- El modelo de Poisson no produce una contraparte natural al  $R^2$  en un modelo de regresión lineal, dado que la media condicional no es lineal y la regresión es heteroscedástica. Por lo tanto, se han propuesto medidas de bondad del ajuste alternativas

- Se puede utilizar una medida de basada en los residuos estandarizados, la cual compara el ajuste del modelo con el ajuste proporcionado por un modelo solo con el término constante. No obstante, esta medida puede ser negativa y puede aumentar cuando se sacan regresores del modelo

$$R_p^2 = 1 - \frac{\sum_{i=1}^n \left[ (y_i - \hat{\lambda}_i) / \sqrt{\hat{\lambda}_i} \right]}{\sum_{i=1}^n \left[ (y_i - \bar{y}) / \sqrt{\bar{y}} \right]}$$

- Para una observación individual, la desviación del modelo permite obtener la suma de las desviaciones, las cuales se pueden utilizar como una medida de bondad del ajuste (siempre que se asuma que se incluye un término constante)

$$d_i = 2 \left[ y_i \ln \left( \frac{y_i}{\hat{\lambda}_i} \right) - (y_i - \hat{\lambda}_i) \right] = 2 \left[ y_i \ln \left( \frac{y_i}{\hat{\lambda}_i} \right) - e_i \right]$$

$$\text{where } 0 \ln(0) = 0$$

$$\Rightarrow G^2 = 2 \sum_{i=1}^n \left[ y_i \ln \left( \frac{y_i}{\hat{\lambda}_i} \right) - e_i \right] = 2 \sum_{i=1}^n y_i \ln \left( \frac{y_i}{\hat{\lambda}_i} \right)$$



- Cameron y Windmeijer, en su investigación de 1993, sugirieron el uso de una medida llamada  $R_d^2$ , la cual se puede expresar en términos de la función de la log-verosimilitud. Tanto el numerador como el denominador miden la mejora del ajuste del modelo sobre el modelo con un solo término constante, y el denominador mide la mejora máxima posible. Esta se define de la siguiente manera:

$$R_d^2 = 1 - \frac{\sum_{i=1}^n \left[ y_i \ln \left( \frac{y_i}{\hat{\lambda}_i} \right) - e_i \right]}{\sum_{i=1}^n y_i \ln \left( \frac{y_i}{\bar{y}} \right)} = \frac{\ln L(\hat{\lambda}_i | y_i) - \ln L(\bar{y} | y_i)}{\ln L(\lambda_i | y_i) - \ln L(\bar{y} | y_i)} \in [0,1]$$

- Algunos programas utilizan una medida de bondad del ajuste parecida a  $R_d^2$ , llamada pseudo- $R^2$  o *likelihood ratio index*, la cual se define de la siguiente manera:

$$R_{pseudo}^2 = 1 - \frac{\ln L(\hat{\lambda}_i | y_i)}{\ln L(\bar{y} | y_i)}$$

- El modelo de Poisson se ha criticado porque hace una suposición implícita de que la varianza de  $y_i$  iguala a la media, el primer paso del análisis es contrastar esta suposición implícita en el contexto de un modelo simple

$$H_0 : Var(y_i) = E(y_i)$$

$$H_1 : Var(y_i) = E(y_i) + \alpha g[E(y_i)]$$

- Un contraste propuesto por Cameron y Trivedi se puede realizar a través de hacer una regresión para  $z_i$ , donde  $\hat{\lambda}_i$  es el valor predicho de la regresión. Una vez hecho esto, un simple contraste  $t$  de si el vector de coeficientes es significativamente diferente de cero

$$z_i = \frac{(y_i - \hat{\lambda}_i)^2 - y_i}{\hat{\lambda}_i \sqrt{2}} \Rightarrow T = \frac{\beta}{Std.Err(\beta)}$$

- Como la igualdad asumida entre la esperanza condicional y la función de varianza es una desventaja del modelo de Poisson, la alternativa más utilizada es la distribución negativa binomial, que nace de una formulación natural de heterogeneidad transversal
- Se generaliza el modelo de Poisson introduciendo un efecto individual no observable en la media condicional, en donde  $\varepsilon_i$  refleja el error de especificación o la heterogeneidad

$$\ln \mu_i = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i = \ln \lambda_i + \ln u_i$$

- La distribución de  $y_i$  condicionada a  $\mathbf{x}_i$  y  $u_i$  sigue siendo una distribución de Poisson con media condicional y varianza  $\mu_i$

$$f(y_i | \mathbf{x}_i, u_i) = \frac{e^{-\lambda_i u_i} (\lambda_i u_i)^{y_i}}{y_i!}$$

- La distribución incondicional  $f(y_i | \mathbf{x}_i)$  es el valor esperado de  $u_i$  de  $f(y_i | \mathbf{x}_i, u_i)$ , de modo que  $u_i$  define la distribución incondicional

$$f(y_i | \mathbf{x}_i) = \int_0^\infty \frac{e^{-\lambda_i u_i} (\lambda_i u_i)^{y_i}}{y_i!} g(u_i) du_i$$

- Por conveniencia matemática, se suele asumir una distribución gamma se asume para  $u_i = \exp(\varepsilon_i)$ . Como la media de la distribución no está identificada si el modelo contiene un término constante (porque las perturbaciones entran multiplicativamente), por lo que se asume que  $E(u_i)$  y esta normalización causa que la función  $g$  sea de la siguiente manera:

$$g(u_i) = \frac{\theta^\theta}{\Gamma(\theta)} e^{-\theta u_i} u_i^{\theta-1}$$

- En consecuencia, la función de densidad para  $y_i$  con esta especificación es una forma de la distribución binomial negativa, la cual tiene media condicional  $\lambda_i$  y varianza  $\lambda_i [1 + (1/\theta)\lambda_i]$

$$\begin{aligned} f(y_i | \mathbf{x}_i) &= \int_0^\infty \frac{e^{-\lambda_i u_i} (\lambda_i u_i)^{y_i}}{y_i!} \frac{\theta^\theta}{\Gamma(\theta)} e^{-\theta u_i} u_i^{\theta-1} du_i = \\ &= \frac{\Gamma(\theta + y_i)}{\Gamma(y_i + 1)\Gamma(\theta)} r_i^{y_i} (1 - r_i)^\theta \quad \text{where } r_i = \frac{\lambda_i}{\lambda_i + \theta} \end{aligned}$$

- Este modelo se denomina binomial negativo y puede ser estimado mediante máxima verosimilitud. Un contraste para la distribución de Poisson normalmente se lleva a cabo para la hipótesis  $\alpha = 1/\theta = 0$  usando el contraste de Wald o de razón de verosimilitud
- Los efectos marginales, en este caso, son los mismos que los efectos marginales en el caso en donde la variable dependiente seguía una distribución de Poisson, debido a que la esperanza condicionada es la misma

$$\frac{\partial E(y_i|x_i)}{\partial x_k} = \frac{\partial \lambda_i}{\partial x_k} = \beta_k \lambda_i = \beta_k e^{x_i' \beta}$$

- No se puede interpretar la magnitud del coeficiente (porque no es el coeficiente únicamente), pero si se puede interpretar el signo de los coeficientes
  - El efecto marginal medio estimado, al ser la media de los efectos marginales estimados para cada individuo, también se puede expresar como el producto entre la media de la variable y y el coeficiente del regresor  $x_k$  para el que se estima el efecto
  - Aunque la magnitud no sea interpretable en el efecto marginal, sí tiene una interpretación para el efecto sobre el valor esperado de  $\lambda_i$  igual que en el modelo de Poisson
- Con tal de poder modelar de manera más precisa el conteo de eventos con la presencia de muchos ceros (en relación a la cantidad esperada) se han utilizado modelos valla o *hurdle* y modelos inflados a cero o *zero-inflated models*

- Algunos investigadores han analizado una extensión para el modelo valla en la que el resultado de cero puede producirse en dos regímenes
  - En un régimen, el resultado siempre es cero, mientras que, en el otro, el proceso de Poisson es el que puede producir ceros u otros valores (por lo que los ceros se pueden producir en ambos regímenes)

$$P(y_i = 0|x_i) = P(\text{regime 1}) + P(y_i = 0|x_i, \text{regime 2})P(\text{regime 2})$$

$$P(y_i = j|x_i) = P(y_i = j|x_i, \text{regime 2})P(\text{regime 2}) \text{ for } j = 1, 2, \dots$$

- Siendo  $z$  un indicador binario del régimen (0 si se está en el primer régimen y 1 si está en el segundo régimen) y  $y^*$  el resultado del proceso de Poisson del segundo régimen, entonces la  $y$  observada es  $z \times y^*$
- Una extensión natural del modelo anterior es permitir que  $z$  se determine por un conjunto de variables que pueden ser iguales o no a las presentes en el proceso de Poisson. Esta extensión se llama modelo inflado a cero o *zero-inflated model* (ZIP)

$$P(z_i = 0|w_i) = F(w_i' \gamma) \text{ (Regime 1)}$$

$$P(y_i = j|x_i, z_i = 1) = \frac{\exp(-\lambda_i) \lambda_i^j}{j!} \text{ (Regime 2)}$$

- El modelo inflado a cero puede verse como un tipo de modelo latente. Las dos clases de probabilidades son  $F(\mathbf{w}_i, \gamma)$  y  $1 - F(\mathbf{w}_i, \gamma)$  y los dos regímenes son  $y = 0$  y el proceso de generación de datos Poisson. Se puede escoger una  $F(\mathbf{w}_i, \gamma)$  *probit* o *logit*
- De manera alternativa, el modelo se puede expresar de la siguiente manera, dado que la formulación anterior no limita el dominio de  $j$

$$P(y_i = 0 | \mathbf{x}_i, \mathbf{w}_i, \gamma) = F(\mathbf{w}_i, \gamma) + [1 - F(\mathbf{w}_i, \gamma)]P(y_i = 0 | \mathbf{x}_i, z_i = 1)$$

$$P(y_i = j | \mathbf{x}_i, \mathbf{w}_i, \gamma) = [1 - F(\mathbf{w}_i, \gamma)]P(y_i = j | \mathbf{x}_i, z_i = 1) \text{ for } j > 0$$

- Se puede hacer una extensión del modelo ZIP para incluir un proceso de generación de datos binomial negativo, llamado ZINB. En este modelo también se puede escoger una  $F(\mathbf{w}_i, \gamma)$  *probit* o *logit*
- La media de la variable aleatoria en el caso de Poisson sería la siguiente:

$$E(y_i | \mathbf{x}_i, \mathbf{w}_i) = F_i \times 0 + (1 - F_i) \times E(y_i^* | \mathbf{x}_i, z_i = 1) = (1 - F_i)\lambda_i$$

- Aunque sería interesante contrastar si hay un régimen divisor o no, el modelo básico anterior y el modelo ZIP no son modelos anidados
  - No se puede contrastar si todos los coeficientes son nulos para comprobarlo, dado que eso no elimina la estructura divisoria, sino que modifica las probabilidades
  - Como se quiere un contraste que verifique si se sigue un proceso de Poisson o no (dado que otra distribución podría generar más ceros que una Poisson), esto hace que el contraste sea un procedimiento difícil y que las hipótesis no sean anidadas, haciendo que el poder del contraste sea una función de la hipótesis alternativa
- Vuong, en su investigación de 1989 propuso un estadístico de contraste para modelos no anidados que se adecúa a la situación en la que se puede especificar la distribución alternativa
  - Siendo  $f_j(y_i | \mathbf{x}_i)$  la probabilidad predicha de que la variable  $Y = y_i$  bajo la suposición de que la distribución es  $f_j(y_i | \mathbf{x}_i)$  para  $j = 1, 2$  y definiendo  $m_i = \ln[f_1(y_i | \mathbf{x}_i)] - \ln[f_2(y_i | \mathbf{x}_i)]$ , entonces el estadístico para contrastar la hipótesis nula del modelo 1 respecto al 2 es el siguiente:

$$v = \frac{\frac{\sqrt{n}}{n} \sum_{i=1}^n m_i}{\sqrt{\frac{1}{n} \sum_{i=1}^n (m_i - \bar{m})^2}} = \frac{\bar{m}\sqrt{n}}{s_m} \sim N(0,1)$$

- Este estadístico es bidireccional y permite contrastar la hipótesis de que  $E(m_i) = 0$ . Si  $|v| < 2$  entonces el contraste no favorece a ningún modelo, pero para valores absolutos grandes, se apoya un modelo u otro (al 1 si es positivo o al 2 si es negativo)
  - Rechazar la hipótesis nula quiere decir rechazar que el proceso generador de datos es una Poisson común, de modo que el modelo ZIP se adecúa mejor a los datos. Se puede hacer del mismo modo con una distribución binomial negativa
  - No obstante, existe un debate sobre el uso de este contraste para contrastar los modelos ZIP contra los de conteo comunes, debido a la definición y prerequisites que determina Vuong para su uso
- En algunos contextos, el resultado cero del proceso generador de datos es cualitativamente diferente de los valores positivos, de modo que se utilizan modelos valla o *hurdle models*
- Este es el caso cuando el resultado proviene de una decisión separada sobre si participar o no en la actividad y, al decidir si participar o no, el individuo hace una decisión sobre cuánto o sobre la intensidad

$$P(y_i = 0 | \mathbf{w}_i) = F(\mathbf{w}_i, \boldsymbol{\gamma})$$

$$P(y_i = j | \mathbf{x}_i, \mathbf{w}_i, y_i > 0) = [1 - F(\mathbf{w}_i, \boldsymbol{\gamma})] \frac{\exp(-\lambda_i) \lambda_i^j}{j! [1 - \exp(-\lambda_i)]}$$

- Como se puede ver en esta formulación, los ceros se determinan en la primera ecuación, pero no en la segunda, de modo que se cambia la probabilidad de obtener un resultado nulo y se escalan las probabilidades que sobran para que sumen a uno
- Igual que antes, se puede escoger una forma funcional para  $F(\mathbf{w}_i, \boldsymbol{\gamma})$  que sea binaria como el *logit* o el *probit* y se puede extender el modelo a utilizar un proceso binomial negativo en la segunda ecuación
- La media condicional del modelo valla con un proceso de Poisson permite obtener fácilmente los efectos parciales

$$E[y_i|x_i, \mathbf{w}_i] = \frac{[1 - F(\mathbf{w}_i'\boldsymbol{\gamma})]\lambda_i}{1 - \exp(-\lambda_i)} \quad \text{where } \lambda_i = \exp(\mathbf{x}_i'\boldsymbol{\beta})$$

$$\frac{\partial E[y_i|x_i, \mathbf{w}_i]}{\partial x_i} = [1 - F(\mathbf{w}_i'\boldsymbol{\gamma})]\delta_i$$

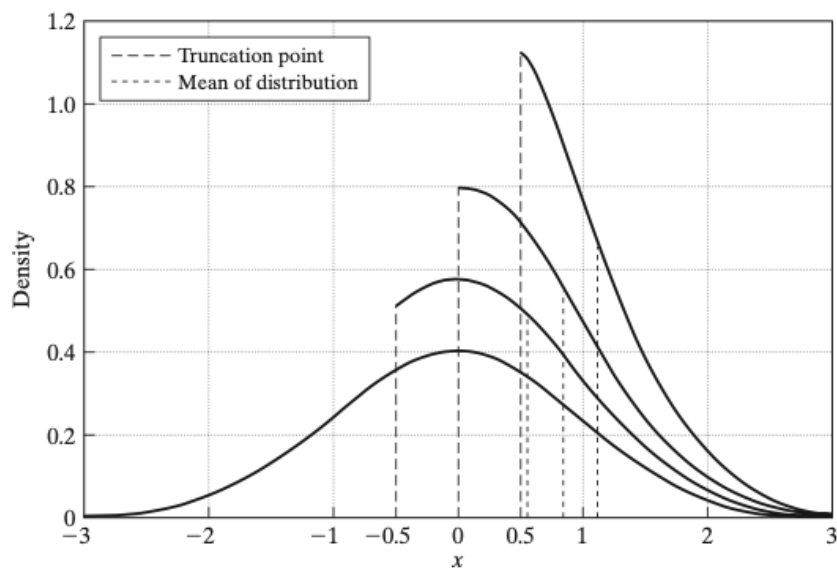
$$\frac{\partial E[y_i|x_i, \mathbf{w}_i]}{\partial \mathbf{w}_i} = \left[ -\frac{f(\mathbf{w}_i'\boldsymbol{\gamma})}{1 - \exp(-\lambda_i)} \right] \boldsymbol{\gamma}$$

- Para aquellas variables que aparecen tanto en  $\mathbf{x}_i$  como en  $\mathbf{w}_i$ , los efectos se suman entre sí para obtener el efecto total de esa variable. Para variables binarias, en cambio, las fórmulas dan un resultado aproximado, dado que se tendrían que tomar las diferencias de las medias condicionales cuando la variable pasa de cero a uno
- Puede ser de interés contrastar los efectos valla, pero igual que con el modelo ZIP, el modelo valla tiene hipótesis no anidadas (produciendo los mismos resultados que antes en la probabilidad), de modo que se suele utilizar el contraste de Vuong visto anteriormente en el contexto del modelo ZIP
- El modelo valla guarda similitudes con el ZIP, pero las implicaciones del comportamiento de los individuos en los modelos son diferentes y las modificaciones alteran la formulación de Poisson o de la binomial negativa
  - Ambos modelos tienen sobredispersión, por lo que no hay una igualdad entre la media y la varianza como en el proceso de Poisson. No obstante, esta sobredispersión no nace de la heterogeneidad, sino del proceso generador de ceros
  - Esto provoca un problema de identificación, dado que es difícil saber si el exceso de ceros en los datos se produce por la heterogeneidad o por el mecanismo divisor

## Los modelos de truncación y censura

- Muchas veces se tiene interés en analizar variables que no pueden ser observadas directamente porque estas muestran información incompleta porque están truncadas o censuradas
  - La truncación y la censura dificultan el análisis de datos que, de otro modo, se podrían analizar a través de modelos econométricos más simples como las regresiones

- Se puede considerar que la censura y la truncación son situaciones en las que no se tiene información completa sobre una variable latente  $Y^*$  de interés (solo de la variable  $Y$ )
- La truncación nace como consecuencia de descartar datos, tales como valores nulos, negativos o inusuales. La censura, en cambio, es una característica del diseño de muestreo en el que se representan ciertos valores con otros valores (de modo que no se saben los valores originales)
  - La truncación es esencialmente una característica de la distribución de la que se coge la muestra, mientras que la censura es un defecto en los datos muestrales, dado que, si no se censurara, los datos serían una muestra representativa de la población
  - La truncación y la censura producen unos sesgos sistemáticos al hacer inferencias sobre la población entera, llamados sesgo de truncación o *truncation bias* y sesgo de censura o *censorship bias*
- La truncación y la censura producen efectos similares en la distribución de valores de las variables aleatorias y en características como la media de esta distribución
- Una distribución truncada es la parte de una distribución no truncada que está por encima o por debajo de un valor específico, por lo que se coge como datos un subconjunto de la distribución entera



- Si una variable aleatoria  $x$  continua tiene una función de densidad de probabilidad  $f(x)$  y  $a$  es una constante, entonces la función de

densidad de la variable  $x$  truncada será una función de densidad condicional para valores mayores a  $a$  será la siguiente:

$$f(x|x > a) = \frac{f(x)}{P(x > a)}$$

- La demostración proviene de la definición de la probabilidad condicional y se basa en escalar la densidad para que al integrar la función se obtenga 1 en el rango de  $a$  (para cumplir con las propiedades de la probabilidad)
  - Una función muy similar aplica para el caso en que la truncación no es por abajo ( $x > a$ ) si no por arriba ( $x < a$ )
- Normalmente se utiliza la distribución normal truncada en aplicaciones recientes. Si la variable  $x$  sigue una distribución normal con media  $\mu$  y desviación estándar  $\sigma$ , entonces:

$$f(x|x > a) = \frac{f(x)}{1 - \Phi(\alpha)} = \frac{\frac{1}{\sigma} \phi\left(\frac{x - \mu}{\sigma}\right)}{1 - \Phi(\alpha)} \quad \text{where } \alpha = \frac{a - \mu}{\sigma}$$

- La probabilidad de que  $x > a$  cuando  $x$  se distribuye normalmente es  $1 - \Phi(\alpha)$

$$P(x > a) = 1 - \Phi\left(\frac{a - \mu}{\sigma}\right) = 1 - \Phi(\alpha)$$

- Se estandarizan la variable  $x$  y el valor  $a$  porque se utiliza la función de distribución y de densidad de una variable normal estándar. Es posible utilizar una distribución normal estándar truncada, por lo que no se estandarizarían los valores

$$P(x > a) = 1 - \Phi(a)$$

$$f(x|x > a) = \frac{f(x)}{1 - \Phi(a)} = \frac{\frac{1}{\sigma} \phi(x)}{1 - \Phi(a)}$$

- Las fórmulas anteriores se alteran ligeramente para variables truncadas por arriba, dado que la probabilidad cambia

$$P(x < a) = \Phi\left(\frac{a - \mu}{\sigma}\right) = \Phi(\alpha)$$

- Normalmente, uno se interesa por la media y la varianza de una variable aleatoria truncada, las cuales se obtendrían con las siguientes fórmulas generales:



$$E(x|x > a) = \int_a^{\infty} xf(x|x > a) dx$$

$$Var(x|x > a) = \int_a^{\infty} [x - E(x|x > a)]^2 f(x|x > a) dx$$

- Los momentos de una distribución normal se pueden obtener del mismo modo, resultando en las siguientes fórmulas:

$$E(x|x > a) = \mu + \sigma\lambda(\alpha)$$

$$Var(x|x > a) = \sigma^2[1 - \delta(\alpha)]$$

$$\text{where } \lambda(\alpha) = \begin{cases} \frac{\phi(\alpha)}{1 - \Phi(\alpha)} & \text{if truncation is } x > a \\ -\frac{\phi(\alpha)}{\Phi(\alpha)} & \text{if truncation is } x < a \end{cases}$$

$$\text{and } \delta(\alpha) = \lambda(\alpha)[\lambda(\alpha) - \alpha]$$

- La variable  $\lambda(\alpha)$  se denomina *ratio* inversa de Mills, y se define como el cociente entre la función de densidad normal y la función de distribución de probabilidad normal. Esta función es positiva en todo su dominio, y por definición, se cumple que:

$$-\frac{1}{2\sqrt{2\pi}} < \lambda(\alpha) < \frac{1}{2\sqrt{2\pi}} \Rightarrow 0 < \delta(\alpha) < 1$$

- La truncación tiene dos efectos importantes en los momentos estadísticos de la distribución de la variable:
  - Si la truncación es por abajo, la media será más alta que la media de la distribución original, mientras que, si la truncación es por arriba, la media será menor
  - Como  $0 < \delta(\alpha) < 1$ , la varianza de la distribución será menor que la varianza original, dado que incluye menos valores extremos
- A partir del concepto de truncación y de las propiedades de la distribución de una variable aleatoria truncada, es posible crear un modelo de regresión para lidiar con esta, llamado modelo de regresión truncada
  - Se asume que hay una variable  $Y_i^*$  latente, la cual está determinada por unos regresores  $X_i$  y un término de error  $\varepsilon_i$

$$Y_i^* = x_i'\beta + \varepsilon_i$$

- Se asume que el término de error está normalmente distribuido con media nula y varianza  $\sigma^2$ , de modo que la variable  $Y_i^*$  también sigue una distribución normal

$$\varepsilon_i \sim N(0, \sigma^2) \Rightarrow Y_i^* \sim N(\mathbf{x}_i' \boldsymbol{\beta}, \sigma^2)$$

- El foco principal está en la distribución de la variable truncada  $Y_i$ , la cual dependerá de la variable latente  $Y_i^*$  y de su esperanza, que es una función no lineal de  $a, \sigma, \mathbf{x}$  y  $\boldsymbol{\beta}$

$$E(Y_i^* | \mathbf{X}_i; Y_i^* > a) = \mathbf{x}_i' \boldsymbol{\beta} + \sigma \lambda(\alpha_i) = \mathbf{x}_i' \boldsymbol{\beta} + \sigma \frac{\phi(\alpha_i)}{1 - \Phi(\alpha_i)}$$

$$\text{where } \alpha_i = \frac{a - \mathbf{x}_i' \boldsymbol{\beta}}{\sigma}$$

- En este caso, la media de  $Y_i^*$  es  $\mathbf{X}_i' \boldsymbol{\beta}$ , por lo que  $\alpha_i$  dependerá del individuo y la esperanza también
- La esperanza se puede expresar de manera similar cuando la variable está truncada por arriba

$$E(Y_i^* | \mathbf{x}_i; Y_i^* > a) = \mathbf{x}_i' \boldsymbol{\beta} + \sigma \lambda(\alpha_i) = \mathbf{x}_i' \boldsymbol{\beta} - \sigma \frac{\phi(\alpha_i)}{\Phi(\alpha_i)}$$

- A través de la esperanza condicional de  $Y_i^*$ , es posible obtener el efecto marginal de los regresores para una variable truncada y el efecto en la varianza

$$\begin{aligned} \frac{\partial E(Y_i^* | \mathbf{x}_i; Y_i^* > a)}{\partial X_{ki}} &= \beta_k + \sigma \frac{\partial \lambda(\alpha_i)}{\partial \alpha_i} \frac{\partial \alpha_i}{\partial X_{ki}} = \\ &= \beta_k + \sigma \left[ \frac{\phi(\alpha_i)^2}{[1 - \Phi(\alpha_i)]^2} - \frac{\alpha_i \phi(\alpha_i)}{1 - \Phi(\alpha_i)} \right] \left( \frac{\beta_k}{\sigma} \right) = \\ &= \beta_k [1 + \lambda(\alpha_i)^2 - \alpha_i \lambda(\alpha_i)] = \beta_k [1 - \delta(\alpha_i)] \end{aligned}$$

$$Var(Y_i^* | \mathbf{x}_i, Y_i^* > a) = \sigma^2 [1 - \delta(\alpha_i)]$$

- Para  $Y_i^* < a$ , se obtiene una expresión equivalente que resulta en el mismo efecto marginal
- El efecto en  $Y_i$  de un incremento unitario de una variable  $X_{ki}$  se debe interpretar como el efecto que tiene sobre aquellos individuos que previamente tenían un valor dentro de la

distribución truncada, dado que solo se tienen esas observaciones en la muestra y no se puede generalizar el efecto para aquellos individuos que no estén dentro (no se tienen en cuenta las otras observaciones)

- Debido a que  $0 < \delta(\alpha_i) < 1$ , el efecto marginal de aumentar una unidad en el regresor  $X_{ki}$  es menor a su coeficiente en valor absoluto, y la varianza es menor a la que tendría la variable latente, por lo que hay una atenuación debido a la truncación

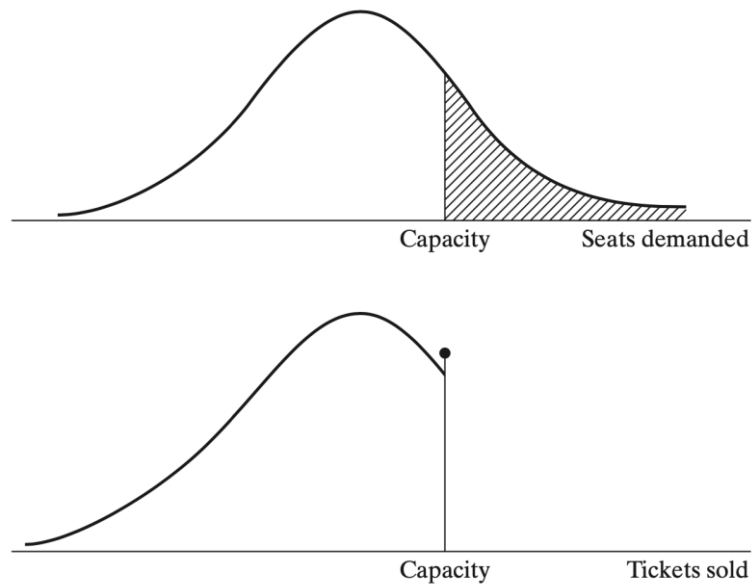
$$|ME_{trun}| \leq |\beta_k|$$

- Cuando un regresor está en el modelo con una forma polinómica o con una interacción, el efecto marginal no siempre tiene el mismo signo que el primer coeficiente del regresor, debido a que tiene dos o más y depende de los signos de estos coeficientes
  - Si el estudio econométrico se centra en la población, entonces los coeficientes  $\beta$  son los elementos de interés, dado que el efecto marginal obtenido es el de la subpoblación, pero el efecto marginal de las variables latentes serían los coeficientes
- Para la subpoblación de la que se obtienen los datos de la variable latente, se puede utilizar un modelo de regresión de la siguiente forma:

$$Y_i = E(Y_i^* | \mathbf{x}_i, Y_i^* > a) = \mathbf{x}_i' \beta + \sigma \lambda(\alpha_i) + u_i$$

$$\text{where } u_i = Y_i^* - E(Y_i^* | \mathbf{x}_i, Y_i^* > a)$$

- Por definición, el término de error  $u_i$  se distribuirá de manera normal con media nula, pero tiene una varianza equivalente a  $Var(Y_i^* | \mathbf{x}_i, Y_i^* > a)$  y es heteroscedástico (dado que la varianza de la variable truncada depende de  $\mathbf{x}_i$ )
  - Si se estima el modelo por mínimos cuadrados ordinarios, no obstante, los estimadores estarán sesgados porque se ignora el término no lineal  $\lambda(\alpha_i)$  y habría un sesgo por omisión de variables
- La teoría relevante para la distribución de una variable censurada es similar a la de una truncada, por lo que los aspectos relevantes son similares



- Cuando los datos están censurados, la distribución que aplica a la muestra de estos datos es una combinación de distribuciones discretas y continuas. Para analizar esta distribución, se define una variable  $Y_i$  transformada a partir de  $Y_i^*$

$$Y_i = \begin{cases} a & \text{if } Y_i^* \geq a \\ Y_i^* & \text{if } Y_i^* < a \end{cases}$$

- La distribución que aplica cuando  $Y_i^* \leq a$  es la función de distribución de probabilidad  $P(Y_i = a) = P(Y_i^* \geq a)$ , mientras que cuando  $Y^* < a$ , la función de densidad será  $f(Y) = f(Y_i^*)$ . De este modo, la probabilidad total es igual a 1 y se le asigna la probabilidad de la parte censurada al punto de censura  $a$
- La censura mostrada en el ejemplo es una censura a la derecha, pero también se puede censurar a la izquierda y censurar tanto a la izquierda como a la derecha

$$Y_i = \begin{cases} Y_i^* & \text{if } Y^* < a \\ a & \text{if } Y_i^* \geq a \end{cases} \quad Y_i = \begin{cases} a & \text{if } Y_i^* < a \\ Y_i^* & \text{if } a \leq Y_i^* \leq b \\ b & \text{if } Y_i^* > b \end{cases}$$

- Para una variable truncada, la única parte relevante para los cálculos de su distribución es la parte donde  $Y_i^* > a$  o  $Y_i^* < a$ , de modo que se escala la función de densidad con la probabilidad para que la integral sea igual a 1
- Se suele asumir que la variable aleatoria censurada sigue una distribución normal, dado que es lo que más se usa en las investigaciones. Por lo tanto, se pueden analizar los momentos estadísticos de la distribución censurada bajo esta suposición

If  $Y_i^* \sim N(\mu, \sigma^2)$  and  $\alpha = \frac{a - \mu}{\sigma}$ :

$$E(Y_i) = \Phi(\alpha)a + [1 - \Phi(\alpha)][\mu + \sigma\lambda(\alpha)]$$

$$Var(Y_i) = \sigma^2[1 - \Phi(\alpha)][1 - \delta(\alpha)] + [\alpha - \lambda(\alpha)]^2\Phi(\alpha)$$

- Las funciones de la *ratio* inversa de Mills y de la distribución normal se definen de la manera ya mencionada anteriormente con las variables truncadas. Para cada tipo de censura, se redefine la forma funcional de la probabilidad (dado que varía)
- La demostración de la esperanza bajo la suposición de distribución normal es la siguiente:

$$\begin{aligned} E(Y_i) &= E(Y_i|Y_i < a)P(Y_i < a) = \\ &= P(Y_i = a)E(Y_i|Y_i = a) + P(Y_i < a)E(Y_i|Y_i < a) = \\ &= P(Y_i^* \geq a)a + P(Y_i^* < a)E(Y_i^*|Y_i^* < a) = \\ &= P(Y_i^* \geq a)a + P(Y_i^* < a)[\mu + \sigma\lambda(\alpha)] \end{aligned}$$

- Para el caso especial en que  $a = 0$ , los momentos se pueden expresar de la siguiente manera:

$$\begin{aligned} E(Y_i) &= \left[1 - \Phi\left(-\frac{\mu}{\sigma}\right)\right]\left[\mu + \sigma\lambda\left(-\frac{\mu}{\sigma}\right)\right] \\ Var(Y_i) &= \sigma^2 \left[1 - \Phi\left(-\frac{\mu}{\sigma}\right)\right] \left[1 - \delta\left(-\frac{\mu}{\sigma}\right) + \left[-\frac{\mu}{\sigma} - \lambda\left(-\frac{\mu}{\sigma}\right)\right]^2 \Phi\left(-\frac{\mu}{\sigma}\right)\right] \end{aligned}$$

- Una variable puede estar censurada y truncada a la vez, de modo que hay valores de  $Y^*$  para los cuales se censura con un valor concreto, pero hay valores que quedan fuera de la muestra:

$$Y = \begin{cases} Y^* & \text{if } b \leq Y^* \leq c \\ b & \text{if } a < Y^* < b \end{cases}$$

- En este caso en concreto, la variable está truncada de  $b$  a  $c$  (valores por encima de  $c$  o por debajo de  $b$  no entran en la muestra), pero está censurada a la izquierda de  $b$  (los valores entre  $a$  y  $b$  tienen un valor  $Y = b$ )
- A partir del concepto de censura y de las propiedades de la distribución de una variable aleatoria censurada, es posible crear un modelo de regresión para lidiar con esta, llamado modelo Tobit o de regresión censurada

- Se asume que hay una variable  $Y_i^*$  latente, la cual está determinada por unos regresores  $\mathbf{X}_i$  y un término de error  $\varepsilon_i$ , y que hay una variable observada  $Y_i$  que está censurada

$$Y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i \quad Y_i = \begin{cases} a & \text{if } Y_i^* \geq a \\ Y_i^* & \text{if } Y_i^* < a \end{cases}$$

- La variable observada se puede censurar de diferentes maneras, de modo que se puede definir de maneras alternativas
- Se asume que el término de error está normalmente distribuido con media nula y varianza  $\sigma^2$ , de modo que la variable  $Y_i^*$  también sigue una distribución normal

$$\varepsilon_i \sim N(0, \sigma^2) \Rightarrow Y_i^* \sim N(\mathbf{x}_i' \boldsymbol{\beta}, \sigma^2)$$

- El foco principal está en la distribución de la variable censurada  $Y_i$ , la cual dependerá de la variable latente  $Y_i^*$  y de su esperanza, que es una función no lineal de  $a$ ,  $\sigma$ ,  $\mathbf{x}$  y  $\boldsymbol{\beta}$

$$E(Y_i | \mathbf{X}_i) = \Phi(\alpha_i) a + [1 - \Phi(\alpha_i)] [\mathbf{x}_i' \boldsymbol{\beta} + \sigma \lambda(\alpha_i)]$$

$$\text{where } \alpha_i = \frac{a - \mathbf{x}_i' \boldsymbol{\beta}}{\sigma}$$

- En este caso, la media de  $Y_i^*$  es  $\mathbf{x}_i' \boldsymbol{\beta}$ , por lo que  $\alpha_i$  dependerá del individuo y la esperanza también
- La definición de la *ratio* inversa de Mills y las probabilidades de cada valor dependerán de la definición de la variable censurada, de modo que el resultado obtenido es generalizable a más casos
- En el caso especial en el que  $a = 0$ , se puede obtener la siguiente expresión, la cual también es generalizable:

$$\begin{aligned} E(Y_i | \mathbf{X}_i) &= \left[ 1 - \Phi\left(\frac{-\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) \right] \left[ \mathbf{x}_i' \boldsymbol{\beta} + \sigma \lambda\left(\frac{-\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) \right] \\ &= \Phi\left(\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) \left[ \mathbf{x}_i' \boldsymbol{\beta} + \sigma \lambda\left(\frac{-\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) \right] \end{aligned}$$

- A través de la esperanza condicional de  $Y_i$ , es posible obtener el efecto marginal de los regresores para una variable censurada

$$\frac{\partial E(Y_i | \mathbf{x}_i)}{\partial X_{ki}} = \beta_k P(a < Y_i^* < b | \mathbf{x}_i) = \beta_k \left[ F\left(\frac{b - \mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) - F\left(\frac{a - \mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) \right]$$

where  $F$  is the cond. prob. distrib. function of  $Y_i$

- Esta expresión es generalizable para casos en los que no se censura entre dos valores y para casos en los que la distribución no sea normal (los errores no se distribuyen normalmente)
- El efecto en  $Y_i$  (siendo censurada) de un incremento unitario de una variable  $X_{ki}$  se debe interpretar como el efecto que tiene sobre aquellos individuos que tienen valores dentro de la distribución censurada y aquellos que tienen valores fuera de ella (el efecto se interpreta para la variable latente)
- Debido a que la probabilidad está entre  $[0,1]$ , el efecto marginal de aumentar una unidad en el regresor  $X_{ki}$  es siempre menor o igual a su coeficiente en valor absoluto

$$|ME_{cens}| \leq |\beta_k|$$

- Cuando un regresor está en el modelo con una forma polinómica o con una interacción, el efecto marginal no siempre tiene el mismo signo que el primer coeficiente del regresor, debido a que tiene dos o más y depende de los signos de estos coeficientes
- McDonald y Moffitt sugieren una descomposición del efecto marginal del modelo Tobit que permite analizar el efecto de un aumento unitario en un regresor

$$\begin{aligned} \frac{\partial E(Y_i | \mathbf{x}_i)}{\partial X_{ki}} &= \beta_k [\Phi(\alpha_i)[1 - \delta(\alpha_i)] + \phi(\alpha_i)[\alpha_i + \lambda(\alpha_i)] = \\ &= P(Y_i < a | \mathbf{x}_i) \frac{\partial E(Y_i | \mathbf{x}_i; Y_i < a)}{\partial X_{ki}} + E(Y_i | \mathbf{x}_i; Y_i < a) \frac{\partial P(Y_i < a | \mathbf{x}_i)}{\partial X_{ki}} \end{aligned}$$

- La primera parte de la expresión es el efecto en el valor esperado condicional en la parte no censurada de la distribución, la cual se llama margen intensivo. La segunda parte es el efecto en la probabilidad de que la observación pertenezca a esa parte de la distribución, llamada margen extensivo

$$ME_{cens} = ME_{int} + ME_{ext}$$

$$ME_{int} = P(Y_i < a | \mathbf{x}_i) \frac{\partial E(Y_i | \mathbf{x}_i; Y_i < a)}{\partial X_{ki}}$$

$$ME_{ext} = E(Y_i | \mathbf{x}_i; Y_i < a) \frac{\partial P(Y_i < a | \mathbf{x}_i)}{\partial X_{ki}}$$

- Como el efecto marginal de un regresor en la variable censurada debe ser menor o igual en valor absoluto al coeficiente, los márgenes tienen que ser ambos menores al valor absoluto del coeficiente

$$|ME_{cens}| \leq |\beta_k| \Rightarrow |ME_{int}| < |\beta_k| \text{ and } |ME_{ext}| < |\beta_k|$$

- A través de analizar el margen intensivo y el margen extensivo, se puede ver como un incremento en el regresor tiene un efecto en la esperanza condicional y en la probabilidad de pertenecer a la distribución censurada con el mismo signo (al ser múltiplos positivos de  $\beta_k$ )

$$\frac{\partial P(Y_i < a | \mathbf{x}_i)}{\partial X_{ki}} = \left( \frac{\beta_k}{\sigma} \right) \phi(\alpha_i)$$

$$\frac{\partial E(Y_i | \mathbf{x}_i; Y_i < a)}{\partial X_{ki}} = \left( \frac{\beta_k}{\sigma} \right) [1 - \delta(\alpha_i)]$$

- La proporción que representa el margen extensivo y el margen intensivo en el efecto marginal censurado debe ser el mismo para todos los individuos, aunque el efecto marginal censurado difiera para cada uno
- Para la población de la que se obtienen los datos de la variable latente, se puede utilizar un modelo de regresión con la esperanza, llamado modelo Tobit:

$$Y_i = E(Y_i | \mathbf{x}_i) = \Phi(\alpha_i)a + [1 - \Phi(\alpha_i)][\mathbf{x}_i' \boldsymbol{\beta} + \sigma \lambda(\alpha_i)] + u_i$$

$$\text{where } u_i = Y_i - E(Y_i | \mathbf{x}_i)$$

- Por definición, el término de error  $u_i$  se distribuirá de manera normal con media nula, pero tiene una varianza equivalente a  $Var(Y_i | \mathbf{x}_i)$  y es heteroscedástico (dado que la varianza de la variable censurada depende de  $\mathbf{x}_i$ )
- Si se estima el modelo por mínimos cuadrados ordinarios, no obstante, los estimadores estarán sesgados porque se ignora el término no lineal  $\lambda(\alpha_i)$  y habría un sesgo por omisión de variables
- El modelo de Tobit tiene algunas limitaciones debido a la suposiciones que se hacen y a su especificación



- Se suelen hacer las suposiciones de normalidad y homoscedasticidad, pero estas pueden no ser apropiadas dada la población y la muestra
- El modelo Tobit propone un modelo para la variable observada  $Y_i$  en base a la variable latente  $Y_i^*$ . No obstante,  $Y_i$  depende de una variable endógena implícita, la cual se puede entender como la pertenencia a la distribución censurada o no (una variable binaria  $D_i$ ), y como no está modelada, se asume que el efecto de las  $x_i$  es el mismo para ambas, lo cual es poco realista (no se participa porque no se quiere)
- Para estimar los coeficientes mediante el método de máxima verosimilitud se tiene que maximizar la función de verosimilitud en relación a los coeficientes de los regresores y al parámetro  $\sigma$ 
  - La función de verosimilitud para una variable censurada a la derecha será la siguiente:

$$L = \prod_{i=1}^{N_1} P(Y_i = a) \prod_{i=N_1+1}^N f(Y_i) = \prod_{i=1}^{N_1} [1 - \Phi(\alpha_i)] \prod_{i=N_1+1}^N \frac{\phi\left[\frac{Y_i - \alpha_i}{\sigma}\right]}{\sigma}$$

$1, \dots, N_1 = \text{censored individuals}$

$N_1 + 1, \dots, N = \text{uncensored individuals}$

- Como en realidad hay dos grupos de individuos, se pueden dividir entre aquellos con un valor positivo de  $Y_i^*$  y aquellos con un valor de  $a$
- En el grupo de aquellos que no consumen solo se puede observar un valor discreto  $a$ , de modo que se puede utilizar la probabilidad puntual (condicionada a los valores de los regresores de cada individuo) para expresar la probabilidad de que  $Y = a$  (de que  $Y_i^* \geq a$ )

$$P(Y_i^* < a) + P(Y_i^* \geq a) = \Phi(\alpha_i) + 1 - \Phi(\alpha_i) = 1$$

$$P(Y_i = a) = P(Y_i^* \geq a) = 1 - \Phi(\alpha_i) = \Phi(-\alpha_i)$$

- En el grupo de aquellos que consumen se puede observar infinitos valores positivos continuos, de modo que se tiene que utilizar la función de densidad de probabilidad condicionada a los valores de los regresores de cada individuo

- La función de verosimilitud para una variable truncada a la derecha, en cambio, será la siguiente:

$$L = \prod_{i=1}^N \frac{f(Y_i)}{P(Y < a)} = \prod_{i=1}^N \frac{\frac{1}{\sigma} \phi \left[ \frac{Y_i^*}{\sigma} - \alpha_i \right]}{\Phi(\alpha_i)}$$

- En este caso, la muestra solo incluye observaciones para las cuales  $Y < a$ , de modo que no hay dos grupos como en el caso anterior y todos los individuos pertenecen al mismo (no es una muestra representativa de la población porque solo representa la parte de la población con valores positivos)
- Como solo se tienen valores superiores a  $a$ , la probabilidad de que se de un valor positivo cualquiera de  $Y_i$  (condicionada a los valores de los regresores de cada individuo) está condicionada a la probabilidad (condicionada) de que  $Y_i^* < a$  (por la representatividad de la muestra), de modo que se divide entre la probabilidad puntual de que  $Y_i^* < a$
- Dado que la probabilidad puntual condicionada de que  $Y_i^* < a$  sumada a la probabilidad puntual condicionada de que  $Y_i = a$  es toda la población, entonces:

$$P(Y_i^* < a) + P(Y_i^* \geq a) = \Phi(\alpha_i) + 1 - \Phi(\alpha_i) = 1$$

$$P(Y_i < a) = P(Y_i^* < a) = \Phi(\alpha_i) = 1 - \Phi(-\alpha_i)$$

- Se maximiza la función de verosimilitud en base a los coeficientes y a la  $\sigma$  con tal de poder obtener los coeficientes que maximizan la función:

$$\max_{\beta, \sigma} L = \prod_{i=1}^{N_1} [1 - \Phi(\alpha_i)] \prod_{i=N_1+1}^N \frac{\phi \left[ \frac{Y_i}{\sigma} - \alpha_i \right]}{\sigma}$$

$$\max_{\beta, \sigma} L = \prod_{i=1}^N \frac{\frac{1}{\sigma} \phi \left[ \frac{Y_i}{\sigma} - \alpha_i \right]}{\Phi(\alpha_i)}$$

- En este caso, el sistema de ecuaciones con las condiciones de primer orden para todos los coeficientes también tiene una ecuación de primer orden para la desviación típica  $\sigma$ , de modo que también se puede estimar

$$\beta, \sigma \Rightarrow \hat{\beta}, \hat{\sigma}$$

- Los estimadores obtenidos por máxima verosimilitud solo serán eficientes si el error se distribuye de manera normal y los errores son homoscedásticos (además de si la especificación del modelo es buena)
- La función de verosimilitud para una variable censurada y truncada a la vez tendrá que dividir entre la probabilidad que representa la muestra, pero teniendo en cuenta la probabilidad de ambos valores de  $Y$ :

$$L = \prod_{i=1}^{N_1} \frac{P(a < Y^* < b)}{P(a < Y^* < c)} \prod_{i=N_1+1}^N \frac{f(Y^*)}{P(a < Y^* < c)}$$

- Para estimar los coeficientes de ambos modelos sin sesgos, se puede usar un método de dos etapas en donde se hace una primera estimación con un modelo *probit* (primera etapa) y una segunda estimación (segunda etapa) mediante mínimos cuadrados ordinarios
  - Los individuos se pueden dividir en dos grupos cuando se trunca o se censura una variable: los que tienen valores  $Y_i^* > a$  y aquellos que no. Por lo tanto, se puede crear una variable binaria para usarse como dependiente en un modelo *probit*

$$D_i = \begin{cases} 1 & \text{if } Y_i^* > a \\ 0 & \text{if } Y_i^* \leq a \end{cases} \Rightarrow d_i^* = \mathbf{x}_i' \boldsymbol{\beta} + u_i \quad \text{where } u_i \sim N(0,1)$$

$$\Rightarrow P(D_i = 1 | \mathbf{x}_i) = \Phi(\alpha_i)$$

- Estos coeficientes son diferentes de las  $\boldsymbol{\beta}$  y se estiman por máxima verosimilitud
- Se escoge una definición de  $D_i$  u otra dependiendo de si la truncación o la censura es por arriba, por abajo, por la derecha, por la izquierda o por ambas
- Como la probabilidad de que  $D_i = 1$  es igual a la probabilidad de que  $d_i^* > a$ , se puede ver que existe una relación entre los coeficientes del modelo de truncación o censura y los del modelo *probit* si el modelo es correcto

$$P(d_i^* > a) = 1 - \Phi(a - \mathbf{x}_i' \boldsymbol{\beta}) \Rightarrow \Phi(\mathbf{x}_i' \boldsymbol{\gamma}) = \Phi(\mathbf{x}_i' \boldsymbol{\beta} - a)$$

- Se puede ver como hay una relación de proporcionalidad entre los coeficientes  $\boldsymbol{\beta}$  y  $\boldsymbol{\gamma}$  de los regresores no constantes: los coeficientes del modelo no son más que los coeficientes del *probit* multiplicados por  $\sigma$

$$\gamma = \frac{\beta}{\sigma} \Rightarrow \beta = \gamma\sigma$$

- Parar el coeficiente constante, la relación será la misma, pero se le añadirá o se le sustraerá  $a$  dependiendo de la especificación de la variable dependiente observada

$$\gamma_0 = \frac{\beta_0 \pm a}{\sigma} \Rightarrow \beta = \gamma\sigma \pm a$$

- La segunda etapa para una variable truncada consiste en hacer una regresión lineal de la variable dependiente  $Y_i$  sobre los regresores, de modo que se estiman los coeficientes de dos etapas por mínimos cuadrados

$$\hat{Y}_i = E(Y_i | x_i; D_i = 1) = x_i' \hat{\beta} + \hat{\sigma} \lambda(\hat{\alpha}_i)$$

$$\text{By MCO: } \hat{\beta}^{2stage}, \hat{\sigma}^{2stage}$$

- El valor de  $\lambda(\hat{\alpha}_i)$  (que funciona como regresor) se obtiene con los coeficientes de la primera parte
- La segunda etapa para una muestra censurada consiste en hacer una regresión lineal de la variable dependiente  $Y_i$  sobre los regresores multiplicados por  $\Phi(\hat{\alpha}_i)$ , de modo que se estiman los coeficientes de dos etapas por mínimos cuadrados

$$\hat{Y}_i = E(Y_i | x_i) = \Phi(\hat{\alpha}_i) a + [1 - \Phi(\hat{\alpha}_i)] [x_i' \hat{\beta} + \hat{\sigma} \lambda(\hat{\alpha}_i)]$$

$$\text{By MCO: } \hat{\beta}^{2stage}, \hat{\sigma}^{2stage}$$

- El valor de  $\Phi(\hat{\alpha}_i)$  y  $\phi(\hat{\alpha}_i)$  se obtienen con los coeficientes de la primera parte, y  $\phi(\hat{\alpha}_i)$  funciona como un regresor más
- Para estimar por mínimos cuadrados, sin embargo, hace falta multiplicar los valores de  $\Phi(\hat{\alpha}_i)$  por los de  $x$  (creando así nuevos regresores) para obtener los coeficientes de  $x$ . El coeficiente constante, por tanto, ahora depende de un regresor que no es  $X_{0i} = 1$ , sino que es  $\Phi(\hat{\alpha}_i)$ , y por tanto es como si no hubiera un coeficiente constante
- Los estimadores obtenidos por este método no requieren que el error se distribuya normal y/o que los errores sean homoscedásticos, pero si se cumplen ambas condiciones, se obtendrán estimadores menos eficientes de los que se obtendrían por máxima verosimilitud

## Los modelos de dos ecuaciones

- Al usar el modelo Tobit, una de las suposiciones que se hacían es que los individuos con un valor de  $Y_i = a$  (la solución de esquina) eran individuos en los cuales tenían un valor de  $Y_i^* > a$  para la variable latente (dependiendo de la definición de la variable observada)
  - Lo que esto significa esta suposición es que el efecto de los regresores  $x_i$  en la probabilidad de que  $Y_i^* > a$  es el mismo que en el valor esperado de la variable  $Y_i^*$  cuando  $Y_i^* > a$  (dado que solo se tienen unos coeficientes  $\beta$  para ambas etapas)
    - Esta restricción es poco realista, dado que el efecto de las variables no siempre es el mismo para ambas etapas, pero el modelo considera que los coeficientes son los mismos
  - Para resolver este problema, se pueden plantear modelos en dos partes o *hurdle models* en los cuales se separan ambas etapas: una para la probabilidad de que  $Y_i^* > a$ , y otra para el valor esperado de  $Y_i^*$  cuando  $Y_i^* > a$ 
    - Se separa la probabilidad de pertenecer a la distribución censurada  $P(Y_i^* > a)$  del valor que se tiene dentro de esta distribución  $E(Y_i^* | x_i; Y_i^* > a)$
    - Aunque no se usen los mismos regresores ni el mismo efecto para cada decisión o parte (al tener diferentes coeficientes), puede haber correlación entre los errores, lo cual haría que las variables se determinaran de manera endógena
- Para poder lidiar con los regresores y la posible correlación entre términos de errores, se pueden plantear modelos aún más generales, como el modelo de Heckman, el *double hurdle model* de Cragg y el *two part model* (2P)
  - La ecuación de participación en los modelos depende de una variable latente  $d_i^*$  que determina una variable binaria  $D_i$

$$d_i^* = \mathbf{z}_i' \boldsymbol{\gamma} + u_i \quad \text{and} \quad D_i = \begin{cases} 1 & \text{if } d_i^* > a \\ 0 & \text{if } d_i^* \leq a \end{cases}$$

- El término de error de la variable latente se distribuye de manera normal estándar (aunque se puede asumir una distribución normal para después estandarizar con la desviación típica, es normal usar la normal estándar)

$$u_i \sim N(0,1)$$

- La ecuación de intensidad en los modelos es general y sirve tanto para las observaciones censuradas como para las no censuradas, de modo que se modela una variable latente  $Y_i^*$

$$Y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i \quad \varepsilon_i \sim N(0, \sigma^2)$$

- El término de error en esta ecuación se distribuye de manera normal
- Existen dos mecanismos de observación: uno en el que se asume que  $Y_i^* = a$  cuando  $D_i = 0$  y/o  $Y_i^* \leq a$ , y otro en el que se asume que  $Y_i$  no se observa cuando  $D_i = 0$

(Cragg & 2P)  $Y_i = Y_i^*$  if  $D_i = 1$  &  $Y_i^* > a$  and  $Y_i = a$  otherwise

(Heckman)  $Y_i = Y_i^*$  if  $D_i = 1$  and  $Y_i$  is unobs. if  $D_i = 0$

- El primer mecanismo de observación permite obtener el *double hurdle model* y el *two part model*, mientras que el segundo permite obtener el modelo de Heckman. El modelo de Heckman es un modelo de selección muestral porque permite estimar cuando hay un problema de selección muestral, mientras que el *double-hurdle* o el *two-part model* no lo son porque tienen en cuenta las observaciones  $Y_i = a$
- En el caso del modelo de dos partes y de Cragg, se asume que se observan los regresores  $\mathbf{x}_i$  y que estos hacen que  $Y_i = a$  (dado que se pone la condición para cualquier otra circunstancia) para los participantes. Esto hace que, si  $\mathbf{x}_i$  se observa, las observaciones  $Y_i = a$  puedan contribuir a la verosimilitud para la muestra completa y que haya dos fuentes para las cuales  $Y_i = a$
- No obstante, en el modelo de Cragg no se puede observar  $D_i$ , mientras que en el de dos partes sí. Esto hace que los valores  $Y_i = a$  en el primer modelo puedan darse por no participar o como una solución de esquina (no se sabe exactamente al no observarse), pero que en el segundo modelo se pueda saber si proviene de uno o de otro
- En el modelo de Heckman se asume que no se observan los regresores  $\mathbf{X}_i$ , por lo que no se puede observar una  $Y_i^*$  ni determinar una  $Y_i$  para  $D_i = 0$ . Esto hace que, si  $\mathbf{X}_i$  no se observa, los individuos que no participan no contribuyan a la verosimilitud

- Los errores pueden estar correlacionados, por lo que se pueden determinar las variables de manera endógena y los errores seguirán una distribución normal bivalente con correlación  $\rho$

$$(u_i, \varepsilon_i) \sim N(0,0,1, \sigma^2, \rho)$$

- La esperanza de  $Y_i$  para los individuos que están dentro de la muestra censurada permite obtener una expresión en la que se ve como se incluye la posibilidad de correlación y endogeneidad en los modelos

$$E(Y_i | \mathbf{x}_i, \mathbf{z}_i; D_i = 1) = \mathbf{x}_i' \boldsymbol{\beta} + E(\varepsilon_i | \mathbf{x}_i, \mathbf{z}_i; D_i = 1)$$

$$\Rightarrow E(Y_i | \mathbf{x}_i, \mathbf{z}_i; D_i = 1) = \mathbf{x}_i' \boldsymbol{\beta} + E(\varepsilon_i | \mathbf{x}_i, \mathbf{z}_i; u_i > a - \mathbf{z}_i' \boldsymbol{\gamma})$$

$$\Rightarrow E(Y_i | \mathbf{x}_i, \mathbf{z}_i; D_i = 1) = \mathbf{x}_i' \boldsymbol{\beta} + \rho \sigma \lambda(a - \mathbf{z}_i' \boldsymbol{\gamma})$$

$$\Rightarrow E(Y_i | \mathbf{x}_i, \mathbf{z}_i; D_i = 1) = \mathbf{x}_i' \boldsymbol{\beta} + \rho \sigma \frac{\phi(a - \mathbf{z}_i' \boldsymbol{\gamma})}{1 - \Phi(a - \mathbf{z}_i' \boldsymbol{\gamma})}$$

$$\Rightarrow Y_i = \mathbf{x}_i' \boldsymbol{\beta} + \rho \sigma \frac{\phi(a - \mathbf{z}_i' \boldsymbol{\gamma})}{1 - \Phi(a - \mathbf{z}_i' \boldsymbol{\gamma})} + \eta_i$$

- El último término (la *ratio* inversa de Mills) puede variar de especificación dependiendo de la definición de la variable  $D_i$ , y se utiliza  $\mathbf{z}_i' \boldsymbol{\gamma}$  porque la variable ya se distribuye de manera normal estándar
- Para el modelo de dos partes, se pueden utilizar formas funcionales tales como una log-normal o una truncada normal
- A partir de los efectos marginales, se puede comprobar que el signo del efecto marginal de un regresor (ya sea  $\mathbf{x}_i$  o  $\mathbf{z}_i$ ) depende del signo del coeficiente. Por lo tanto, se puede interpretar el signo del coeficiente, pero no la magnitud
  - No obstante, los coeficientes  $\boldsymbol{\beta}$  son el efecto marginal latente (de la variable latente  $Y_i^*$ ), por lo que son directamente interpretables. De este modo, se puede interpretar tanto la magnitud como el signo (dado que es una regresión lineal)

$$ME_k^{lat} = \frac{\partial E(Y_i^* | \mathbf{x}_i)}{\partial X_{ki}} = \beta_k$$

- Si  $\mathbf{X}_i$  y  $\mathbf{Z}_i$ , y son diferentes, el efecto marginal de incrementar una unidad en el regresor  $X_{ik}$  es el coeficiente  $\beta_k$ , de modo que

el efecto marginal latente coincide con el efecto marginal de la variable dependiente observable. No obstante,  $\mathbf{z}_i$  no afecta

$$ME_k = \frac{\partial E(Y|\mathbf{x}_i; D_i = 1)}{\partial X_{ki}} = \beta_k \quad \text{only if } \mathbf{x} \neq \mathbf{z}$$

- En el modelo de Heckman, el efecto en  $Y_i$  (siendo censurada) de un incremento unitario de una variable  $X_{ki}$  se debe interpretar como el efecto que tiene sobre aquellos individuos que tienen valores dentro de la distribución truncada, pero no para aquellos que la tienen fuera, dado que no se observan valores para  $D_i = 0$
- En el *double hurdle model* y en el *two-part model*, el efecto en  $Y_i$  (siendo censurada) de un incremento unitario de una variable  $X_{ki}$  se debe interpretar como el efecto que tiene sobre aquellos individuos que tienen valores dentro de la distribución censurada y fuera
- A partir de la definición de los modelos, se puede ver como la función de densidad de la variable censurada  $Y_i$  es una combinación entre funciones continuas y discretas
- Para el *double hurdle model* de Cragg, la función de verosimilitud no tendrá en cuenta  $D_i = 0$  porque no se observa, de modo que solo tiene en cuenta las condiciones de las observaciones (que  $D_i = 1$ )

$$L = \prod_{Y_i=a} [1 - P(D_i = 1; Y_i^* > a)] \prod_{Y_i < a} P(D_i = 1; Y_i^* = Y_i)$$

$$1 - P(D_i = 1; Y_i^* > a) = 1 - P(u_i > a - \mathbf{z}_i'\boldsymbol{\gamma}; \varepsilon_i > a - \mathbf{x}_i'\boldsymbol{\beta}) =$$

$$= 1 - N(\mathbf{z}_i'\boldsymbol{\gamma} - a, \mathbf{x}_i'\boldsymbol{\beta} - a, 1, \sigma, \rho)$$

$$P(D_i = 1; Y_i^* = Y_i) = P(u_i > a - \mathbf{z}_i'\boldsymbol{\gamma}; \varepsilon_i = Y_i - \mathbf{x}_i'\boldsymbol{\beta}) =$$

$$= \int_{-\mathbf{z}_i'\boldsymbol{\gamma}}^{\infty} N(u, \mathbf{x}_i'\boldsymbol{\beta} - a, 1, \sigma, \rho) du =$$

$$= \frac{\phi\left(\frac{Y_i - \mathbf{x}_i'\boldsymbol{\beta}}{\sigma}\right)}{\sigma} \left[ 1 - \Phi\left(\frac{-\mathbf{z}_i'\boldsymbol{\gamma} - \rho\left(\frac{Y_i - \mathbf{x}_i'\boldsymbol{\beta}}{\sigma}\right)}{\sqrt{1 - \rho^2}}\right) \right]$$



- Para el modelo de dos partes, la función de verosimilitud si tendrá en cuenta  $D_i = 0$  y por tanto lo puede utilizar para la estimación:

$$f(Y_i|\mathbf{z}_i) = \begin{cases} P(D_i = 1|\mathbf{z}_i)f(Y_i|\mathbf{x}_i; D_i = 1) & \text{when } Y_i > a \\ P(D_i = 0|\mathbf{z}_i) & \text{when } Y_i = a \end{cases}$$

$$f(Y_i|\mathbf{x}_i; D_i = 1) = \text{density of the rand. var. when } < a$$

$$P(D_i = 1 \text{ or } 0|\mathbf{z}_i) = \text{probit or logit specification}$$

- Para el modelo de Heckman, la función de verosimilitud será la siguiente (todo condicionado a los regresores correspondientes):

$$L = \prod_{Y_i=a} P(D_i = 0) \prod_{Y_i < a} P(D_i = 1; Y_i^* = Y_i)$$

$$P(D_i = 0) = P(u_i \leq a - \mathbf{z}_i'\boldsymbol{\gamma}) = 1 - \Phi(a - \mathbf{z}_i'\boldsymbol{\gamma})$$

$$P(D_i = 1; Y_i^* = Y_i) = P(u_i > a - \mathbf{z}_i'\boldsymbol{\gamma}; \varepsilon_i = Y_i - \mathbf{x}_i'\boldsymbol{\beta}) =$$

$$= \int_{-\mathbf{z}_i'\boldsymbol{\gamma}}^{\infty} N(u, \mathbf{x}_i'\boldsymbol{\beta} - a, 1, \sigma, \rho) du =$$

$$= \frac{\phi\left(\frac{Y_i - \mathbf{x}_i'\boldsymbol{\beta}}{\sigma}\right)}{\sigma} \left[ 1 - \Phi\left(\frac{-\mathbf{z}_i'\boldsymbol{\gamma} - \rho\left(\frac{Y_i - \mathbf{x}_i'\boldsymbol{\beta}}{\sigma}\right)}{\sqrt{1 - \rho^2}}\right) \right]$$

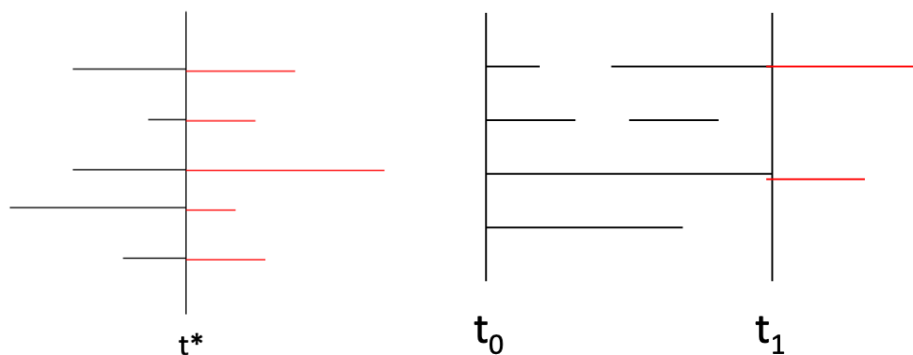
- Cuando la correlación es nula, entonces se usa la probabilidad de dos variables independientes (se multiplican las probabilidades para cada variable y no se usa la distribución bivalente)

## Los modelos de duración

- Un tipo de modelos en donde se utilizan variables censuradas son los modelos de duración. En estos modelos se considera el tiempo hasta que algún tipo de transición como la duración, y la transición misma como un evento
  - El rol que juega la censura en este tipo de modelos es que en casi todos los casos en los que se estudia datos de duración, algunos o muchos de los *spells* que se observan no acaban en transiciones
  - El término *spell* se refiere a las diferentes variables de duración que se pueden medir, y los datos de estos *spells* están, en efecto,

censurados, por lo que los modelos de duración tienen esto en cuenta explícitamente

- Los aspectos más interesantes a analizar con estos modelos son aspectos relacionados con la probabilidad y el tiempo hasta que ocurre una transición
  - Aunque es interesante la duración de un evento, también es interesante la probabilidad de que un evento realice una transición en un periodo futuro condicional a la duración que ha tenido hasta ahora
  - El análisis de tiempo hasta el fallo o tiempo de supervivencia (*time until failure* o *survival time*), definido como el tiempo que un evento desde que comienza hasta que ocurre una transición, también es interesante de analizar
- La variable de interés es la duración de un evento desde que comienza hasta que termina o se mide, de modo que la medida puede preceder la terminación
  - Las observaciones de las duraciones suelen tener una forma de datos transversales o *cross-sectional data*, o con una forma de datos longitudinales o *longitudinal data*



- Como se puede ver, la duración es el valor de la variable desde que comienza el proceso (puede haber diferentes comienzos para diferentes individuos) hasta  $t^*$  para datos transversales, mientras que la duración es el valor de la variable desde  $t_0$  hasta  $t_1$  (hay un solo comienzo para cada individuo)
- Normalmente se representa el largo de estos *spells* con una variable  $T$ , mientras que la duración realizada (observada en los datos) se representa con  $t$
- La censura es un problema general que comúnmente es debido a que se suele medir cuando el proceso aún sigue, de modo que las

observaciones de individuos cuyo evento sigue en marcha están censuradas necesariamente

- En consecuencia, la duración o supervivencia de los individuos en una muestra para datos transversales es como mínimo  $T \geq t^*$ , pero no igual a  $t^*$  (dado que no se suele observar cuando hay una transición), por lo que se considera que todos los *spells* están incompletos y censurados. Para los datos longitudinales, en cambio, pueden haber *spells* completos e incompletos, de modo que no hay una duración mínima
- La estimación de este tipo de modelos tiene que tener en cuenta esta naturaleza censurada de los datos por el efecto de esta censura en la media y la varianza y en la inferencia (sesgo por omisión de variables)
- Representando la duración (el *spell*) como una variable aleatoria  $T$ , hay aspectos de su distribución de probabilidad que permitirán estudiar y hacer un análisis regresivo a partir de modelos

- Suponiendo que la variable  $T$  tiene una función de densidad  $f(t)$ , se puede obtener la función de distribución de probabilidad acumulada y la función de supervivencia

- Integrando la función de densidad de la variable de 0 a  $t$ , se puede obtener la función de probabilidad acumulada

$$F(t) = \int_0^t f(s) ds = P(T \leq t)$$

- La función de supervivencia o *survival function* se define como la probabilidad de que un *spell* sea de una duración de al menos  $t$ , por lo que es la probabilidad a la derecha de  $t$

$$S(t) = 1 - \int_0^t f(s) ds = 1 - P(T \leq t) = P(T \geq t)$$

- Para poder saber cuál es la probabilidad de que un *spell* finalice en un momento futuro, se utiliza la función de distribución y la noción de tasa de riesgo

- La probabilidad de que un *spell* finalice en un momento futuro  $t + \Delta t$  dado que ha durado hasta el momento  $t$  (como mínimo) es la siguiente:

$$l(t, t + \Delta t) = P(t \leq T \leq t + \Delta t \mid T \geq t) = \frac{F(t + \Delta t) - F(t)}{S(t)}$$

- Una función útil para caracterizar este aspecto de la distribución de probabilidad es la función de riesgo o *hazard function*, que se entiende como la velocidad a la que los *spells* sufren una transición después de durar  $t$ , dado que duran como mínimo  $t$ . Se puede demostrar, además, que la tasa es la semielasticidad negativa de la función de supervivencia, por lo que la función de densidad es el producto entre  $S(t)$  y  $\lambda(t)$

$$\lambda(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T \leq t + \Delta t \mid T \geq t)}{\Delta t} = \lim_{\Delta t \rightarrow 0} \frac{F(t + \Delta t) - F(t)}{\Delta t S(t)} = \frac{f(t)}{S(t)}$$

$$\lambda(t) = -\frac{d \ln S(t)}{dt} = -\frac{d \ln S(t)}{dS(t)} \frac{dS(t)}{t} = -\frac{-f(t)}{S(t)} = \frac{f(t)}{S(t)}$$

$$\Rightarrow f(t) = \lambda(t)S(t)$$

- Una función útil que también se usa en los modelos es la función de riesgo integrada, la cual permite obtener una expresión alternativa para la función de supervivencia

$$\Lambda(t) = \int_0^t \lambda(s) ds \Rightarrow \Lambda(t) = \int_0^t -\frac{d \ln S(s)}{ds} ds = -\ln S(t)$$

$$\Rightarrow S(t) = e^{-\Lambda(t)}$$

- Además, la derivada de la función de riesgo permite saber cómo los *spells* dependen de la duración  $t$ , de modo que, si es positiva, más duración realizada (una  $t$  mayor) significa que la probabilidad de que finalice el *spell* es mayor, mientras que una negativa significa lo opuesto

$$\frac{d\lambda(t)}{dt} > 0 \text{ (positive dependence)}$$

$$\frac{d\lambda(t)}{dt} < 0 \text{ (negative dependence)}$$

- Dependiendo de si los datos son transversales o longitudinales, la función de máxima verosimilitud varía

- Para los datos transversales, la función de máxima verosimilitud y su logaritmo vienen dadas por las siguientes expresiones:

$$L = \prod_{i=1}^N \frac{S(t)}{E(T)} = \prod_{i=1}^{N_1} \frac{1 - F(t)}{\int_0^\infty [1 - F(t)] dt}$$

$$\Rightarrow \ln L = \sum_{i=1}^{N_1} \ln S(t) + \sum_{i=N_1+1}^N \ln E(t)$$

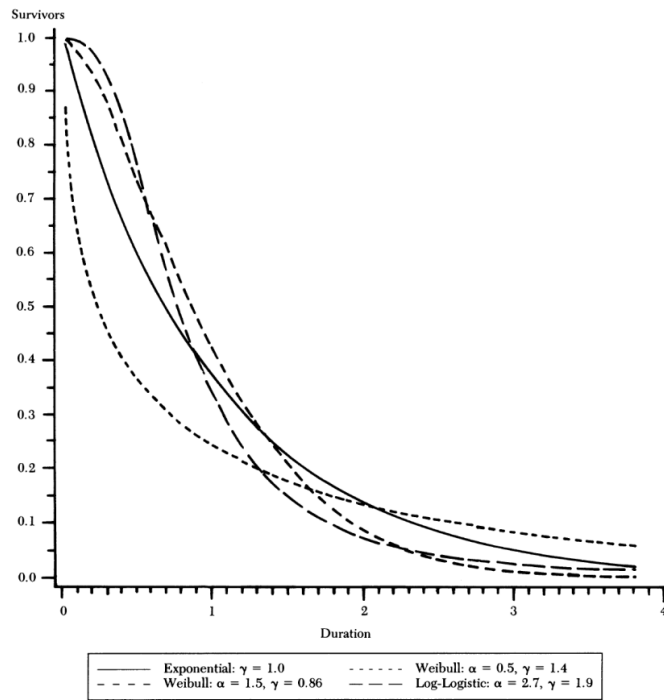
- Para los datos longitudinales, la función de máxima verosimilitud y su logaritmo vienen dadas por las siguientes expresiones:

$$L = \prod_{i=1}^{N_1} P(T_i = t_i) \prod_{i=N_1+1}^N P(T_i > t_i) = \prod_{i=1}^{N_1} f(t) \prod_{i=N_1+1}^N S(t)$$

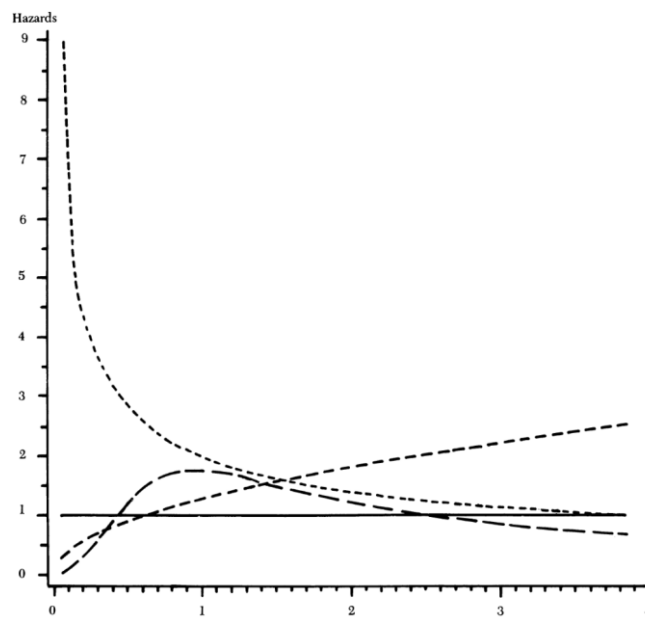
$$\Rightarrow \ln L = \sum_{i=1}^{N_1} \ln f(t) + \sum_{i=N_1+1}^N \ln S(t)$$

where  $1, \dots, N_1$  are completed and  $N_1 + 1, \dots, N$  are incompleting spells

- Representando la duración (el *spell*) como una variable aleatoria  $T$ , se puede aplicar un análisis regresivo para la muestra de *spells* observados, por lo que se puede caracterizar la duración esperada condicionada a unos regresores y la distribución de probabilidad
  - La función de riesgo es más útil en el análisis que no la función de supervivencia o la densidad, por lo que es mejor modelar la función de riesgo
    - Como el interés está en si hay más o menos probabilidad de que un *spell* termine en el momento  $t$  dado que ha durado hasta el momento  $t$ , la dependencia de la duración es lo más relevante y eso implica que también lo es la función de riesgo
    - La elección del modelo dependerá de cómo se ajusten los datos a cada uno de ellos. Mirando la función de supervivencia, se puede ver que los modelos principales son muy similares, y, por tanto, no es muy útil modelar esta función



- No obstante, la función de riesgo para cada modelo exhibe un comportamiento diferenciado, de modo que es útil modelar esta función con tal de ajustarse a los datos



- Los modelos más importantes y más usados para el análisis de regresión son el modelo exponencial, el modelo Weibull, el modelo log-normal y el modelo log-logístico

<b>Distribution</b>	<b>Hazard Function, <math>\lambda(t)</math></b>	<b>Survival Function, <math>S(t)</math></b>
Exponential	$\lambda,$	$S(t) = e^{-\lambda t}$
Weibull	$\lambda p(\lambda t)^{p-1},$	$S(t) = e^{-(\lambda t)^p}$
Lognormal	$f(t) = (p/t)\phi[p \ln(\lambda t)]$ [ln $t$ is normally distributed with mean $-\ln \lambda$ and standard deviation $1/p$ .]	$S(t) = \Phi[-p \ln(\lambda t)]$
Loglogistic	$\lambda(t) = \lambda p(\lambda t)^{p-1}/[1 + (\lambda t)^p],$ [ln $t$ has a logistic distribution with mean $-\ln \lambda$ and variance $\pi^2/(3p^2)$ .]	$S(t) = 1/[1 + (\lambda t)^p]$

- El modelo exponencial es un modelo en el que se asume que la función de riesgo es constante, por lo que la dependencia de la duración es nula

- En este caso, se asume que la función de riesgo es una constante  $\lambda$ , lo cual implica que la función integrada será una función del tiempo y que la dependencia con la duración es nula, por lo que la probabilidad no se altera por la duración

$$\lambda(t) = \gamma \quad \Lambda(t) = \gamma t \quad \frac{d\lambda(t)}{dt} = 0$$

- A partir de la función de riesgo, se puede ver que la función de supervivencia y las funciones de distribución y densidad son las siguientes:

$$S(t) = e^{-\gamma t} \quad F(t) = 1 - e^{-\gamma t} \quad f(t) = \lambda e^{-\gamma t}$$

- Adicionalmente, se puede demostrar que los momentos de la distribución de probabilidad de  $T$  son los siguientes:

$$E(T) = \gamma^{-1} \quad Var(T) = \gamma^{-2}$$

- El modelo Weibull es un modelo similar al exponencial, pero en el que la función de riesgo no es constante, haciendo que exista una dependencia de la duración no nula

- En este caso, se asume que la función de riesgo depende de la duración realizada o tiempo  $t$  y de un parámetro  $\alpha$

$$\lambda(t) = \gamma \alpha t^{\alpha-1} \quad \Lambda(t) = \gamma t^\alpha \quad \frac{d\lambda(t)}{dt} = \gamma \alpha (\alpha - 1) t^{\alpha-2}$$

- La dependencia de la duración dependerá de estos parámetros. Si  $\alpha > 1$ , la dependencia es positiva y la probabilidad de falla aumenta cuanto más haya durado el *spell*, mientras que se da lo opuesto si  $\alpha < 1$ . Si  $\alpha = 1$ , entonces el modelo es equivalente al exponencial y no hay dependencia de la duración

$$\frac{d\lambda(t)}{dt} = \gamma\alpha(\alpha - 1)t^{\alpha-2} \Rightarrow \begin{cases} \frac{d\lambda(t)}{dt} > 0 \text{ if } \alpha > 1 \\ \frac{d\lambda(t)}{dt} = 0 \text{ if } \alpha = 1 \\ \frac{d\lambda(t)}{dt} < 0 \text{ if } \alpha < 1 \end{cases}$$

- A partir de la función de riesgo, se puede ver que la función de supervivencia y las funciones de distribución y densidad son las siguientes:

$$S(t) = e^{-\gamma t^\alpha} \quad F(t) = 1 - e^{-\gamma t^\alpha} \quad f(t) = \gamma\alpha t^{\alpha-1} e^{-\gamma t^\alpha}$$

- Adicionalmente, se puede demostrar que la esperanza de la distribución de probabilidad de  $T$  es la siguiente:

$$E(T) = \gamma^{-\frac{1}{\alpha}} \Gamma\left(1 + \frac{1}{\alpha}\right)$$

- El modelo log-logístico es un modelo similar al Weibull, pero en el que la integral de la función de riesgo es un logaritmo, lo cual provoca que la dependencia de la duración no sea constante y pueda variar con la duración

- En este caso, se asume que la función de riesgo depende de la duración realizada o tiempo  $t$  y de un parámetro  $\alpha$

$$\lambda(t) = \frac{\gamma\alpha t^{\alpha-1}}{1 + \gamma t^\alpha} \quad \Lambda(t) = \ln(1 + \gamma t^\alpha)$$

- La dependencia de la duración dependerá de estos parámetros. Si  $\alpha > 1$ , la dependencia es positiva en un principio y después es negativa, haciendo que la probabilidad de falla aumente y después se reduzca, mientras que la dependencia será negativa si  $\alpha \leq 1$ . La duración  $t^{max}$  en la que cambia el signo de la dependencia se puede encontrar a partir de su expresión

$$\begin{aligned} \frac{d\lambda(t)}{dt} &= \frac{\gamma\alpha(\alpha - 1)t^{\alpha-2}}{1 + \gamma t^\alpha} - \frac{(\gamma\alpha t^{\alpha-1})^2}{(1 + \gamma t^\alpha)^2} \\ \Rightarrow \begin{cases} \frac{d\lambda(t)}{dt} > 0 \text{ and } \frac{d\lambda(t)}{dt} < 0 \text{ after if } \alpha > 1 \\ \frac{d\lambda(t)}{dt} < 0 \text{ if } \alpha \leq 1 \end{cases} \end{aligned}$$

$$\frac{d\lambda(t)}{dt} = 0 \Rightarrow t^{max} = \left(\frac{\alpha - 1}{\gamma}\right)^{\frac{1}{\alpha}}$$



- A partir de la función de riesgo, se puede ver que la función de supervivencia y las funciones de distribución y densidad son las siguientes:

$$S(t) = \frac{1}{1 + \gamma t^\alpha} \quad F(t) = 1 - \frac{1}{1 + \gamma t^\alpha} \quad f(t) = \frac{\gamma \alpha t^{\alpha-1}}{(1 + \gamma t^\alpha)^2}$$

- Adicionalmente, se puede demostrar que la esperanza de la distribución de probabilidad de  $T$  es la siguiente:

$$E(T) = \frac{\gamma^{-1}}{\sin\left(\frac{\pi}{\alpha}\right)}$$

- Dos limitaciones de los modelos anteriores es que no hay lugar para factores externos y que, aunque se incluyan, puede haber diferencias individuales en las distribuciones remanentes por una especificación incompleta. Estas se pueden superar a través de añadir regresores y de tener en cuenta esta heterogeneidad
  - La adición de regresores a los modelos de duración es bastante simple, aunque la interpretación de los coeficientes no lo es tanto. Para introducir la heterogeneidad, se asume que el parámetro  $\gamma$  es una función que depende de  $x_i$

$$\gamma_i = e^{x_i' \beta}$$

- Normalmente se asume que estos regresores  $x_i$  tienen un valor constante. Para datos transversales, se puede asumir que se escoge el valor de  $x_i$  en  $t^*$ , mientras que, para datos longitudinales, se puede asumir que los regresores toman un valor fijo de  $t_0$  a  $t_1$
- Además, se suele asumir que la duración realizada  $t$  sigue una distribución normal, aunque esta suposición no suele ser la mejor y se pueden asumir otro tipo de distribuciones
- En este caso, la función de riesgo no depende del tiempo porque se asume que los regresores son constantes, pero sí depende del individuo
- Como hacer que  $\lambda_i$  sea una función de sus regresores es equivalente a realizar un cambio de unidades en el eje del tiempo o de las ordenadas, estos modelos se suelen llamar modelos de duración acelerada

- El efecto de un cambio en  $X_{ki}$  a partir de un modelo exponencial o Weibull permite ver como el signo es el mismo que el del coeficiente, y se suele expresar a través de la semielasticidad para una mejor interpretación

$$(exponential) \quad \frac{d\lambda_i}{dX_{ki}} = \beta_k e^{x_i' \beta} \Rightarrow \frac{d \ln \lambda_i}{dX_{ki}} = \beta_k$$

$$(Weibull) \quad \frac{d\lambda_i}{dX_{ki}} = \beta_k e^{x_i' \beta} \alpha t^{\alpha-1} \Rightarrow \frac{d \ln \lambda_i}{dX_{ki}} = \beta_k$$

- Dado que para el modelo exponencial y Weibull la semielasticidad es el coeficiente del regresor, se puede interpretar el efecto de aumentar un regresor como en una regresión con logaritmos en la variable dependiente
- El efecto de un incremento marginal en un regresor para el modelo log-logístico y la semielasticidad de este es más complejo que para los otros modelos

$$\frac{d\lambda_i}{dX_{ki}} = \beta_k e^{x_i' \beta} \frac{\alpha t^{\alpha-1}}{(1 + e^{x_i' \beta t^\alpha})^2} \Rightarrow \frac{d \ln \lambda_i}{dX_{ki}} = \frac{\beta_k}{1 + e^{x_i' \beta t^\alpha}}$$

$$\frac{d\lambda_i}{dX_{ki}} = 0 \Rightarrow \beta_k e^{x_i' \beta} \frac{\alpha t^{\alpha-1}}{(1 + e^{x_i' \beta t^\alpha})^2} = 0 \Rightarrow \beta_k e^{x_i' \beta} \alpha t^{\alpha-1}$$

- En este caso, el signo de la dependencia será el mismo al de los coeficientes, de modo que aumentar en una unidad un regresor  $X_{ki}$  hace que la probabilidad de que el *spell* finalice en un momento futuro aumente, y viceversa
- Otra forma de expresar el efecto en la probabilidad de falla de los regresores es a través de la función de riesgo proporcional, definida como la *ratio* de funciones de riesgo

$$\frac{\lambda_i}{\lambda_j} = \frac{e^{x_i' \beta}}{e^{x_j' \beta}} = e^{(x_i - x_j)' \beta}$$

- El efecto de aumentar una unidad en un regresor  $X_{ki}$  es la reducción en  $\lambda_i$  con respecto a  $\lambda_j$
- Tanto el modelo exponencial como el Weibull tienen una función de riesgo proporcional, mientras que el modelo log-logístico no. Esto hace que, para los dos primeros modelos, se pueda interpretar el exponencial del coeficiente  $e^\beta$  como el

cambio en  $\lambda$  al variar una unidad de un regresor (como *log-level*). No obstante, esto no aplica al modelo log-logístico

- Para lidiar con la heterogeneidad que no se observa a partir de los datos individuales, se suele incluir una variable aleatoria  $v_i$  que siga una distribución en la función de riesgo multiplicando

- De este modo, esta variable aporta aleatoriedad que no es observada en la muestra en la función de riesgo. No obstante, su introducción hace que a veces se den resultados contradictorios al sentido común

- Asumiendo un modelo Weibull, se puede ver el efecto en la función de supervivencia

$$\lambda_i(t|v_i) = \gamma_i \alpha t^{\alpha-1} v_i \Rightarrow S_i(t|v_i) = e^{-\int_0^t \lambda_i(s|v_i) ds} = e^{-\gamma_i t^\alpha v_i}$$

- La distribución comúnmente asumida es una distribución gamma con umbral  $\theta$  y factor de forma 1

$$v_i \sim \Gamma(1, \theta)$$

- Los modelos paramétricos pueden distorsionar las funciones de riesgo estimadas, por lo que se podría hacer una representación menos restrictiva
  - El estimador de la función de supervivencia para la muestra cuando no hay censura es la siguiente:

$$\hat{S}(t) = \frac{n^o \text{ of spells such that } T \geq t}{n^o \text{ of total spells}}$$

- Esto se da porque, al no haber censura, entonces el mejor estimador para la probabilidad de que la duración sea mayor a  $t$  es la proporción de los spells (observaciones) que son mayores a  $t$  dentro de la muestra
- El estimador producto límite de Kaplan-Meier es un estimador estrictamente empírico y no paramétrico para la estimación de la función de riesgo y de supervivencia

$$\hat{\lambda}(t) = \frac{h_j}{n_j} \quad \hat{S}(t_j) = \prod_{i=1}^j \frac{n_i - h_i}{n_i} = \prod_{i=1}^j 1 - \hat{\lambda}(t_i)$$

- Se asume que hay  $K$  duraciones completas que corresponden a *spells* concretos están ordenadas en orden ascendente, y que  $K < N$  porque hay *spells* censurados y vínculos

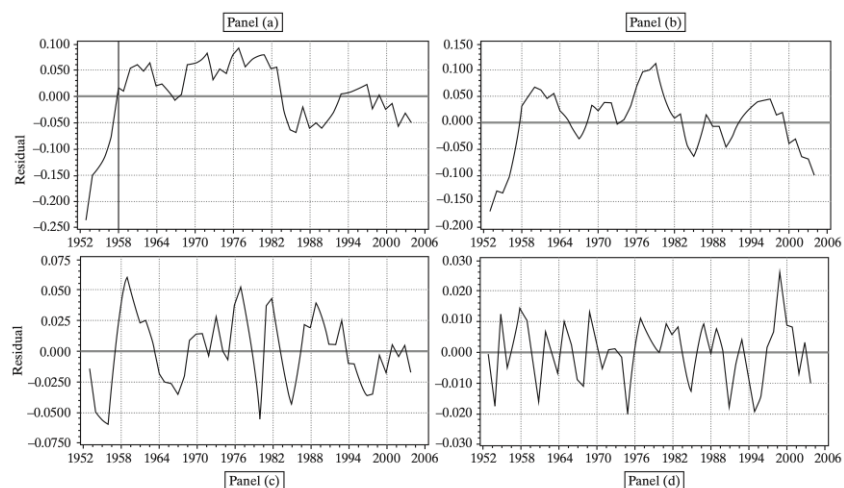
$$t_1 < t_2 < t_3 < \dots < t_K$$

- En este caso,  $h_j$  representa el número de observaciones o duraciones completas con una duración  $t_j$ ,  $m_j$  representa el número de observaciones censuradas entre  $t_j$  y  $t_{j+1}$ ,  $m_K$  representa el número de observaciones con una duración mayor a  $t_K$  y  $n_j$  representa el número de *spells* que no se han censurado ni completado antes de  $t_j$

$$n_j = \sum_{i \geq j}^K (m_i + h_i)$$

## La correlación serial

- Las series temporales normalmente presentan autocorrelación o correlación serial en las perturbaciones cuando se miran diferentes periodos
  - Una explicación para la autocorrelación es que los factores relevantes omitidos de la regresión de series temporales, como aquellos incluidos, están correlacionados en el tiempo
  - Este hecho se puede deber a que la correlación serial de los factores tendría que haberse incluido en la regresión. Por lo tanto, puede haber autocorrelación inducida por una mala especificación del modelo



- Los problemas para la estimación y la inferencia causados para la autocorrelación son parecidos (aunque más complejos) a aquellos causados por la heteroscedasticidad

- Los estimadores de mínimos cuadrados ordinarios son ineficientes y la inferencia basada en estos estimadores se ve afectada adversamente
  - Dependiendo del proceso subyacente, los estimadores de mínimos cuadrados generalizados y los de mínimos cuadrados generalizados factibles se pueden usar para solucionar este problema
- Un modelo de series temporales normalmente describe el camino de una variable  $y$  en términos de factores  $x_t$  contemporáneos, las perturbaciones  $\varepsilon_t$  (llamadas innovaciones) y sus valores pasados  $y_{t-1}, y_{t-2}, \dots$ 
  - Una serie temporal es una sola ocurrencia de un evento aleatorio (de la variable  $y$  en el momento  $t$ ), de modo que el historial entero sobre el periodo examinado en las series temporales constituye una sola realización del proceso
    - En relación a economía, el proceso no se puede repetir, y no hay contraparte al muestreo con repetición con datos transversales o una replicación del experimento que involucre una serie temporal en física o ingeniería
    - Sin embargo, podría haberse dado una realización completamente diferente de toda la serie temporal
  - La secuencia de observaciones  $\{y_t\}_{t=-\infty}^{t=\infty}$  es un proceso de serie temporal, el cual se caracteriza por su orden temporal y su correlación sistemática entre las observaciones de esta secuencia
    - La característica más importante de un proceso de series temporales es que, empíricamente, el mecanismo de generación de datos produce exactamente una realización de esta secuencia
    - Los resultados estadísticos basados en características muestrales no tienen relación con el muestreo aleatorio de una población, sino del muestreo aleatorio de las distribuciones de estadísticos contruidos a partir de conjuntos de observaciones tomadas de esta realización en una ventana temporal  $t = 1, 2, \dots, T$
    - En consecuencia, la teoría asintótica en este contexto está relacionada con el comportamiento de los estadísticos contruidos de una muestra con un tamaño incremental en esta secuencia
  - Las propiedades de  $y_t$  como variable aleatoria en datos transversales son directas y resumidas convenientemente en una proposición sobre

su media y su varianza o de la distribución generadora de  $y_t$ . En este caso, eso es menos obvio

- Es común asumir que las innovaciones se generan independientemente de un periodo al siguiente y que la media es nula y la varianza es homoscedástica. Bajo estas suposiciones, la distribución de  $\varepsilon_t$  sería débilmente estacionaria

$$E(\varepsilon_t|X) = E(\varepsilon_t) = 0 \quad \text{Var}(\varepsilon_t|X) = \text{Var}(\varepsilon_t) = \sigma_\varepsilon^2$$

$$\text{Cov}(\varepsilon_t, \varepsilon_s) = 0 \quad \text{for } t \neq s$$

- Bajo las suposiciones, aunque la noción de muestreo aleatorio se tiene que extender a las series temporales  $\varepsilon_t$ , los resultados matemáticos basados en la noción aplican aquí
- No obstante, hay una diferencia obvia entre las series temporales de  $y_t$  y de  $\varepsilon_t$ : las observaciones de  $y_t$  en diferentes puntos en el tiempo están necesariamente correlacionadas

- Suponiendo que la serie  $y_t$  es débilmente estacionaria y un modelo  $y_t = \beta_0 + \beta_1 y_{t-1} + \varepsilon_t$ , entonces se puede ver como la varianza de  $y_t$  es homoscedástica siempre que se cumpla que  $|\beta_1| < 1$  (para que la varianza sea positiva y finita)

$$E(y_t) = \beta_0 + \beta_1 E(y_{t-1}) + E(\varepsilon_t) \Rightarrow E(y_t) = \frac{\beta_0}{1 - \beta_1}$$

$$\text{Var}(y_t) = \beta_1^2 \text{Var}(y_{t-1}) + \text{Var}(\varepsilon_t) \Rightarrow \text{Var}(y_t) = \frac{\sigma_\varepsilon^2}{1 - \beta_1^2}$$

- Alternativamente, se puede visualizar sabiendo que la realización en un momento pasado depende de un momento anterior a ese, por lo que se puede expresar  $y_t$  como una suma infinita que depende de los coeficientes (como un proceso que genera los datos en un pasado infinito). Por lo tanto, estos tienen que ser  $|\beta_1| < 1$  para que la suma no diverja

$$y_t = \beta_0 + \beta_1(\beta_0 + \beta_1(\beta_0 + \beta_1 y_{t-3} + \varepsilon_t) + \varepsilon_t) + \varepsilon_t$$

$$\Rightarrow y_t = \sum_{i=0}^{\infty} \beta_1^i (\beta_0 + \varepsilon_{t-i})$$

- También se podría llegar a la misma conclusiones que con los últimos dos ejemplos si se asumen condiciones iniciales para  $y_0$  y  $\varepsilon_0$ , de modo que se asume que la observación de esta serie temporal se inicia en un momento  $t = 0$  en el que el proceso

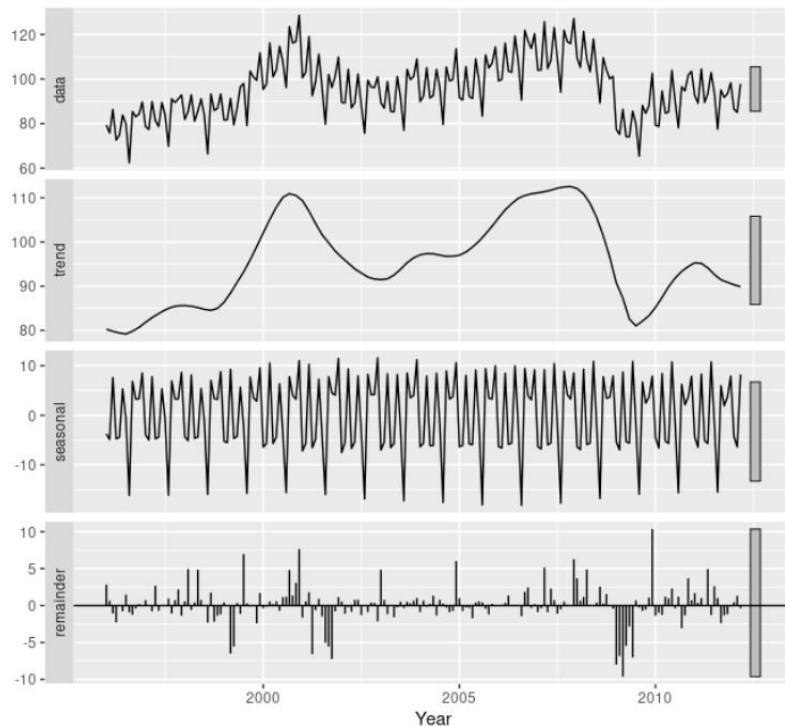
subyacente ha alcanzado un estado en donde la distribución de la variable ya no cambia

- Excepto en casos muy especiales, se espera que todos los elementos en el vector aleatorio  $y_t = (y_1, y_2, \dots, y_T)$  estén correlacionados (llamada autocorrelación), por lo que no se pueden utilizar los métodos de estimación para observaciones no correlacionadas vistos anteriormente no funcionan
  - Hay una contraparte a la estimación paramétrica para datos transversales, pero bajo suposiciones que hagan que los parámetros (en el sentido familiar) existan
  - Aún con estacionariedad, los resultados de muestra finita vistos anteriormente no son útiles
- Las series temporales pueden exhibir una variedad de patrones, por lo que ayuda mucho descomponer estas series en varios componentes
  - Al describir patrones de series temporales, se suelen utilizar tres conceptos básicos: la tendencia, la estacionalidad y la ciclicidad
    - La tendencia se puede entender como un crecimiento o decrecimiento prolongado en los datos
    - La estacionalidad se puede entender como un patrón que ocurre debido a factores estacionales como el momento del año o el día de la semana. Esta siempre tiene una frecuencia fija y conocida
    - La ciclicidad se puede entender como las caídas y subidas que no son de una frecuencia fija. Estas se suelen dar por factores económicos y por el ciclo comercial
  - A partir de estos tres tipos de patrones, se proponen dos tipos de descomposiciones para los datos: la descomposición aditiva y la descomposición multiplicativa

$$y_t = T_t + S_t + C_t + R_t \quad (\text{additive dec.})$$

$$y_t = T_t S_t C_t R_t \quad (\text{multiplicative dec.})$$

- En este caso,  $T_t$  representa el componente de tendencia,  $S_t$  representa el componente de estacionalidad, y  $R_t$  representa el componente residual, que son las fluctuaciones aleatorias o sistémicas. Debido a que la ciclicidad es muy difícil de identificar, se suele omitir en los paquetes más básicos de *software*



- La descomposición aditiva es más útil si la magnitud de las fluctuaciones del componente estacional no varía con el nivel de la serie temporal (la tendencia no es proporcional). Si, en cambio, la variación en el componente estacional es proporcional al componente tendencial, entonces la descomposición multiplicativa es más apropiada
- Una manera de utilizar la descomposición multiplicativa es hacer una transformación logarítmica de los datos hasta que la variación parezca ser estable en el tiempo, y así aplicar la descomposición aditiva (ya que serán descomposiciones equivalentes)

$$y_t = T_t S_t C_t R_t \Rightarrow \log y_t = \log T_t + \log S_t + \log C_t + \log R_t$$

- Ignorando el componente estacional, es posible hacer una estimación de la tendencia de una serie temporal a través de una regresión de mínimos cuadrados para el tiempo (el número de observaciones)

$$y_t = a + bt + \varepsilon_t$$

- Obviamente, la tendencia no tiene por qué ser perfectamente lineal, por lo que se puede generalizar a un modelo polinómico que se puede estimar por mínimos cuadrados ordinarios

$$y_t = a + b_1 t + b_2 t^2 + \dots + b_p t^p + \varepsilon_t$$



- Aunque la tendencia sea un predictor natural de los valores de  $y_t$ , debido a que las innovaciones  $\varepsilon_t$  están autocorrelacionadas, es posible obtener mejores estimadores para  $y_t$  teniendo en cuenta este hecho
- Una vez vistas las nociones básicas sobre series temporales, se pueden obtener resultados teóricos utilizando conclusiones sobre la autocorrelación y proponer modelos para el análisis de series temporales
  - En las series temporales comunes, las perturbaciones son homoscedásticas, pero correlacionadas entre observaciones. Por lo tanto, siendo  $\sigma^2 \mathbf{\Omega}$  una matriz completa definida positiva con constante  $Var(\varepsilon_t | \mathbf{X}) = \sigma^2$  en la diagonal, se obtiene el siguiente resultado:

$$E(\varepsilon \varepsilon' | \mathbf{X}) = \sigma^2 \mathbf{\Omega}$$

- Se debe asumir que  $\mathbf{\Omega}_{t,s}$  es una función de  $|t - s|$  (lo cual es una suposición de estacionariedad), pero no de  $t$  o de  $s$  únicamente. Esto implica que la covarianza entre observaciones  $t$  y  $s$  son una función solo de  $|t - s|$ , la distancia temporal entre las observaciones
- Debido a que  $\sigma^2$  no está restringida, se normaliza  $\mathbf{\Omega}_{t,t} = 1$
- A partir de esto, se puede definir la noción de autocovarianza y la de autocorrelación
  - La autocovarianza se define como la covarianza condicional de los errores a  $\mathbf{X}$ . La autocovarianza de dos errores en el mismo momento será la varianza de  $\varepsilon_t$

$$\gamma_s \equiv Cov(\varepsilon_t, \varepsilon_{t-s} | \mathbf{X}) = Cov(\varepsilon_{t+s}, \varepsilon_t | \mathbf{X}) = \sigma^2 \mathbf{\Omega}_{t,t-s}$$

$$\gamma_0 \equiv Cov(\varepsilon_t, \varepsilon_t | \mathbf{X}) = \sigma^2 \mathbf{\Omega}_{t,t} = Var(\varepsilon_t | \mathbf{X}) = \sigma^2$$

- La autocorrelación se define como la correlación condicional de los errores a  $\mathbf{X}$

$$\rho_s \equiv Corr(\varepsilon_t, \varepsilon_{t-s} | \mathbf{X}) = \frac{Cov(\varepsilon_t, \varepsilon_{t-s} | \mathbf{X})}{\sqrt{Var(\varepsilon_t | \mathbf{X})} \sqrt{Var(\varepsilon_{t-s} | \mathbf{X})}} = \frac{\gamma_s}{\gamma_0}$$

- Por lo tanto, se puede escribir la covarianza en términos de la esperanza, resultando en una matriz  $\mathbf{\Gamma}$  de autocovarianzas, que es igual a la matriz de autocorrelaciones  $\mathbf{R}$  multiplicada por  $\gamma_0$

$$E(\varepsilon' \varepsilon | \mathbf{X}) = \sigma^2 \mathbf{\Omega}_{t,t-s} = \mathbf{\Gamma} = \gamma_0 \mathbf{R}$$

- Diferentes procesos implican diferentes patrones para la matriz **R**, pero principalmente se pueden identificar tres procesos: el proceso autorregresivo, el proceso de media móvil y el proceso autorregresivo con media móvil

- El caso más estudiado para los errores es el modelo autorregresivo de primer orden  $AR(1)$ , en donde  $u_t$  es estacionario no autocorrelacionado (es un ruido blanco) y  $\rho$  es un parámetro

$$\varepsilon_t = \rho \varepsilon_{t-1} + u_t$$

- Los procesos autorregresivos  $AR(p)$  de mayor orden implican unos patrones más complejos, incluyendo, para algunos valores de los parámetros, un comportamiento cíclico en el comportamiento de las autocorrelaciones

$$\varepsilon_t = \gamma_1 \varepsilon_{t-1} + \gamma_2 \varepsilon_{t-2} + \dots + \gamma_p \varepsilon_{t-p} + u_t$$

$$C(L)\varepsilon_t = u_t \quad \text{for} \quad C(L)x_t \equiv (1 - \gamma_1 L - \gamma_2 L^2 - \dots - \gamma_p L^p)x_t$$

- Las autorregresiones estacionarias están estructuradas de tal manera que la influencia de una perturbación dada se desvanece a medida que retrocede al pasado más distante, pero solo desaparece asintóticamente
- En contraste, los procesos de media móvil tienen una memoria corta, y el proceso de media móvil más utilizado es el proceso  $MA(1)$ . Para este modelo, la memoria del proceso es de solo un periodo, lo cual se puede ver calculando las autocovarianzas

$$\varepsilon_t = u_t - \theta_1 u_{t-1}$$

$$\gamma_0 = \text{Var}(\varepsilon_t, \varepsilon_t) = \sigma_u^2(1 + \lambda_1^2)$$

$$\gamma_1 = \text{Cov}(\varepsilon_{t+1}, \varepsilon_t) = \text{Cov}(-\lambda_1 u_t, u_t) = -\lambda_1 \sigma_u^2$$

$$\gamma_{s>1} = \text{Cov}(\varepsilon_{t+s}, \varepsilon_t) = 0$$

- Los procesos de media móvil  $MA(p)$  de mayor orden implican unos patrones más complejos, y estos se pueden expresar de la siguiente manera:

$$\varepsilon_t = u_t - \theta_1 u_{t-1} - \theta_2 u_{t-2} - \dots - \theta_q u_{t-q}$$

$$\varepsilon_t = D(L)u_t \text{ for } D(L)x_t \equiv (1 - \theta_1 L - \theta_2 L^2 - \dots - \theta_q L^q)x_t$$

- Este modelo también se puede expresar con coeficientes con signo positivo, aunque eso no cambia las propiedades teóricas del modelo (solo el signo)

$$\varepsilon_t = u_t + \theta_1 u_{t-1} + \theta_2 u_{t-2} - \dots + \theta_q u_{t-q}$$

$$\varepsilon_t = -D(L)u_t \text{ for } D(L)x_t \equiv (1 + \theta_1 L + \dots - \theta_q L^q)x_t$$

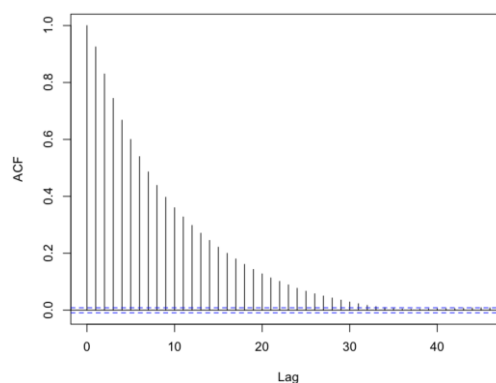
- Finalmente, también es posible combinar los dos modelos mostrados anteriormente en un modelo autorregresivo con media móvil, denotado por  $ARMA(p, q)$

$$\varepsilon_t = \gamma_1 \varepsilon_{t-1} + \dots + \gamma_p \varepsilon_{t-p} + u_t - \theta_1 u_{t-1} - \dots - \theta_q u_{t-q}$$

$$C(L)\varepsilon_t = D(L)u_t$$

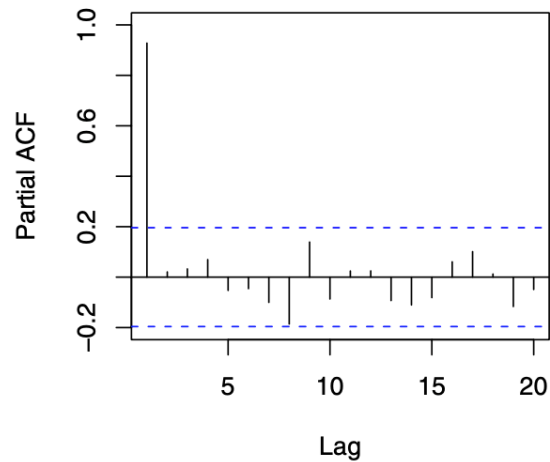
- Para poder seleccionar el orden de estos modelos, se necesita tener en cuenta la función de autocorrelación y la función de autocorrelación parcial a la vez

- La función de autocorrelación parcial es similar a la de autocorrelación, pero esta no tiene en cuenta la correlación existente en los valores entre  $t$  y  $t - s$ , sino solo la del valor en  $t$  y en  $t - s$
- Para un modelo autorregresivo, la función de autocorrelación decrece de manera exponencial a cero, dado que las observaciones pasadas influyen en las observaciones presentes en el modelo, pero esta influencia disminuye cada vez más

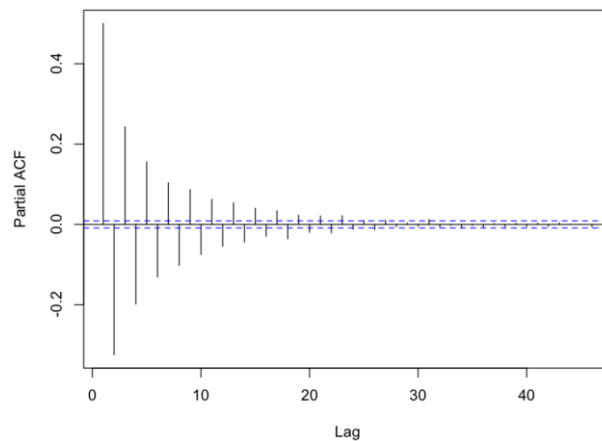


- El número de retrasos significativos en la función de autocorrelación parcial (independientemente de la frecuencia de los datos) determinará el orden del proceso  $AR$ , ya que esos son

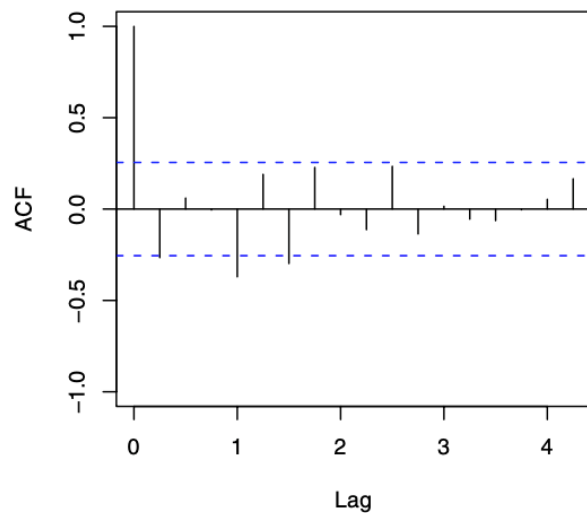
los retrasos que determinan significativamente más la autocorrelación



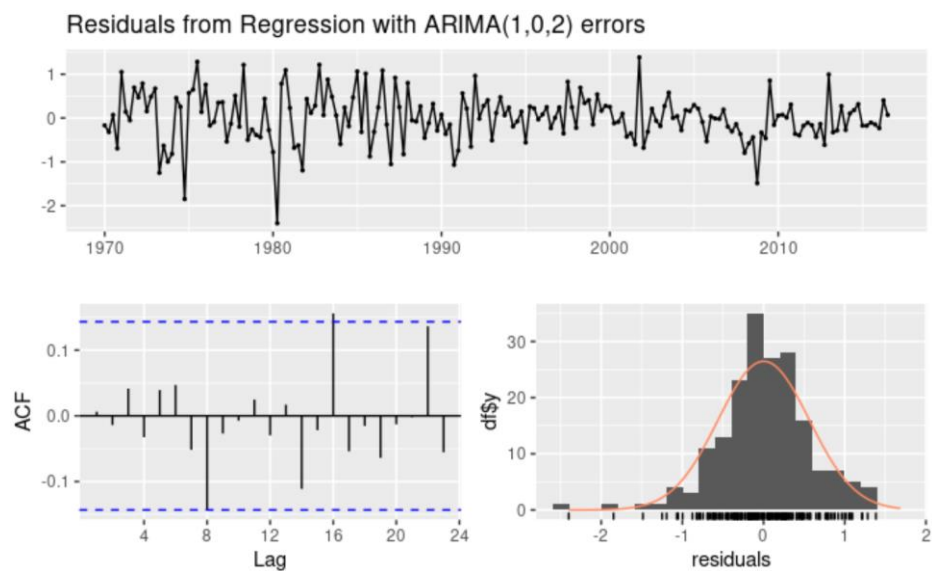
- Para un modelo de media móvil, la función de autocorrelación parcial decrece de manera exponencial a cero, dado que las observaciones pasadas no influyen en las observaciones presentes en el modelo solo se determina por los *shocks*, aunque estos influyen cada vez menos



- El número de retrasos significativos en la función de autocorrelación ((independientemente de la frecuencia de los datos) determinará el orden del proceso *MA*, ya que esos son los retrasos que determinan significativamente más la autocorrelación



- Con un proceso mixto como el *ARMA* o el *ARIMA*, ambas funciones exhiben un descenso exponencial, por lo que no se pueden identificar un *AR* o un *MA* por separado. El orden, en este caso, vendrá determinado de otros modos (a través del uso de criterios de información u otros)
- Si a través del modelaje del término de error se consigue crear una serie temporal estacionaria, entonces los residuos del modelo estimado deberían comportarse como un proceso de ruido blanco



- En este caso, no se puede apreciar ningún patrón de tendencia o de estacionalidad en los residuos. Además, la mayoría de retrasos no son significativos cuando se mira la función de autocorrelación (solo algunos muy grandes, lo cual es común) y se distribuye como una normal (aunque haya algunas observaciones extremas debido a la muestra)

- Los modelos anteriores pueden ser caracterizados por su orden, el valor de sus parámetros y el comportamiento de sus autocorrelaciones, por lo que se pueden considerar varias formas en diferentes puntos. No obstante, el modelo  $AR(1)$  es el modelo más utilizado con diferencia

- Los procesos más complejos suelen ser muy difíciles de analizar, pero una razón más práctica para el uso de este modelo es que es muy optimista pensar que se puede saber precisamente la forma correcta para el modelo de las perturbaciones para cualquier situación
  - La autorregresión de primer orden ha demostrado a lo largo del tiempo que es un modelo razonable para el proceso subyacente de las perturbaciones, el cual en realidad seguramente es impenetrable
- El modelo autorregresivo de primer orden es representado en su forma autorregresiva de la siguiente manera:

$$\varepsilon_t = \rho \varepsilon_{t-1} + u_t$$

- En este caso, se asume que el término de error en este proceso cumple las siguientes condiciones:

$$E(u_t | \mathbf{X}) = 0 \quad \text{Var}(u_t | \mathbf{X}) = \sigma_u^2$$

$$\text{Cov}(u_t, u_s | \mathbf{X}) = 0 \quad \text{for } t \neq s$$

- Debido a que  $u_t$  es un ruido blanco, los momentos condicionales son equivalentes a los momentos incondicionales, por lo que  $E(u_t | \mathbf{X}) = E(u_t)$  y así
- Por sustitución repetitiva, se puede ver que el proceso autorregresivo se puede expresar de la siguiente forma:

$$\varepsilon_t = u_t + \rho u_{t-1} + \rho^2 u_{t-2} + \dots$$

- A partir del proceso de media móvil, es evidente que cada perturbación  $\varepsilon_t$  carga con el historial entero de  $u$ , en donde las observaciones más recientes son las que tienen más peso y las más alejadas tienen menos. Los signos en este caso son positivos, aunque también se puede usar el signo opuesto sin cambiar las propiedades teóricas del modelo
- Dado que estos valores de  $u_t$  no están correlacionados, la varianza de  $\varepsilon_t$  se puede expresar de la siguiente forma:

$$\text{Var}(\varepsilon_t) = \sigma_u^2 + \rho^2 \sigma_u^2 + \rho^4 \sigma_u^2 + \dots$$

- Se impone la restricción de que  $|\rho| < 1$  debido a que, de otra manera, la suma de la parte derecha de la ecuación anterior tiende a infinito (suposición de estacionariedad discutida anteriormente). Esto, a su vez, implica las siguientes propiedades para el proceso autorregresivo de primer orden:

$$\lim_{s \rightarrow \infty} \rho^s = 0 \quad E(\varepsilon_t) = 0 \quad Var(\varepsilon_t) = \frac{\sigma_u^2}{1 - \rho^2}$$

- Con la suposición de estacionariedad, hay una forma fácil de obtener la varianza y la covarianza para  $t$  y  $t - 1$ , debido a que  $Cov(u_t, \varepsilon_s) = 0$  para  $t > s$  y  $Var(\varepsilon_t) = Var(\varepsilon_{t-1})$

$$\begin{aligned} Cov(u_t, \varepsilon_s) &= Cov(u_t, u_t + \rho u_{t-1} + \rho^2 u_{t-2} + \dots) = \\ &= Cov(u_t, u_s) + \rho Cov(u_t, u_{s-1}) + \rho^2 Cov(u_t, u_{s-2}) + \dots = 0 \end{aligned}$$

$$\begin{aligned} \Rightarrow Var(\varepsilon_t) &= \rho^2 Var(\varepsilon_{t-1}) + Var(u_t) + 2Cov(u_t, \varepsilon_s) = \\ &= \rho^2 Var(\varepsilon_{t-1}) + \sigma_u^2 = \frac{\sigma_u^2}{1 - \rho^2} \end{aligned}$$

$$\begin{aligned} \Rightarrow Cov(\varepsilon_t, \varepsilon_{t-1}) &= E(\varepsilon_t \varepsilon_{t-1}) = E[\varepsilon_{t-1}(\rho \varepsilon_{t-1} + u_t)] \\ &= \rho Var(\varepsilon_{t-1}) = \frac{\rho \sigma_u^2}{1 - \rho^2} \end{aligned}$$

- Sustituyendo repetidamente, se puede obtener la siguiente expresión de  $\varepsilon_t$  en términos de las perturbaciones  $\varepsilon$  y de  $u$  para cualquier  $s$ :

$$\varepsilon_t = \rho^s \varepsilon_{t-s} + \sum_{i=0}^{s-1} \rho^i u_{t-i}$$

- A partir de esta, se puede obtener una expresión para la covarianza y para la correlación para cualquier  $t$  y  $t - s$ . Con la suposición de estacionariedad, se puede ver como la covarianza y la correlación se desvanece con el tiempo

$$\begin{aligned} Cov(\varepsilon_t, \varepsilon_{t-s}) &= E[\varepsilon_t \varepsilon_{t-s}] = \\ &= \rho^s Cov(\varepsilon_{t-s}, \varepsilon_{t-s}) + \sum_{i=0}^{s-1} \rho^i Cov(u_{t-i}, \varepsilon_{t-s}) = \rho^s Var(\varepsilon_{t-s}) = \frac{\rho^s \sigma_u^2}{1 - \rho^2} \end{aligned}$$

$$\text{Corr}(\varepsilon_t, \varepsilon_{t-s}) = \frac{\text{Cov}(\varepsilon_t, \varepsilon_{t-s})}{\text{Var}(\varepsilon_{t-s})} = \rho_s = \rho^s$$

- Dependiendo del signo de  $\rho$ , las autocorrelaciones pueden decrecer siguiendo una progresión geométrica o alternarse en signo si  $\rho$  es negativa

$$\sigma^2 \mathbf{\Omega} = \frac{\sigma_u^2}{1 - \rho^2} \begin{bmatrix} 1 & \rho & \rho^2 & \rho^3 & \cdots & \rho^{T-1} \\ \rho & 1 & \rho & \rho^2 & \cdots & \rho^{T-2} \\ \rho^2 & \rho & 1 & \rho & \cdots & \rho^{T-3} \\ \vdots & \vdots & \vdots & \vdots & \cdots & \rho \\ \rho^{T-1} & \rho^{T-2} & \rho^{T-3} & \cdots & \rho & 1 \end{bmatrix}$$

## Los datos no estacionarios

- La mayoría de variables económicas que exhiben tendencias son no estacionarias y, por tanto, no se puede utilizar el análisis anteriormente planteado porque presenta varias complicaciones
  - Los tres operadores más importantes a la hora de modelar series temporales es el operador de retraso  $L$ , el de primeras diferencias  $\Delta$ , y el de diferencias estacionales  $\Delta_s$ , los cuales simplifican las matemáticas del análisis. Estos se definen por las siguientes operaciones:

$$Ly_t \equiv y_{t-1}$$

$$\Delta y_t \equiv (1 - L)y_t = y_t - y_{t-1} \quad \Delta_s y_t \equiv (1 - L^s)y_t = y_t - y_{t-s}$$

- A partir de la definición del operador  $L$ , se puede ver como este cumple las siguientes propiedades:

$$L^p y_t = L \left( L \left( L \left( \dots (Ly_t) \right) \right) \right) = L \left( L \left( L \left( \dots (y_{t-1}) \right) \right) \right) = y_{t-p}$$

$$(L^p)^q y_t = L^{pq} y_t = y_{t-pq}$$

$$(L^p)(L^q)y_t = L^{p+q}y_t = L^p y_{t-q} = y_{t-q-p}$$

- A partir de la definición del operador  $\Delta$ , se puede ver como este cumple la siguiente propiedad:

$$\begin{aligned} \Delta^p &= \Delta \left( \Delta \left( \dots (\Delta(\Delta y_t)) \right) \right) = \Delta \left( \Delta \left( \dots (\Delta(y_t - y_{t-1})) \right) \right) = \\ &= (y_t - y_{t-1}) - (y_{t-1} - y_{t-2}) - \cdots - (y_{t-p+1} - y_{t-p}) \end{aligned}$$



- A partir de la definición del operador  $\Delta_s$ , se puede ver como este cumple la siguiente propiedad:

$$\Delta_s Y_t = (1 + L + \dots + L^{s-1}) \Delta Y_t = y_t - y_{t-s}$$

- Finalmente, en una serie autorregresiva de primer orden con  $|\beta| < 1$ , se puede expresar  $y_t$  de la siguiente manera:

$$y_t = \beta y_{t-1} + \varepsilon_t \Rightarrow (1 - \beta L) y_t = \varepsilon_t$$

$$\Rightarrow y_t = \left( \frac{1}{1 - \beta L} \right) \varepsilon_t = (1 + \beta L + \beta^2 L^2 + \dots) \varepsilon_t = \sum_{s=0}^{\infty} \beta^s \varepsilon_{t-s}$$

- Un proceso que figura comúnmente en situaciones de series temporales no estacionarias es el camino aleatorio o *random walk* con una tasa de deriva o *drift*, el cual tiene la siguiente forma:

$$y_t = \mu + y_{t-1} + \varepsilon_t$$

- Por sustitución directa, se puede ver como  $y_t$  es la suma infinita de perturbaciones y de la tasa de deriva (con una posible media diferente de cero). Si las innovaciones  $\varepsilon_t$  se generan por el mismo proceso con media nula y varianza constante, entonces la varianza de  $y_t$  sería infinita

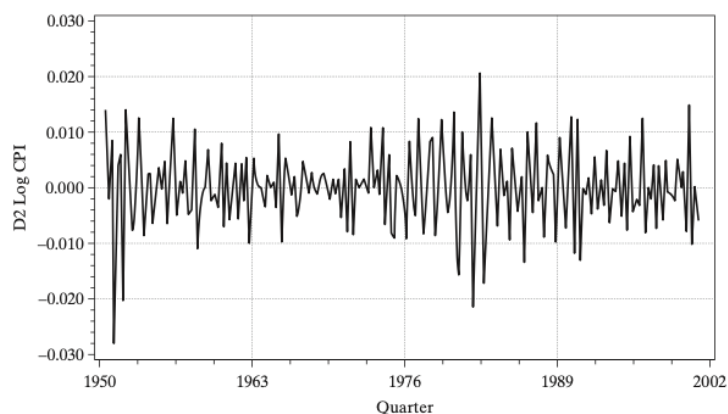
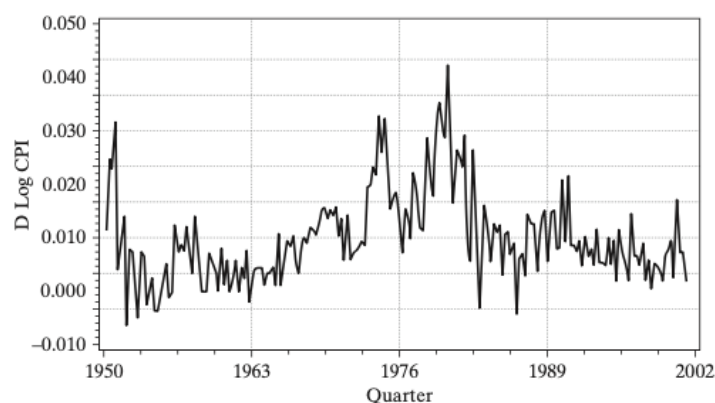
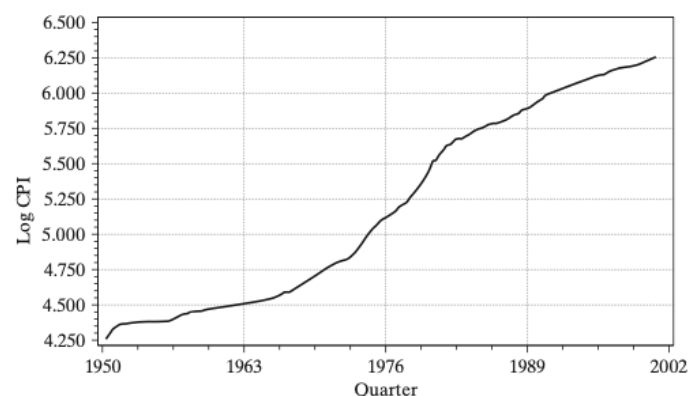
$$y_t = \frac{\mu + \varepsilon_t}{1 - L} = \sum_{s=0}^{\infty} (\mu + \varepsilon_{t-s})$$

- Debido a que la varianza sería infinita, se puede ver como el proceso no sería un proceso estacionario, aunque  $\mu = 0$ . No obstante, el proceso de la primera diferencia  $\Delta y_t$  es el término de innovaciones más la media de la diferencia, de modo que es estacionario (media y varianza constante)

$$\Delta y_t = y_t - y_{t-1} = \mu + \varepsilon_t$$

$$E(\Delta y_t) = \mu \quad Var(\Delta y_t) = \sigma_\varepsilon^2$$

- Por lo tanto, la serie temporal de  $y_t$  se denomina serie integrada de primer orden, denotado por  $I(1)$ , dado que tomar la primera diferencia produce un proceso estacionario. Una serie no estacionaria es una serie integrada de  $d$ -ésimo orden, denotado por  $I(d)$ , si tomar la primera diferencia  $d$  veces produce una serie estacionaria



- Una generalización del modelo *ARMA* sería el modelo *ARIMA*( $p, d, q$ ), el modelo autorregresivo integrado con media móvil, el cual se puede escribir de la siguiente forma:

$$\Delta^d y_t = (1 - L)^d y_t = \mu + \gamma_1 \Delta^d y_{t-1} + \dots + \gamma_p \Delta^d y_{t-p} + \varepsilon_t - \theta_1 \varepsilon_{t-1} - \dots - \theta_q \varepsilon_{t-q}$$

- Este resultado se puede escribir de manera compacta utilizando los polinomios  $C(L)$  y  $D(L)$  en el operador de retraso

$$C(L)[(1 - L)^d y_t] = \mu + D(L)\varepsilon_t$$

$$C(L)x_t \equiv (1 - \gamma_1 L - \gamma_2 L^2 - \dots - \gamma_p L^p)x_t$$

$$D(L)x_t \equiv (1 - \theta_1 L - \theta_2 L^2 - \dots - \theta_q L^q)x_t$$

- En su investigación de 1974, Granger y Newbold discutieron sobre como la alta autocorrelación entre los residuos en los modelos de regresión convencionales pueden hacer que los contrastes de hipótesis fueran engañosos. Esto es causado por la relación entre la existencia de tendencia, la raíz unitaria y el camino aleatorio
  - Con este tipo de datos, los contrastes convencionales tenderán resultar en relaciones significativas entre variables cuando, en verdad, no hay ninguna relación
    - El resultado general es que la regresión lineal convencional de un proceso de camino aleatorio sobre otro sugeriría que existe una relación significativa, aunque estas dos sean realmente independientes. A este tipo de regresiones se las llama regresiones espurias
  - Para la mayoría de series temporales, el proceso de camino aleatorio con tasa de deriva y el proceso de tendencia estacionaria, en donde  $\varepsilon_t$  es ruido blanco en cada caso, son buenas caracterizaciones

$$y_t = \mu + y_{t-1} + \varepsilon_t \quad (\text{random walk with drift})$$

$$y_t = \mu + \beta t + \varepsilon_t \quad (\text{trend stationary})$$

- Claramente, ambas regresiones producen series no estacionarias y con una tendencia fuerte, de modo que no es sorprendente que regresiones que involucren este tipo de variables casi siempre produzcan resultados significativos
- Esta fuerte correlación parece ser debida a la tendencia subyacente, pero la intuición de este razonamiento es menos clara en el caso de un camino aleatorio puro, donde no hay una tasa de deriva que marque una tendencia clara (ni analíticamente ni gráficamente) pero las relaciones parecen persistir aún con series no relacionadas

$$y_t = y_{t-1} + \varepsilon_t \quad (\text{pure random walk})$$

- El factor común en estos tres tipos de procesos no es la tendencia, sino que el proceso generador de datos de los tres procesos se puede escribir como una ecuación con raíz unitaria:

$$(1 - L)y_t = \alpha + v_t$$

- Esta ecuación característica tiene una sola raíz equivalente a uno (cuando se considera  $L$  como variable), por lo que se le llama raíz unitaria
- El uso de datos caracterizados por raíces unitarias tiene el potencial de llevar a errores de inferencia serios, y estos datos pertenecen a series temporales no estacionarias
- Para los tres procesos, lo más natural sería comenzar diferenciando con tal de hacer a la serie estacionaria. No obstante, no es inmediatamente obvio como proceder correctamente, dado que no necesariamente se transforma la serie en estacionaria si se aplica un enfoque incorrecto
  - Hacer la primera diferencia para un proceso de camino aleatorio con tasa de deriva y/o puro permite obtener series estacionarias, pero hacer la primera diferencia para un proceso de tendencia estacionaria hace que se intercambie la tendencia por autocorrelación en forma de un proceso de media móvil de primer orden

$$\Delta y_t = \mu + \varepsilon_t \quad (\text{random walk with drift})$$

$$\Delta y_t = \varepsilon_t \quad (\text{pure random walk})$$

$$\Delta y_t = \beta + \varepsilon_t - \varepsilon_{t-1} \quad (\text{trend stationary})$$

- Por otro lado, extraer la tendencia (calcular los residuos de una regresión sobre el tiempo  $t$ ) es contraproducente, mientras que para un proceso de tendencia estacionaria sería correcto

$$y_t - \hat{y}_t = \varepsilon_t - u_t \quad \text{where} \quad \hat{y}_t \equiv \alpha + \beta t + u_t$$

- Dado que ninguno de estos enfoques es preferible de manera obvia, se puede considerar la primera diferencia en un modelo mixto que considere tanto una tendencia temporal como el primer retraso, y se puede ver como esta ecuación proporciona una base para tests de raíces unitarias con datos económicos

$$y_t = \mu + \beta t + \gamma y_{t-1} + \varepsilon_t$$

$$\Rightarrow y_t - y_{t-1} = \mu + \beta t + (\gamma - 1)y_{t-1} + \varepsilon_t$$

- En principio, un test de contraste de hipótesis para  $\gamma - 1 = 0$  o, equivalentemente, para  $\alpha_2 = 0$ , da una confirmación de que el proceso es un camino aleatorio con tasa de deriva (si  $\alpha_1 = 0$  también) y se puede aplicar la primera diferencia. Si  $\gamma - 1 < 0$ , entonces la evidencia favorece a un proceso de tendencia

estacionaria (o a otro modelo), y la extracción de la tendencia (o alguna alternativa) es preferible

- Las implicaciones de las raíces unitarias son profundas, dado que dependiendo de como es la serie temporal no estacionaria, el modelo permite inferir una u otra cosa
  - Si el modelo real de una variable es el de una serie integrada de primer orden (un camino aleatorio), entonces los *shocks* de las innovaciones son permanentes
  - Si el modelo real de una variable es el de una serie que sigue un proceso de tendencia estacionaria, los *shocks* de las innovaciones son temporales
- El único problema práctico con el tipo de test propuesto es que los procedimientos de inferencia estándar no son válidos cuando las series son estacionarias (las convergencias no se dan), de modo que se necesitan métodos alternativos
  - Considerando un modelo  $AR(1)$  con media nula y con innovaciones que siguen un proceso de ruido blanco, se puede obtener que el sesgo del estimador de mínimos cuadrados cuando  $\gamma \rightarrow 1$  es negativo

$$y_t = \gamma y_{t-1} + \varepsilon_t$$

$$\hat{\gamma}_{OLS} = \frac{\sum_{t=2}^T y_t y_{t-1}}{\sum_{t=2}^T y_{t-1}^2} \Rightarrow E(\hat{\gamma}_{OLS}) - \gamma < 0$$

- Para  $|\gamma| < 1$ , en cambio, el estimador de mínimos cuadrados converge en probabilidad al verdadero parámetro  $\gamma$  y, por el teorema central del límite, se puede ver que la distribución del estadístico para  $\gamma$  sigue una distribución normal de media nula y varianza  $1 - \gamma^2$

$$\hat{\gamma}_{OLS} = \frac{\sum_{t=2}^T y_t y_{t-1}}{\sum_{t=2}^T y_{t-1}^2} \Rightarrow \lim_{T \rightarrow \infty} P(|\hat{\gamma}_{OLS} - \gamma| > c) = 0$$

$$\Rightarrow \sqrt{T}(\hat{\gamma}_{OLS} - \gamma) \xrightarrow{D} N(0, 1 - \gamma^2)$$

- Este resultado, es dudoso para el caso en el que  $\gamma = 1$ , ya que implicaría una distribución con varianza nula. No obstante, Dickey y Fuller, en su investigaciones de 1979 y 1981, mostraron la siguiente convergencia en distribución cuando  $\gamma = 1$ :

$$T(\hat{\gamma}_{OLS} - \gamma) \xrightarrow{D} v$$

- Esta convergencia expresa que el estadístico tiende a una variable  $v$  con varianza finita y positiva y, para muestras finitas, con un sesgo negativo, dado que  $E(\hat{\gamma}_{OLS}) < 1$
  - De este modo, los resultados obtenidos por Dickey y Fuller muestran como el estimador tiene un sesgo a la baja y como la convergencia de este estimador es más rápido a su valor poblacional, dado que la varianza bajo la hipótesis nula es  $O(1/T^2)$  y no  $O(1/T)$
  - Las medidas convencionales tenderán a esconder el valor real de  $\gamma$ , dado que  $\hat{\gamma}_{OLS}$  estará sesgado y, debido a la varianza tan pequeña que tendrá el estimador, los tests de hipótesis con el estadístico  $t$  tenderán a rechazar la hipótesis nula de que  $\gamma = 1$
  - Por lo tanto, ambos propusieron seguir usando este estadístico, pero revisar los valores para los cuales se rechaza la hipótesis de que  $\gamma = 1$ , los cuales ambos propusieron al obtenerlos a través de simulación de Monte Carlo
- Para poder contrastar la hipótesis de raíces unitarias, se pueden utilizar diferentes versiones del test de Dickey-Fuller dependiendo del tipo de serie temporal

- En todas las versiones siguientes se asume que el término de error se distribuye normalmente y que no hay covarianza intertemporal entre dos realizaciones del término

$$\varepsilon \sim N(0, \sigma^2) \quad \text{Cov}(\varepsilon_t, \varepsilon_{ss}) = 0 \text{ for } \forall t \neq s$$

- La versión más simple del modelo a analizar es la de un camino aleatorio

$$y_t = \gamma y_{t-1} + \varepsilon_t$$

- No obstante, un camino aleatorio simple suele ser inadecuado para muchas series, de modo que se planea utilizar un camino aleatorio con tasa de deriva

$$y_t = \mu + z_t \text{ where } z_t = \gamma z_{t-1} + \varepsilon_t$$

$$\Rightarrow y_t = \mu(1 - \gamma) + \gamma z_{t-1} + \varepsilon_t$$

- No obstante, un camino aleatorio simple suele ser inadecuado para muchas series, de modo que se planea utilizar un camino aleatorio con tasa de deriva

$$y_t = \mu + \beta t + z_t \text{ where } z_t = \gamma z_{t-1} + \varepsilon_t$$

$$\Rightarrow y_t = [\mu(1 - \gamma) + \gamma\beta] + \beta(1 - \gamma)t + \gamma y_{t-1} + \varepsilon_t$$

- Bajo la hipótesis nula de que  $\gamma = 1$ , se puede utilizar el estadístico  $t$  para  $\gamma$  para cualquiera de las versiones vistas, solo que se tiene que tener en cuenta los valores críticos correspondientes para cada caso

$$DF = \frac{\hat{\gamma} - 1}{Est.Std.Error(\hat{\gamma})}$$

	Sample Size			
	25	50	100	$\infty$
F ratio (D-F) <sup>a</sup>	7.24	6.73	6.49	6.25
F ratio (standard)	3.42	3.20	3.10	3.00
AR model <sup>b</sup> (random walk)				
0.01	-2.66	-2.62	-2.60	-2.58
0.025	-2.26	-2.25	-2.24	-2.23
0.05	-1.95	-1.95	-1.95	-1.95
0.10	-1.60	-1.61	-1.61	-1.62
0.975	1.70	1.66	1.64	1.62
AR model with constant (random walk with drift)				
0.01	-3.75	-3.59	-3.50	-3.42
0.025	-3.33	-3.23	-3.17	-3.12
0.05	-2.99	-2.93	-2.90	-2.86
0.10	-2.64	-2.60	-2.58	-2.57
0.975	0.34	0.29	0.26	0.23
AR model with constant and time trend (trend stationary)				
0.01	-4.38	-4.15	-4.04	-3.96
0.025	-3.95	-3.80	-3.69	-3.66
0.05	-3.60	-3.50	-3.45	-3.41
0.10	-3.24	-3.18	-3.15	-3.13
0.975	-0.50	-0.58	-0.62	-0.66

- Los contrastes de Dickey-Fuller asumen que las perturbaciones del modelo son ruido blanco, por lo que, si no se cumple la suposición, este contraste no es adecuado. No obstante, se puede extender al contraste aumentado de Dickey-Fuller, el cual se basa en el siguiente modelo:

$$y_t = \mu + \beta t + \gamma y_{t-1} + \gamma_1 \Delta y_{t-1} + \dots + \gamma_p \Delta y_{t-p} + \varepsilon_t$$

- El estadístico de contraste para este modelo es el mismo que antes, de modo que solo se tienen que considerar los valores críticos correctos

$$DF = \frac{\hat{\gamma} - 1}{Est.Std.Error(\hat{\gamma})}$$

- El camino aleatorio simple se obtiene imponiendo  $\mu = 0$  y  $\beta = 0$ , mientras que si  $\beta = 0$ , el camino aleatorio tiene una tasa de deriva

- Una alternativa sugerida para mejorar las propiedades del contraste para muestras finitas o de generalización es el contraste de Phillips-Perron (PP). Este se basa en el siguiente modelo, en donde  $\delta_t$  puede ser 0,  $\mu$  o  $\mu + \beta t$ :

$$y_t = \delta_t + \gamma y_{t-1} + \gamma_1 \Delta y_{t-1} + \dots + \gamma_p \Delta y_{t-p} + \varepsilon_t$$

- El procedimiento modifica los estadísticos de Dickey-Fuller para obtener el siguiente estadístico, llamado estadístico PP:

$$Z = \sqrt{\frac{c_0}{a}} \left( \frac{\hat{\gamma} - 1}{v} \right) - \frac{1}{2} (a - c_0) \frac{Tv}{\sqrt{as^2}}$$

where

$$s^2 = \frac{\sum_{t=1}^T e_t^2}{T - K} \quad \& \quad v^2 = \text{estim. asympt. variance of } \hat{\gamma}$$

$$c_j = \frac{1}{T} \sum_{s=j+1}^T e_t e_{t-s} \quad \text{for } j = 0, \dots, L \quad \& \quad c_0 = \left( \frac{T-K}{T} \right) s^2$$

$$a = c_0 + 2 \sum_{j=1}^L \left( 1 - \frac{j}{L+1} \right) c_j$$

- Como se puede ver, este tipo de contraste utiliza las mismas ecuaciones que el ADF pero propone una corrección no paramétrica del estadístico para permitir la autocorrelación de los residuos
- Igual que el ADF, tampoco sigue una distribución t-Student, por lo que los valores críticos también se tienen que derivar a partir de simulaciones de Monte Carlo. No obstante, como las ecuaciones y las hipótesis nulas son las mismas, la distribución es la misma y se utilizan los mismos valores críticos que con el ADF
- Otra alternativa también utilizada es el contraste KPSS, el cual permite contrastar la hipótesis nula de estacionariedad en un modelo con  $\varepsilon_t$  estacionaria y  $Z_t$  es una serie estacionaria con media nula y varianza uno en el siguiente modelo:

$$y_t = \alpha + \beta t + \gamma Z_t + \varepsilon_t \quad \text{for } t = 1, 2, \dots, T$$

- Si  $\gamma = 0$ , entonces el proceso es estacionario si  $\beta = 0$  y es estacionario con tendencia si  $\beta \neq 0$ . Como  $Z_t$  es  $I(1)$ , entonces  $y_t$  es no estacionaria si  $\gamma$  no es cero



- El test KPSS para la hipótesis nula de que  $\gamma = 0$  contra la alternativa de que no lo es invierte la estrategia del estadístico de Dickey-Fuller, lo cuál permite que se tenga una mayor potencia en el contraste de hipótesis cuando  $\gamma$  está muy cerca de 0 dado que se usa un estadístico estacionario y no uno no estacionario como el de ADF o PP
- Bajo la hipótesis nula,  $\alpha$  y  $\beta$  solo se pueden estimar por MCO
- El estadístico de contraste KPSS es el siguiente, donde  $L$  se escoge a conveniencia:

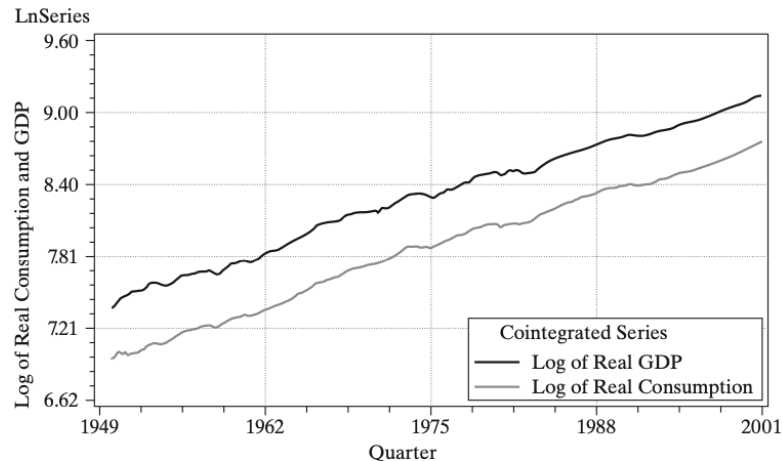
$$KPSS = \frac{\sum_{t=1}^T E_t^2}{T^2 \hat{\sigma}^2} \quad \text{where} \quad E_t = \sum_{s=1}^t e_s$$

$$\hat{\sigma}^2 = \frac{\sum_{t=1}^T e_t^2}{T} + 2 \sum_{j=1}^L \left(1 - \frac{j}{L+1}\right) r_j \quad r_j = \frac{\sum_{s=j+1}^T e_s e_{s-j}}{T}$$

- Los estudios económicos con series temporales casi siempre involucran series no estacionarias y tendenciosas, de modo que se puede utilizar la diferencia o transformaciones de otros tipos para poder obtener estacionariedad. Sin embargo, existen otros métodos más apropiados para analizar este tipo de variables
  - En un modelo completamente especificado  $y_t = \beta x_t + \varepsilon_t$ , se suele asumir que el proceso de las perturbaciones  $\varepsilon_t$  es un ruido blanco estacionario
    - No obstante, de manera general, si dos series están integradas en diferentes órdenes, entonces las combinaciones lineales estarán integradas al mayor de ambos órdenes
    - Por lo tanto, si  $y_t$  y  $x_t$  son series integradas  $I(1)$ , entonces normalmente se esperaría que  $y_t - \beta x_t$  también fuera  $I(1)$  independientemente del valor de  $\beta$ , y no  $I(0)$
    - Si  $y_t$  y  $x_t$  tienen una tasa de deriva (en la misma dirección) para sus propias tendencias, a no ser que haya una relación entre sus tendencias, la diferencia entre estas debe ser también creciente con otra tendencia, y esto causaría que la suposición sobre  $\varepsilon_t$  fuera inconsistente
  - No obstante, también puede darse el caso de que si  $y_t$  y  $x_t$  son series integradas  $I(1)$ , de modo que puede existir un valor  $\beta$  tal que  $\varepsilon_t = y_t -$

$\beta x_t$  sea  $I(0)$  (existe una diferencia concreta  $y_t - \beta x_t$  que cree una serie estacionaria)

- Intuitivamente, si las dos series son  $I(1)$ , entonces esta diferencia parcial entre ellas puede ser estable alrededor de una media fija. La implicación sería que las series van acercándose a la misma velocidad (más o menos)



- Dos series que satisfacen estos requerimientos se dicen que están cointegradas (su combinación lineal está cointegrada), y el vector  $(1, -\beta)$  (o cualquier múltiplo de este) es un vector de cointegración
  - En este caso, se distingue entre una relación a largo plazo entre  $y_t$  y  $x_t$  (la manera en que dos variables se desplazan juntas hacia arriba) y las dinámicas a corto plazo (la relación entre desviaciones de  $y_t$  de su tendencia a largo plazo y de  $x_t$  de su tendencia a largo plazo)
  - De estar cointegradas, la diferenciación de los datos sería contraproducente, dado que oscurecería la relación a largo plazo entre  $y_t$  y  $x_t$
- Si dos variables  $I(1)$  están cointegradas, entonces alguna combinación lineal ( $y_t - \beta x_t$ ) de estas variables produce una serie estacionaria  $I(0)$  (la combinación lineal está cointegrada). Esta combinación lineal no crea misteriosamente una variable estacionaria, sino que hay algo presente en las variables originales que ya no debe de estar en la nueva
- Suponiendo que hay dos variables  $I(1)$  que tienen una tendencia lineal (donde  $u_t$  y  $v_t$  es ruido blanco), una combinación lineal de  $y_{1t}$  y  $y_{2t}$  con el vector  $(1, \theta)$  produce la siguiente variable  $z_t$  que es, generalmente,  $I(1)$ :

$$y_{1t} = \alpha + \beta t + u_t \quad y_{2t} = \gamma + \delta t + v_t$$

$$\Rightarrow z_t = (\alpha + \theta\gamma) + (\beta + \theta\delta)t + u_t + \theta v_t$$

- La única manera de que  $z_t$  pueda hacerse estacionaria es que  $\theta = -\beta/\delta$ . De ser así, entonces el efecto de una combinación lineal de ambas variables es que se elimine la tendencia lineal común entre ellas
- No obstante, esto implica que la única manera de que  $y_{1t}$  y  $y_{2t}$  estén cointegradas es que tengan una tendencia común de algún tipo. Si, además de la tendencia lineal se incluyera una tendencia estocástica a través de un camino aleatorio  $w_{it} = w_{it-1} + \eta_{it}$  en donde  $\eta_t$  para  $i = 1, 2$  es ruido blanco, entonces las series solo estarían cointegradas si  $w_{1t} = w_{2t}$  (dado que una combinación lineal no implicaría caminos aleatorios diferentes, solo uno)

$$y_{1t} = \alpha + \beta t + \lambda w_t + u_t \quad y_{2t} = \gamma + \delta t + \pi w_t + v_t$$

- Se puede demostrar, sin embargo, que no es posible encontrar una combinación lineal que esté cointegrada aunque compartan la misma tendencia estocástica. Por lo tanto, el resultado final es que si  $y_{1t}$  y  $y_{2t}$  están cointegradas, entonces tienen que compartir una sola tendencia común (no dos comunes, como en este caso)
- Stock y Watson determinaron que esta última conclusión era la más importante en el caso de las variables económicas. Un conjunto de  $M$  variables que están cointegradas se puede escribir como un componente estacionario sumado a combinaciones lineales de un conjunto más pequeño de tendencias comunes. Si el rango de cointegración de un sistema es  $r$ , solo puede haber como mucho  $M - r$  tendencias lineales y  $M - r$  tendencias estocásticas comunes. El efecto de la cointegración es purgar estas tendencias comunes de las variables resultantes
- Para analizar la cointegración, el primer paso es establecer si de verdad esta característica está presente en los datos. Hay dos enfoques para ello: el enfoque de Engle-Granger y el del VAR
  - El enfoque de Engle-Granger consiste en contrastar la estacionariedad de los errores de equilibrio de estimaciones de una sola ecuación. El contraste que se utiliza en este enfoque es el contraste aumentado de Dickey-Fuller de Engle-Granger
  - Siendo  $y_t$  un conjunto de  $M$  variables que se cree que están cointegradas, el primer paso es contrastar el orden de

integración de cada una a través de un contraste de Dickey-Fuller estándar o de uno aumentado

- Si el rango de cointegración de un sistema es  $r$ , entonces hay  $r$  vectores independientes  $\gamma_i = (1, -\theta_i)$ , en donde cada vector se distingue por ser normalizado en una variables diferente. Si se supone que hay un conjunto de variables exógenas  $I(0)$ , incluyendo un término constante en el modelo, entonces cada vector cointegrante produce la siguiente relación de equilibrio:

$$y'_i \gamma_i = x'_i \beta + \varepsilon_{it} \Rightarrow y_{it} = Y'_{it} \theta_i + x'_i \beta + \varepsilon_{it}$$

- Se pueden obtener estimadores por mínimos cuadrados ordinarios de esta regresión, y si la teoría es correcta y los estimadores son consistentes, entonces los residuos de esta regresión deberían estimar los errores de equilibrio
- Bajo estas suposiciones, los residuos de la regresión son estimadores de los errores de equilibrio y, por tanto, tendrían que ser  $I(0)$ . Aunque el enfoque natural sería aplicar los contrastes de Dickey-Fuller a estos residuos, se tienen que utilizar otros valores especiales, dado que la distribución no es normal y depende del número de variables  $I(1)$  usadas como regresores