



# Regresión Lineal\_

Sesión Presencial 1



## Alcances de la lectura asignada

- Reconocer la terminología asociada a la modelación estadística.
- Conocer la regresión lineal y sus fundamentos.
- Interpretar los parámetros estimados en la regresión.
- Conocer y ser capaz de interpretar estadísticos de bondad de ajuste y coeficientes.
- Reconocer los supuestos en los que la regresión tiene sustento teórico.
- Implementar un modelo de regresión con *statsmodels*
- Utilizar transformaciones simples en las variables independientes.
- Implementar un modelo predictivo con *scikit-learn*

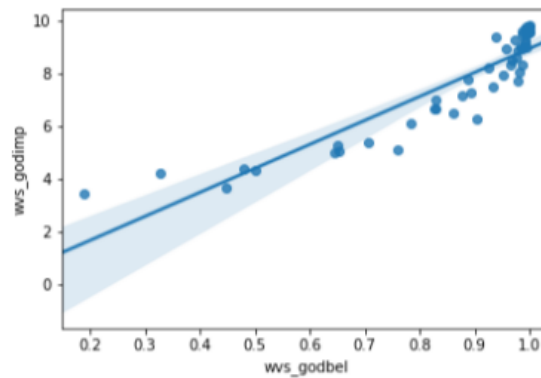
## Activación de Conceptos

- En la unidad anterior aprendimos sobre la inferencia estadística y correlación.
- ¡Pongamos a prueba nuestros conocimientos!

Aproximadamente, ¿Cuál es la correlación entre ambas variables?

```
In [2]: sns.regplot(df['wvs_godbel'], df['wvs_godimp'])
```

```
Out[2]: <matplotlib.axes._subplots.AxesSubplot at 0x103bb09e8>
```



{desafío}  
latam\_

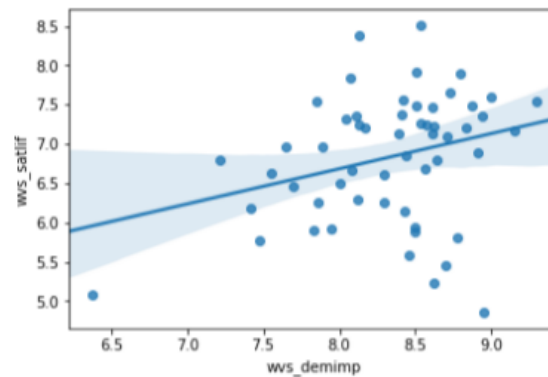
1. 0.14

- -0.3
- 0.9
- No hay suficiente información.

Aproximadamente, ¿Cuál es la correlación entre ambas variables?

```
In [3]: sns.regplot(df['wvs_demimp'], df['wvs_satlif'])
```

```
Out[3]: <matplotlib.axes._subplots.AxesSubplot at 0x1a0b6a7a90>
```



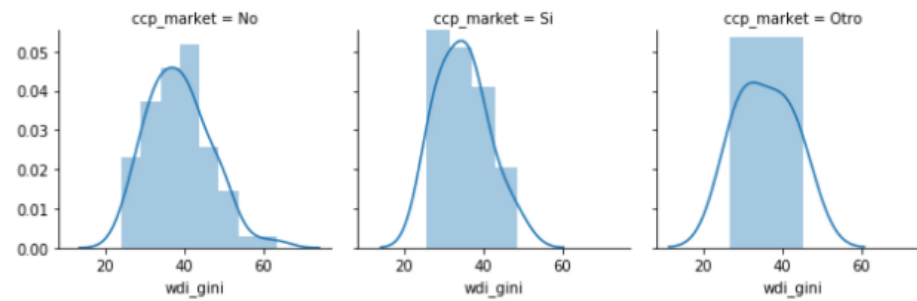
{desafío}  
latam\_

1. 0.20

- -0.7
- .6
- No hay suficiente información.

¿Cuál es la mejor manera para realizar el siguiente gráfico?

In [5]: `question()`



A

```
for i, e in enumerate(ccp_market.unique()):  
    plt.subplot(1, 3, i+1)  
    plt.hist(df[df['ccp_market']==e]['wdi_gini'])
```

B

```
g = sns.PairPlot(df['ccp_market'], df['wdi_gini'])
```

C

```
g = sns.FacetGrid(df, col='ccp_market', col_wrap=3)  
g = g.map(sns.distplot, 'wdi_gini')
```

D

Ninguna de las anteriores

## ¿Qué son los P-Value?

---

1. Es nuestro nivel de confianza en encontrar una hipótesis alternativa
  - Es la probabilidad de encontrar resultados observados más extremos cuando  $H_0$  es verdadera.
  - Es el puntaje de corte entre la región de aceptación y rechazo
  - Ninguna de las anteriores.



Si comparamos la media del índice de Gini frente a una constante de 38 y obtenemos el siguiente resultado.  
¿Qué podemos decir?

---

```
In [6]: stats.ttest_1samp(df['wdi_gini'].dropna(), 38)
```

```
Out[6]: Ttest_1sampResult(statistic=-0.7331483806078863, pvalue=0.46497285823686796)
```

1. La media del índice de Gini no es estadísticamente distinta de 38
- La media del índice de Gini es estadísticamente distinta de 38

Si comparamos el nivel de creencia en Dios entre países del medio oriente y el resto del mundo, obtenemos el siguiente resultado. ¿Qué podemos decir?

```
In [7]: df['middle_east'] = np.where(df['ht_region']==3, 1, 0)
stats.ttest_ind(df.query('ht_region == 3')['wvs_godbel'].dropna(),
                df.query('ht_region !=3')['wvs_godbel'].dropna())

Out[7]: Ttest_indResult(statistic=2.054023852143011, pvalue=0.0454412759545386
6)
```

- 
1. Las medias entre las regiones son estadísticamente distintas al 95% de confianza
    - La media de los países del medio oriente es mayor que la media del resto del mundo.
    - Tenemos un 95% de confianza en que las medias son distintas.

## Regresión Lineal (Desde la Econometría)

## Rudimentos de la Regresión

- Desde la econometría, la regresión lineal implica atributos explicativos a un fenómeno en específico. Su nombre viene de la *ecuación de la recta*:

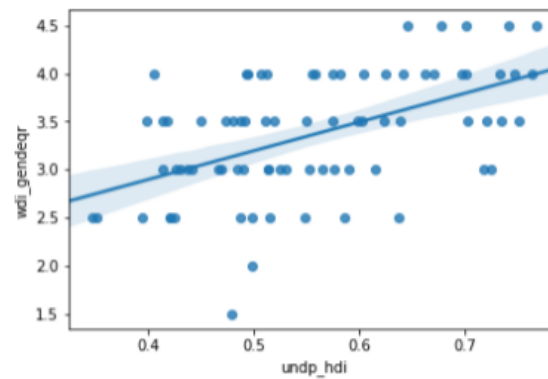
$$y_i = \beta_0 + \beta_1 \cdot X_i + \varepsilon_i$$

- Donde buscamos un punto de partida para nuestra función lineal ( $\beta_0$ ).
- Y una pendiente que indica su movimiento cuando X incrementa en 1 unidad ( $\beta_1 \cdot X$ )
- A diferencia de la ecuación de la recta, acá agregamos un error ( $\varepsilon_i$ )

## La recta de ajuste cuando $\beta_1 > 0$

```
In [8]: sns.regplot(df['undp_hdi'], df['wdi_gendeqr'])
```

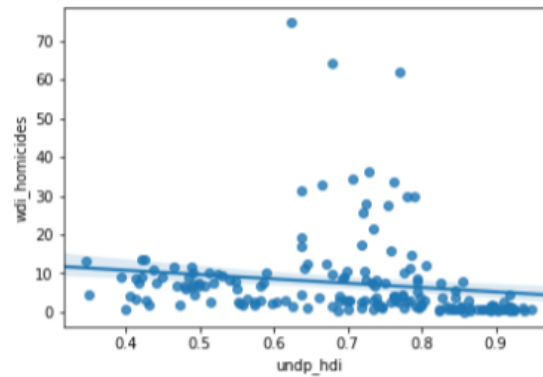
```
Out[8]: <matplotlib.axes._subplots.AxesSubplot at 0x1a17b40438>
```



## La recta de ajuste cuando $\beta_1 < 0$

```
In [9]: sns.regplot(df['undp_hdi'], df['wdi_homicides'])
```

```
Out[9]: <matplotlib.axes._subplots.AxesSubplot at 0x1a178a54e0>
```



{desafío}  
latam\_

## La regresión en Python

- Para implementar la regresión, utilizamos `statsmodels`. Si deseamos ver el efecto del desarrollo humano en las tasas de equidad de género:

$$\text{wdi\_genderqr} = \beta_0 + \beta_1 \cdot \text{undp\_hdi} + \varepsilon_i$$

**{desafío}**  
latam\_

*Importamos la librería*

```
In [10]: import statsmodels.api as sm  
import statsmodels.formula.api as smf
```

*Generamos un objeto con nuestra ecuación*

```
In [11]: model1 = smf.ols('wdi_gendeqr ~ undp_hdi', df)
```

*Ordenamos estimar el modelo*

```
In [12]: model1 = model1.fit()
```

**{desafío}**  
latam\_



## ¿Qué significa todo esto?: Bondad de ajuste

```
In [13]: results = model1.summary()
```

```
In [14]: results.tables[0]
```

Out[14]: OLS Regression Results

<b>Dep. Variable:</b>	wdi_gendeqr	<b>R-squared:</b>	0.269
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.259
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	27.92
<b>Date:</b>	Tue, 30 Oct 2018	<b>Prob (F-statistic):</b>	1.17e-06
<b>Time:</b>	16:53:08	<b>Log-Likelihood:</b>	-63.048
<b>No. Observations:</b>	78	<b>AIC:</b>	130.1
<b>Df Residuals:</b>	76	<b>BIC:</b>	134.8
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

## ¿Qué significa todo esto?: Parámetros estimados

In [15]: results.tables[1]

Out[15]:

	coef	std err	t	P> t	[0.025	0.975]
Intercept	1.6987	0.320	5.303	0.000	1.061	2.337
undp_hdi	2.9911	0.566	5.284	0.000	1.864	4.118

**{desafío}**  
latam\_

## ¿Y cómo llegamos a esto?

- Para estimar los parámetros utilizamos el **método de mínimos cuadrados**, donde buscamos resolver:

$$\begin{aligned}\beta &= \arg \min_{\beta} \mathbb{E} [(y_i - X_i' \beta)^2] \\ &= \sum_{i=0}^N (Y_i - (\beta_0 + \beta_1 X))^2\end{aligned}$$

- Este método será **ELIO** cuando se satisfagan las condiciones *Gauss-Markov*.