



PROBABILIDAD Y ESTADÍSTICA II

UNIDAD N°2

**Licenciatura en Enseñanza de la Matemática
Año 2011
Mg. Lucía C. Sacco**

UNIDAD N°2

Variables aleatorias bidimensionales discretas. Correlación y regresión.

Correlación. Diagramas de dispersión.

Parámetros de una distribución bidimensional. Esperanza matemática y varianza. Covarianza.

Correlación lineal. Coeficiente de correlación lineal. Coeficientes de determinación.

Regresión lineal. Recta de regresión mínimo cuadrática. Fiabilidad de la recta de regresión.

Propósitos:

Brindar oportunidades para la construcción de herramientas que permitan:

- Utilizar los diagramas de dispersión para representar conjuntos de datos de dos variables.
- Aprender el significado de correlación estadística.
- Medir la dependencia estadística con ayuda del coeficiente de correlación lineal.
- Calcular la recta de regresión y emplearla para hacer estimaciones.



Correlación

Cambios en una de las variables influyen en los cambios de la otra

Sentido

Grado

Correlación directa

Correlación inversa

Correlación débil

Correlación fuerte

Ejemplos

- La correlación entre el número de zapato y la estatura de las personas es directa y fuerte. (Las fábricas de zapatos hacen tallas de zapatos más grandes en Suecia que en Japón, pues, en general, los suecos son más altos que los japoneses. No obstante, nadie se compra los zapatos por su estatura, todo el mundo se los prueba!!).
- Las variables temperatura y el número de enfermos de gripe están inversamente correlacionadas: a menor temperatura más enfermos de gripe. Quizás se trate, también, de una correlación fuerte.
- Las variables altura y cociente intelectual de las personas no están correlacionadas.
- La correlación entre el número de errores cometidos y tiempo empleado en realizar una tarea por un grupo de personas no sabemos cómo es; para determinarla habría que tener datos concretos.



Diagrama de dispersión

El primer paso para determinar el sentido y el grado de la correlación entre dos variables consiste en representar gráficamente, en el plano cartesiano, los pares de valores conocidos. Estos gráficos, que reciben el nombre de **diagramas de dispersión**, permiten visualizar la posición de los datos en el plano. La forma de la **nube de puntos** asociada a cada diagrama nos permitirá establecer conjeturas sobre la correlación existente entre las variables estudiadas.

Si una de las variables se puede considerar como la variable que causa, o explica los cambios observados en la otra, a esa variable se la denomina **explicativa** y se la representa sobre el eje x. En este caso, a la otra variable se la denomina **variable respuesta** y se la representa sobre el eje y. Si no se quiere distinguir entre variable explicativa y variable respuesta, cualquiera de las dos puede representarse en el eje de las abscisas.

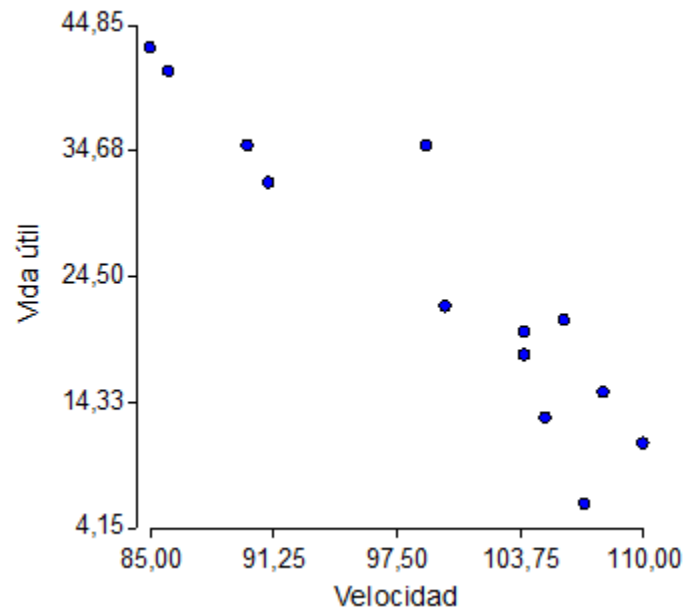
Ejemplo 1

Los siguientes datos corresponden a la vida útil y a la velocidad de corte de una herramienta:

Velocidad de corte:	86	104	100	85	107	104	106	99	90	110	108	105	91
Vida útil:	41	18	22	43	6	20	21	35	<u>35</u>	11	15	13	32



Diagrama de dispersión



En general, dependiendo de la forma de la nube de puntos, puede asegurarse:

- Una nube de puntos alargada indica **correlación lineal**: los puntos se distribuyen en torno a una línea recta. La estrechez de la nube expresa que la correlación es fuerte.
- Si la recta que se ajusta a la nube tiene pendiente positiva, **la correlación será directa**; al crecer la variable X, lo hace también la variable Y.
- Una recta con pendiente negativa, indica que la **correlación es inversa**, al crecer X, disminuye Y.



Diagrama de dispersión

Actividades:

1. Ocho personas, con similar destreza en mecanografía, teclearon 40 líneas de texto en un ordenador. El tiempo empleado, en minutos, y el número de errores cometidos, fueron:

Tiempo (X)	6	7	8	9	9	10	12	12
Errores (Y)	22	15	12	17	21	13	9	6

¿Qué tipo de correlación se da entre las variables estudiadas?

2. Los siguientes datos corresponden al consumo de combustible de un auto a medida que aumenta su velocidad.

Velocidad	10	20	30	40	50	60	70	80	90	100	120	130	140	150
Consumo	21	13	10	8	7	5.9	6.3	6.95	7.57	8.27	9.87	10.79	11.77	12.83

- Dibuja un diagrama de dispersión. ¿Cuál consideras que es la variable explicativa?
- Describe la forma de la relación.



Parámetros de una distribución bidimensional

Media de cada una de las variables (centro geométrico)

$$\bar{x} = \frac{\sum x_i}{n} \quad \bar{y} = \frac{\sum y_i}{n}$$

Varianza y desviaciones típicas

$$s_x^2 = \frac{\sum (x_i - \bar{x})^2}{n} = \frac{\sum x_i^2}{n} - \bar{x}^2$$
$$s_y^2 = \frac{\sum (y_i - \bar{y})^2}{n} = \frac{\sum y_i^2}{n} - \bar{y}^2$$

s_x y s_y

Covarianza

$$s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n} \Leftrightarrow \frac{\sum x_i y_i}{n} - \bar{x} \bar{y}$$

- Si $s_{xy} > 0$, la correlación es directa.
- Si $s_{xy} < 0$, la correlación es inversa.



Parámetros de una distribución bidimensional

Actividad:

3. Si consideramos la distribución que mide la altura de niños en cm de diferentes años:

Años (X)	1	2	3	6
Estatura (Y)	63	80	88	101

- Realizar el diagrama de dispersión
- Hallar su covarianza.
- ¿Qué tipo de correlación presentan las variables años y estatura? ¿qué implica esto?



Coeficiente de correlación lineal

El valor del **coeficiente de correlación lineal** es el criterio que se utiliza para medir la fuerza de la correlación entre dos variables.

Este coeficiente, denotado ρ , se define así
$$\rho = \frac{s_{xy}}{s_x s_y}$$

Esto es, la razón entre la covarianza de las variables X e Y y el producto de sus desviaciones típicas marginales.

Las propiedades fundamentales del coeficiente de correlación son:

1. El valor de ρ no cambia al hacerlo la escala de medición, pues la covarianza y el producto de las desviaciones típicas varían en la misma proporción.
2. El signo de ρ es el mismo que el de la covarianza, pues las desviaciones típicas siempre son positivas. Luego:
 - Si $\rho > 0$, la correlación es directa.
 - Si $\rho < 0$, la correlación es inversa.



Coeficiente de correlación lineal

3. El valor de ρ está entre -1 y $+1$: $-1 \leq \rho \leq 1$
4. Si ρ toma valores cercanos a 0 por izquierda, la correlación es débil (e inversa).
5. Si ρ toma valores cercanos a $+1$, la correlación es fuerte (y directa).
El signo de ρ no determina la fuerza de correlación: sólo el sentido. Tampoco indica la mayor o menor pendiente de la recta asociada a la nube de puntos.
6. Si $|\rho| = 1$, la correlación es perfecta. Hay dependencia lineal entre las variables X e Y .
7. Si ρ toma valores cercanos a 0 , la correlación es débil.

Actividades:

4. Hallar el coeficiente de correlación entre el tiempo empleado y el número de errores cometidos por ocho personas al realizar un trabajo de mecanografía (Actividad 1).
5. Hallar el coeficiente de correlación de la distribución dada por la siguiente tabla:

X	4	7	3	9
Y	3	6	7	5



Coeficiente de determinación

Su valor indica la proporción de la variación en la variable Y que puede ser explicada por los cambios de la variable X

Ejemplos

- a) El coeficiente de correlación entre la edad y la altura de niños (Actividad 3), vale $\rho = 0,94$. Por tanto, el coeficiente de determinación será $\rho^2 = 0,94^2 = 0,88$. Esto significa que, en los niños de nuestro ejemplo, el 88% de su altura se explica por la edad; el resto, hasta el 100%, será debido a otras causas: altura de sus padres, dieta, etc.
- b) El coeficiente de determinación entre el tiempo empleado y el número de errores cometidos por ocho personas (Actividad 1), vale $\rho^2 = (-0,764)^2 = 0,584$. O sea, el tiempo empleado explica el 58,4% de las diferencias de errores. El 41,6% de la variación restante se debe a causas desconocidas por nosotros.



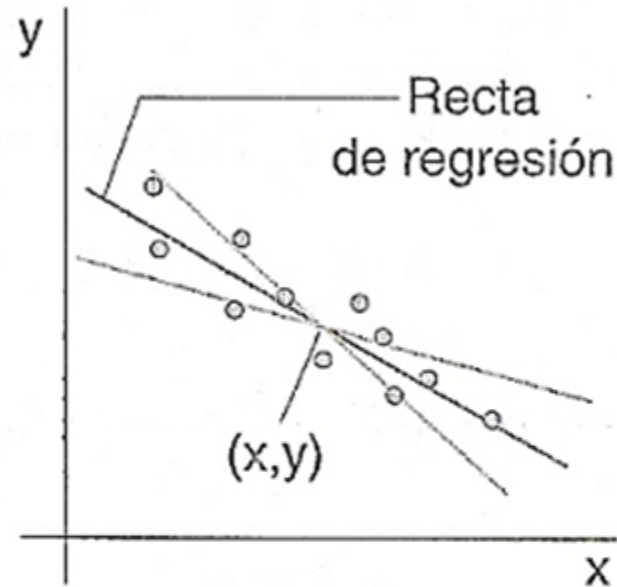
Regresión lineal

Recta de regresión mínimo cuadrática

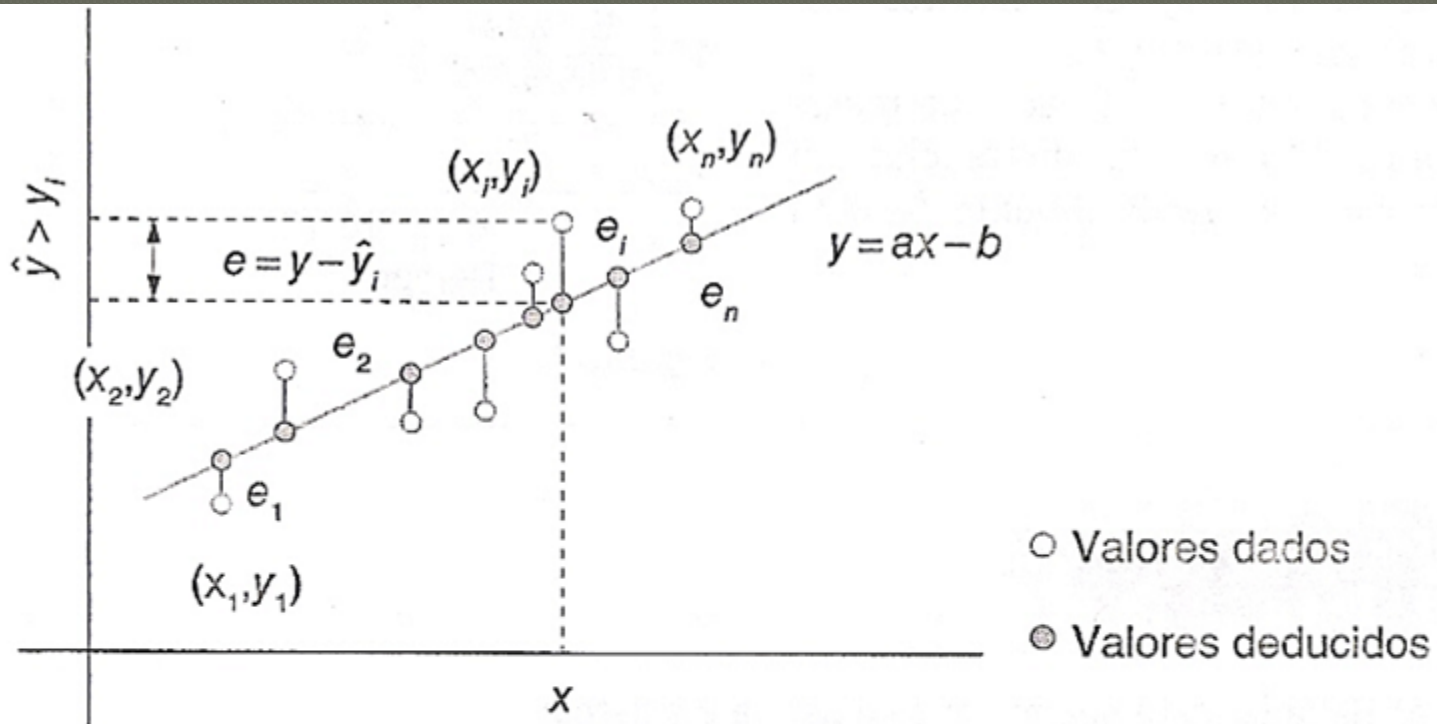
Al ser una recta ideal no tiene por que pasar por ninguno de los puntos dados, pero sí lo más cerca posible de todos ellos. De cualquier manera, siempre se cometerán errores en la estimación. Obviamente, lo que se pretende es que estos errores sean lo más pequeños posibles.

La recta que mejor se ajusta a estos propósitos es la **recta de regresión mínimo cuadrática**, que es aquella que minimiza la suma de los cuadrados de los errores. Esto es, si la ecuación de esta recta es $y = ax + b$, y los puntos dados son $(x_1, y_1), (x_2, y_2), \dots$ y (x_n, y_n) , debe cumplirse que la suma $e_1^2 + e_2^2 + \dots + e_n^2$ sea mínima, siendo e_i la diferencia entre el valor dado, y_i , correspondiente a x_i , y el valor deducido \hat{y}_i por la recta para ese mismo x_i :

$$\hat{y}_i = ax_i + b \quad \text{y} \quad e_i = |\hat{y}_i - y_i|$$



Regresión lineal



Con estas condiciones, los valores de la pendiente a y de la ordenada al origen b de esa recta valen:

$$a = \frac{s_{xy}}{s_x} \quad \text{y} \quad b = \bar{y} - \frac{s_{xy}}{s_x} \bar{x}$$

Luego, **la ecuación de la recta de regresión** es:

$$y - \bar{y} = \frac{s_{xy}}{s_x} (x - \bar{x})$$



Regresión lineal

Actividades:

6. Con los datos de la Actividad 1:

- hallar la ecuación de la recta de regresión que permita estimar los errores a partir del tiempo.
- representar la recta y los puntos dados.
- estimar el número de errores de una persona que tardase 11 minutos en teclear las 40 líneas.
- Ídem para otra persona que tardase 9 minutos. Compararlo con datos del problema.

7. a) Hallar la recta que mejor se ajuste a la distribución dada por la siguiente tabla:

X	1	3	4	5	6
Y	3	4	6	6	8

- mediante la recta obtenida, estimar el valor para Y para $x=2$ y $x=7$.



Fiabilidad de la recta de regresión

La fiabilidad de las estimaciones hechas a partir de la recta de regresión depende fundamentalmente de:

1. El valor del coeficiente de correlación ρ

Una correlación alta (ρ próximo a 1), asegura estimaciones fiables

2. El número de datos considerados

La fiabilidad aumenta al aumentar los datos. Una recta obtenida a partir de pocos datos genera grandes riesgos, aunque ρ sea muy alto.

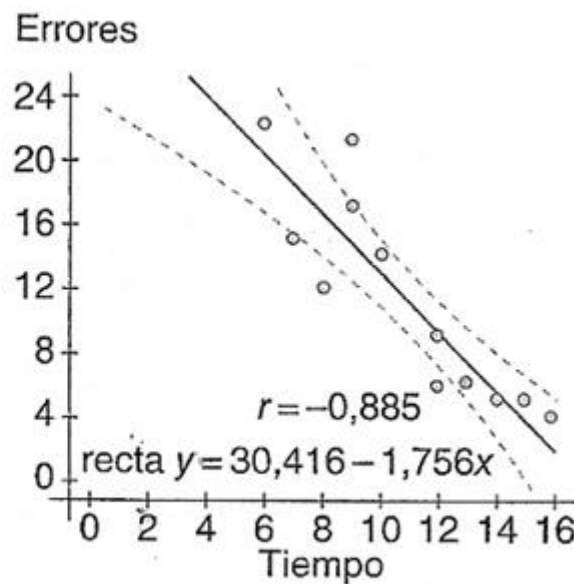
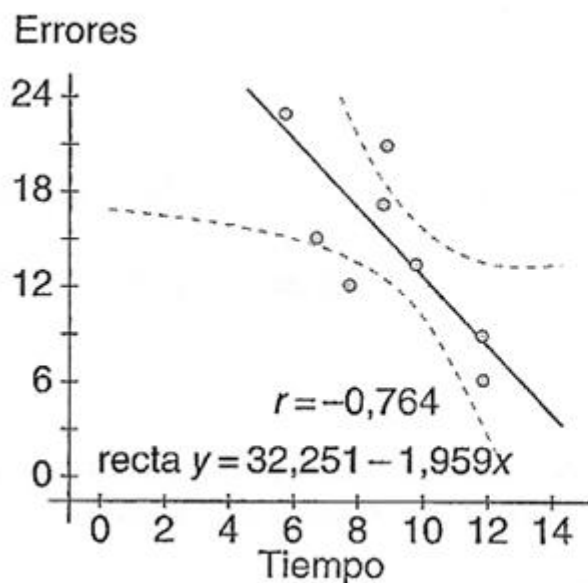
3. La proximidad del valor x_0 , para el que se quiere hacer la estimación, a la media \bar{x}

La estimación de y_0 para un x_0 dado es más fiable cuando x_0 está próximo a \bar{x} ; a medida que x_0 se aleja de \bar{x} la estimación se hace más arriesgada.



Fiabilidad de la recta de regresión

Se observa en las siguientes figuras como las líneas de puntos determinan una banda alrededor de la recta de regresión. Esta banda indica los márgenes del valor estimado \hat{y}_0 , para cada x_0 dado. En este caso, la banda se ha generado para una probabilidad de acierto del 95%.



Fiabilidad de la recta de regresión

Observaciones sobre la recta de regresión

1. La banda se ensancha a medida que nos alejamos de la media \bar{x} . Esto indica que si la estimación desea hacerse para valores alejados de la media, la probabilidad de acierto es menor.
2. La banda se ensancha más rápidamente cuando el coeficiente de correlación es menor.
3. Para la misma distribución, la estimación será más fiable si aumentan los datos de la muestra considerada.

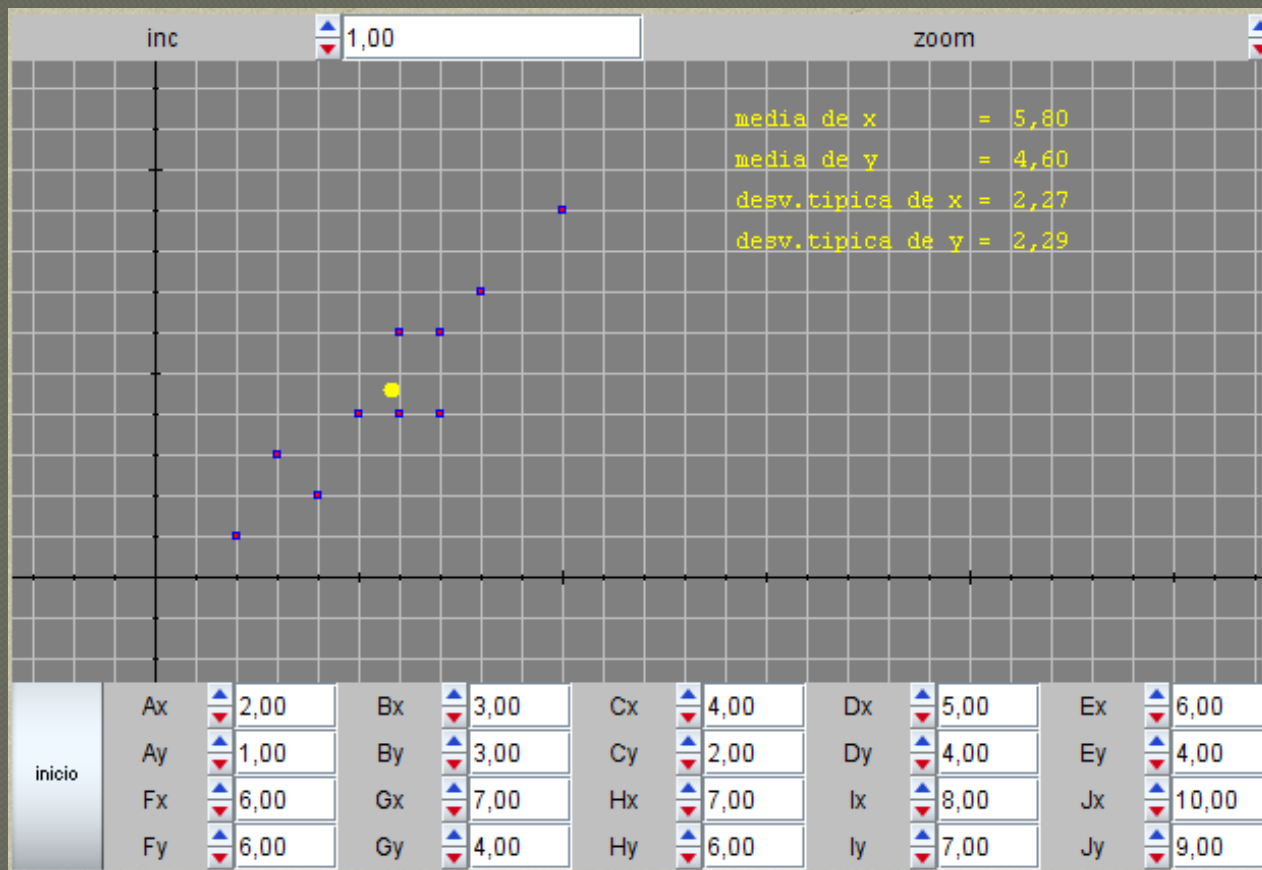
Limitaciones de la recta de regresión

La recta de regresión debe usarse para hacer estimaciones en valores próximos a los considerados. Pretender una estimación en puntos lejanos puede conducir a soluciones absurdas. Por ejemplo, si con la recta que obtuvimos para la Actividad 1 estimamos el número de errores que cometería una persona que tardase 30 minutos en teclear 40 líneas, se obtendría:

$$y = 32,251 - 1,959 \cdot 30 = -26,519 \quad \text{¿-27 errores? ¡Absurdo!}$$



Distribuciones bidimensionales



http://recursostic.educacion.es/descartes/web/materiales_didacticos/distrib_bidimensionales/distribuciones_bidimensionales.htm



Situaciones problemáticas

<http://www.diazdesantos.es/wwwdat/pdf/SP0410004046.pdf>