

# Introducción a la Ciencia de Datos

# ENTRENAR Y TESTEAR

1 

```
from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test= train_test_split(X,y,test_size = 0.15,random_state=1)
```

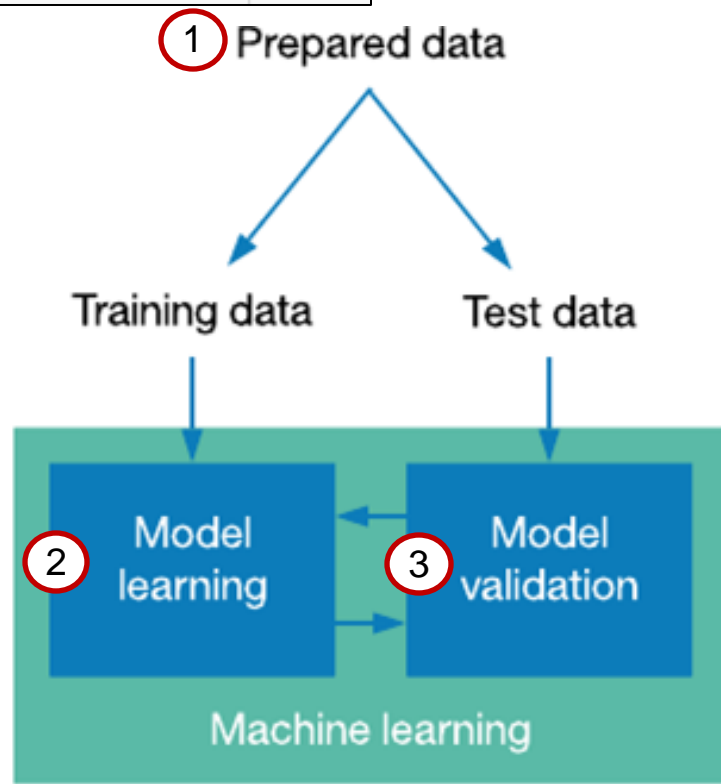
2 

```
logreg.fit(X_train, y_train)
```

3 

```
logreg_pred = logreg.predict(X_test)  
from sklearn.metrics import accuracy_score  
accuracy_score(y_test,logreg_pred)
```

Recuerda que los datos del **entrenamiento** deben tener las **mismas características** que los datos del **test**: número de variables predictoras, tipo de dato, etc.

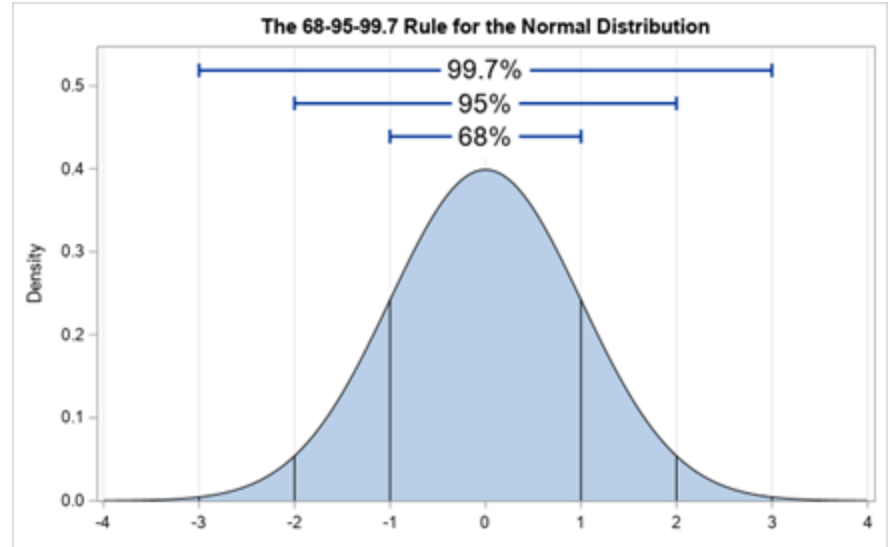


# DISTRIBUCION NORMAL

También llamada distribución Gaussiana.

Nos permite **modelar y trabajar bajo supuestos** de que nuestros datos **siempre van a tender a la media**.

Varios modelos como los de **regresión** asumen que los datos tienen una distribución normal, y es nuestro deber transformar los datos para ello.

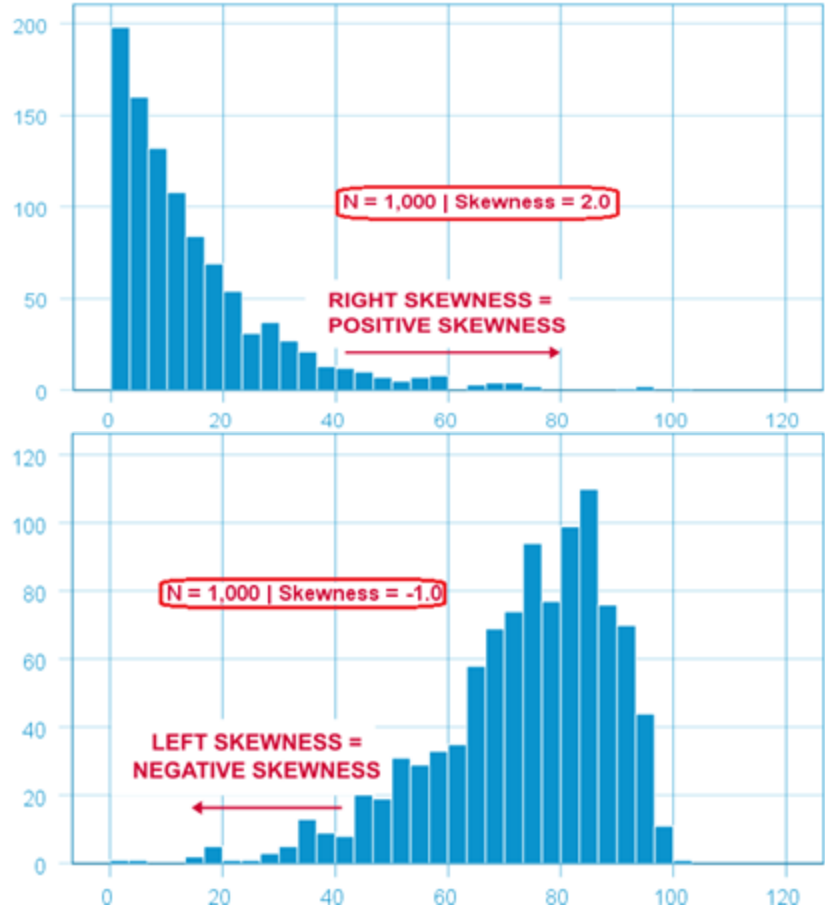


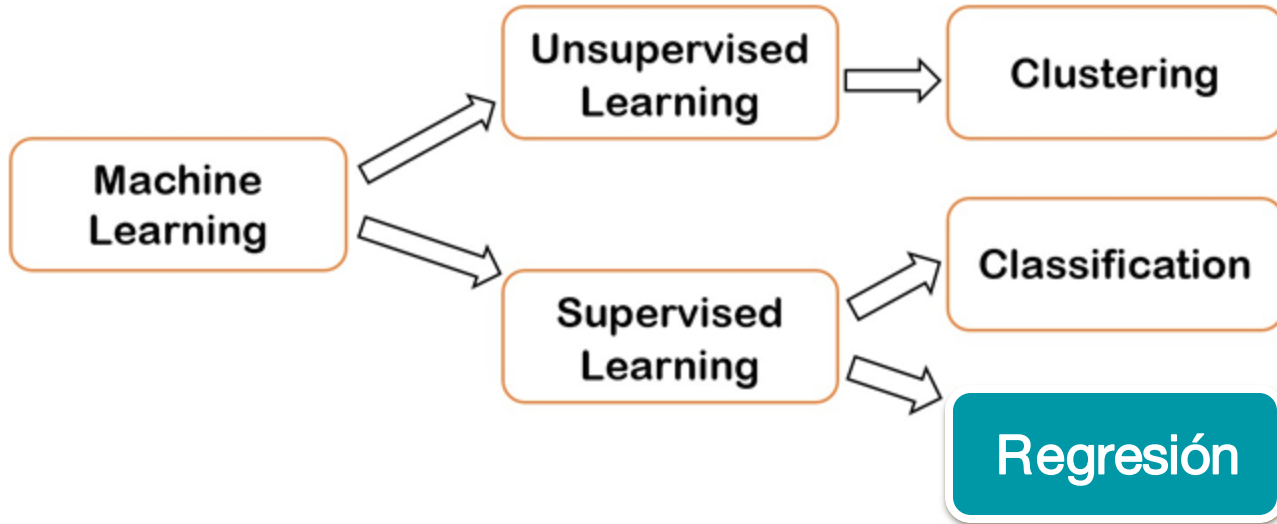
# SKEWNESS

Skewness es una variable matemática que cuantifica una **asimetría** en la distribución gráfica de los datos.

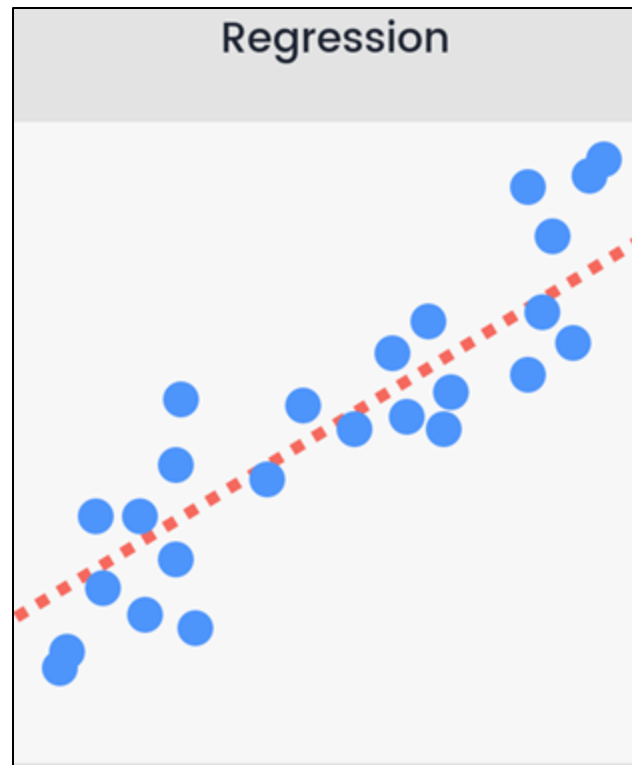
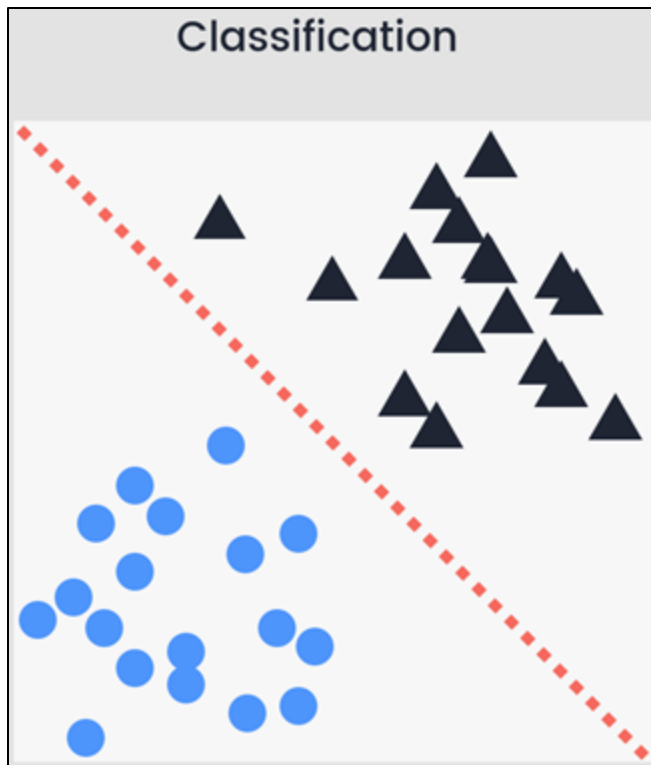
Esta medida **evidencia la presencia de outliers** debido a que **siempre ocasiona una asimetría** a la derecha (positiva) o a la izquierda (negativa).

Una medida para **corregir** esta asimetría es **aplicar un logaritmo** sobre toda variable que presente un **skewness mayor a 1** (valor absoluto).





# MODELOS DE APRENDIZAJE SUPERVISADO



# Aprendizaje Supervisado

Modelos de Regresión

## CASO PRÁCTICO

Supongamos que deseo vender un diamante porque mi abuela me dejó como herencia un anillo engarzado con un **diamante de 1,35 quilates**, y quiero tener una idea de cuánto me pagarán.

Tomo un lápiz y un cuaderno, voy a una joyería y escribo el **precio de todos los diamantes de la vitrina** y cuántos quilates tienen.

<u>Carats</u>	<u>price</u>
1.01	7,366
.49	985
.31	544
1.51	9,140
.37	493
.73	3,011
1.53	11,413
.56	1,814
.41	876
.74	2,690
.63	1,190
.6	4,172
2.06	11,764
1.1	4,682
1.31	6,171

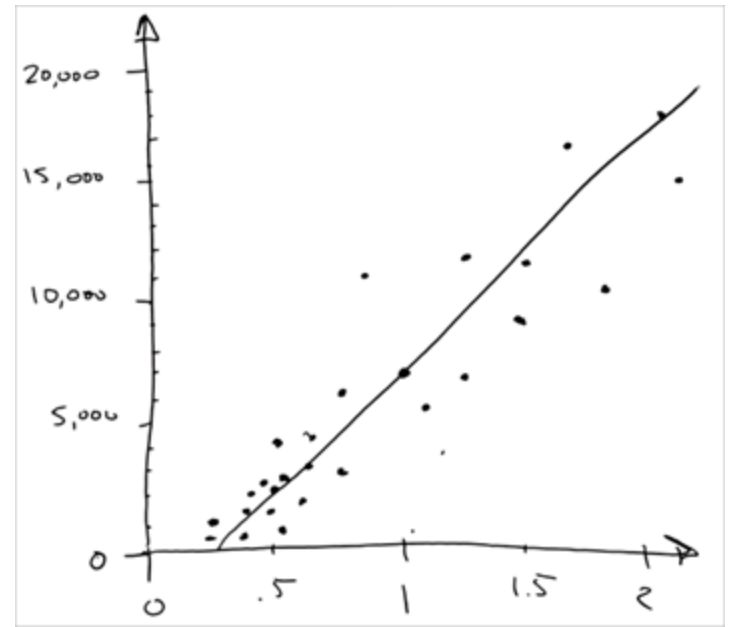
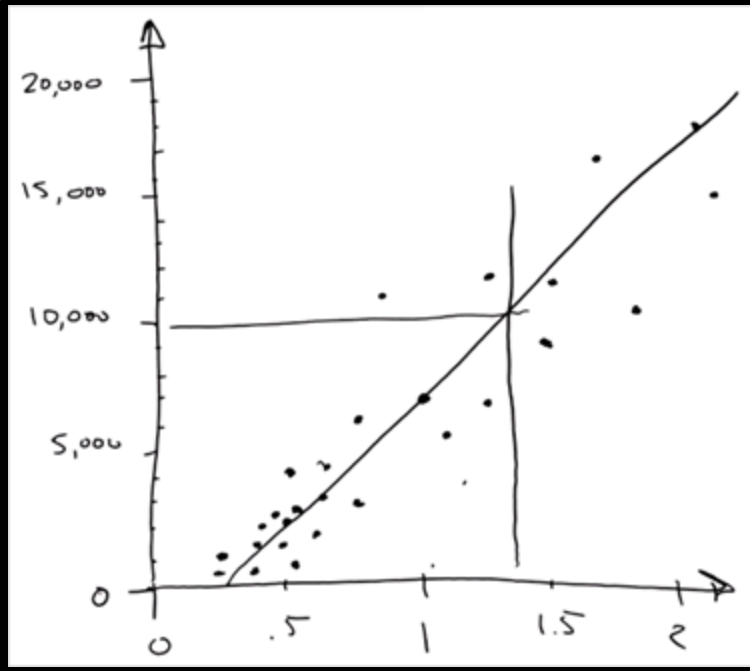


## CASO PRÁCTICO

Ahora plantearemos nuestra pregunta de forma directa: ¿cuánto costará comprar un diamante 1,35 quilates?

Nuestra lista **no contiene ningún diamante de 1,35 quilates**, pero podemos utilizar el **resto de nuestros datos para obtener una respuesta**.

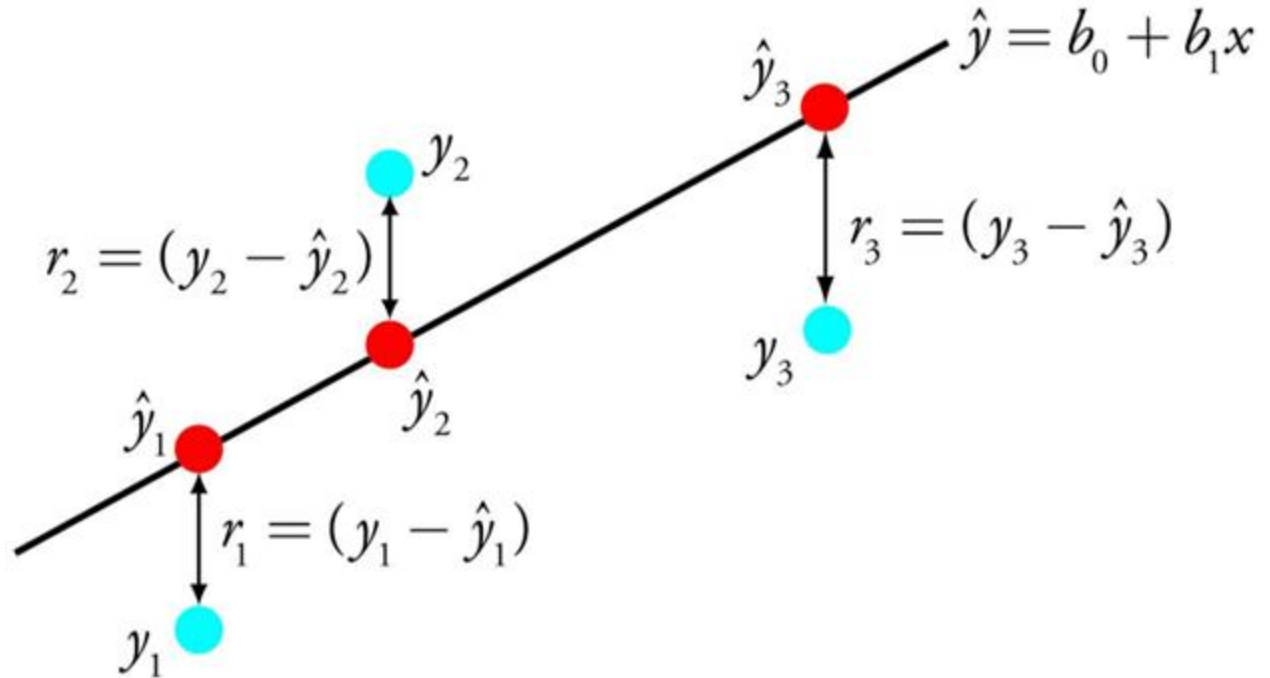
Así que procedemos a dibujar un plano con dos variables: quilates en el eje X y el precio en el eje Y. Luego **ubicamos nuestros datos** sobre el plano según los valores  $x$ ,  $y$  que les correspondan. Finalmente **trazamos una línea recta**



Podemos concluir por la gráfica que nos pagarán por el diamante de 1.35 quilates como máximo 10,000 dólares

# REGRESIÓN LÍNEAL

- Es óptimo para problemas de regresión **con datos de distribución normal**
- Busca **minimizar el error** (función de coste)



# COEFICIENTES DE REGRESIÓN

$$\boxed{\hat{y}} = \beta_0 + \beta_1 \boxed{x}$$

Valor predicho      Coeficientes      Entradas

**Función de coste:**

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y})^2$$

Regresión lineal: múltiples variables

$$\boxed{\hat{y}} = \beta_0 + \beta_1 \boxed{x_1} + \dots + \beta_p \boxed{x_p}$$

# FUNCIONES DE COSTE

$$L_{OLS}(\hat{\beta}) = \sum_{i=1}^n (y_i - x_i' \hat{\beta})^2 = \|y - X\hat{\beta}\|^2$$

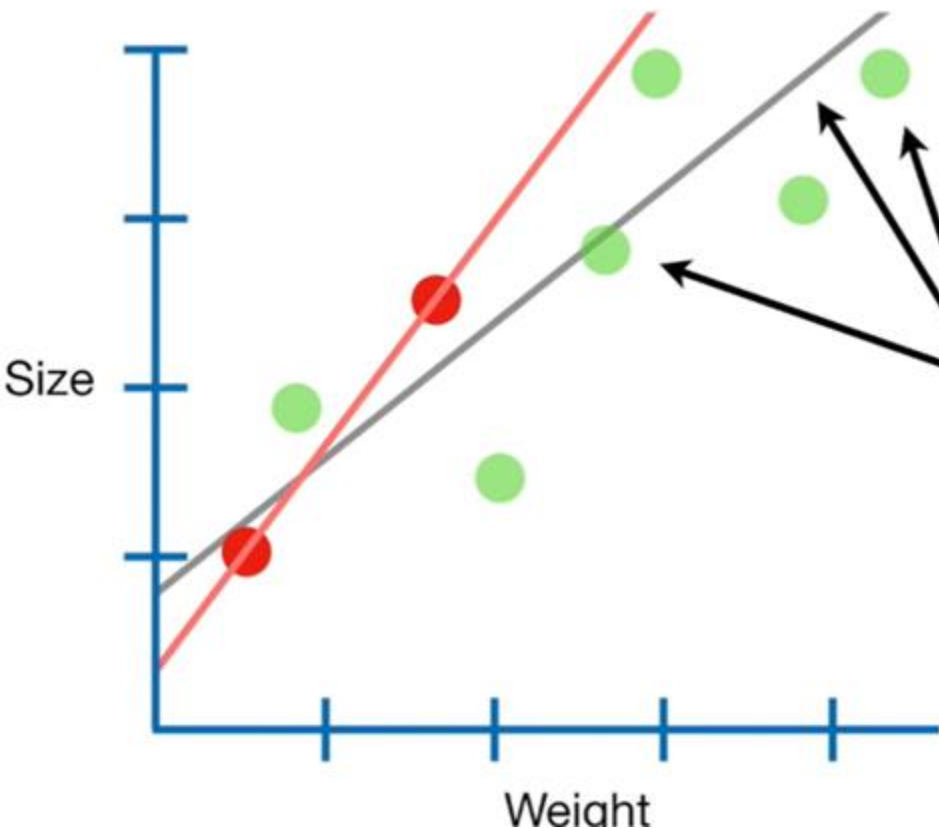
$$L_{ridge}(\hat{\beta}) = \sum_{i=1}^n (y_i - x_i' \hat{\beta})^2 + \lambda \sum_{j=1}^m \hat{\beta}_j^2 = \|y - X\hat{\beta}\|^2 + \lambda \|\hat{\beta}\|^2$$

$$L_{lasso}(\hat{\beta}) = \sum_{i=1}^n (y_i - x_i' \hat{\beta})^2 + \lambda \sum_{j=1}^m |\hat{\beta}_j|$$

- Mayor información en la siguiente ruta URL:

<https://www.analyticsvidhya.com/blog/2016/01/ridge-lasso-regression-python-complete-tutorial/>

# REGRESIÓN RIDGE



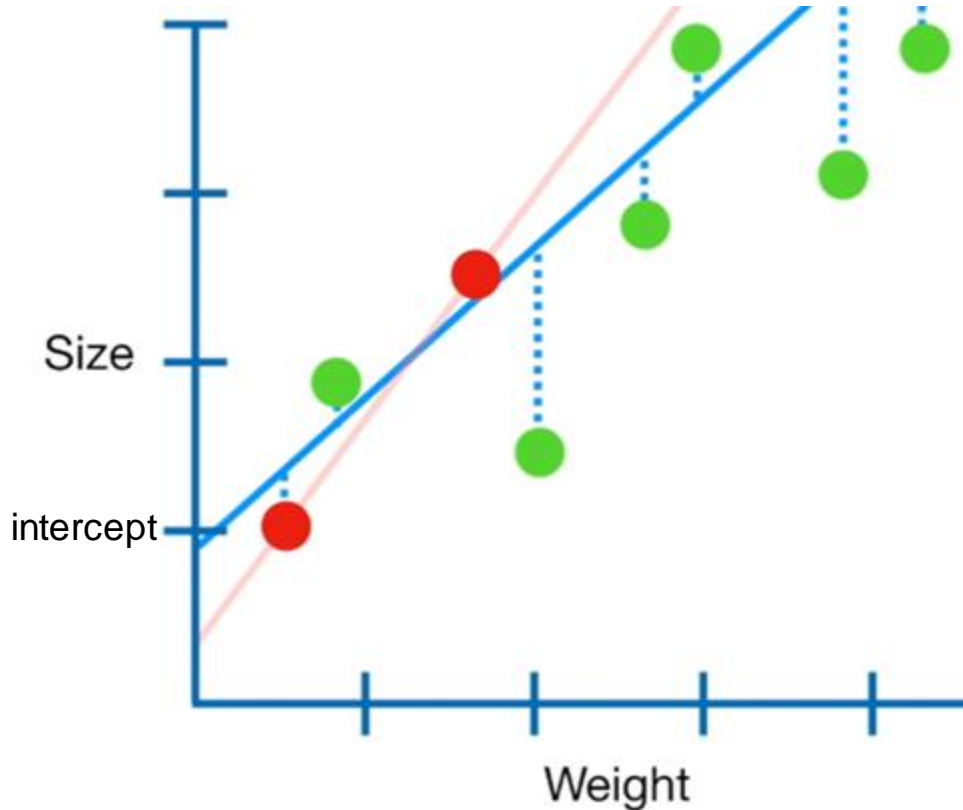
La primera aproximación son

**dos puntos aleatorios**  
De esta manera obtenemos un error  
**respecto a estos dos puntos = 0**

Sin embargo el error acumulado  
**respecto al resto de los puntos es muy grande.**

Entonces el objetivo principal de la Regresión Ridge es encontrar una línea recta diferente y más cercano al total de los puntos

# REGRESIÓN RIDGE



En la Regresión Ridge se busca **minimizar el error total** considerando también el  $\lambda$  y el coeficiente<sup>2</sup>

**Size** = y-axis intercept + slope  $\times$  **Weight**

...it minimizes...

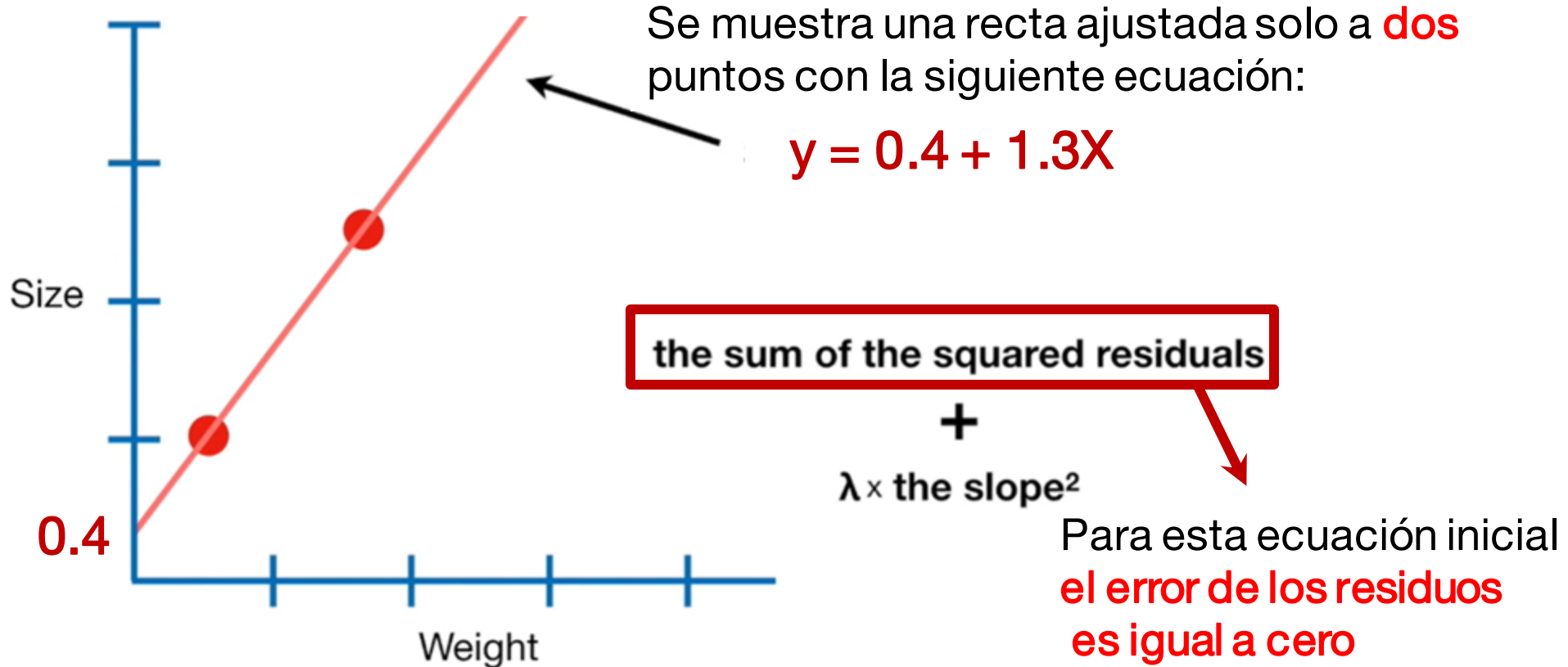
the sum of the squared residuals

+

$\lambda \times \text{the slope}^2$

“slope” es el **coeficiente** que acompaña a cada variable (X) de la ecuación

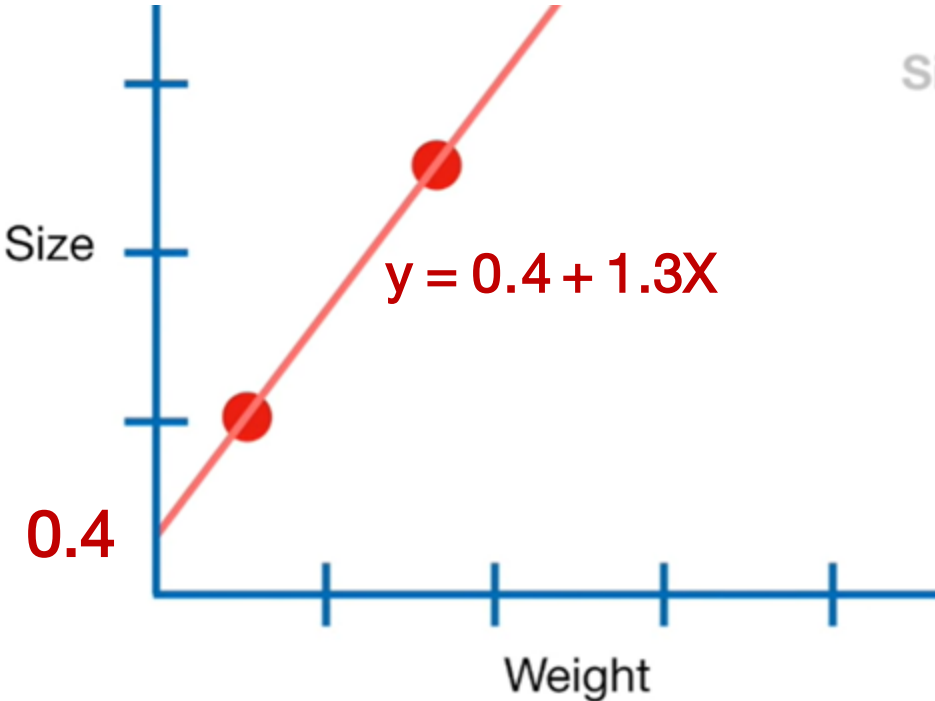
# REGRESIÓN RIDGE





# REGRESIÓN RIDGE

el coeficiente de la variable **Weight** es igual a 1.3



$$\text{Size} = 0.4 + \boxed{1.3} \times \text{Weight}$$

0

+

$$\lambda \times \boxed{1.3^2}$$

0

+

$$= 0 + 1.69 = \boxed{1.69}$$

para este ejemplo **lambda** es igual a 1  $\rightarrow 1 \times 1.3^2$

# REGRESIÓN RIDGE

Ahora creamos una nueva recta con nuevo intercepto y un nuevo coeficiente

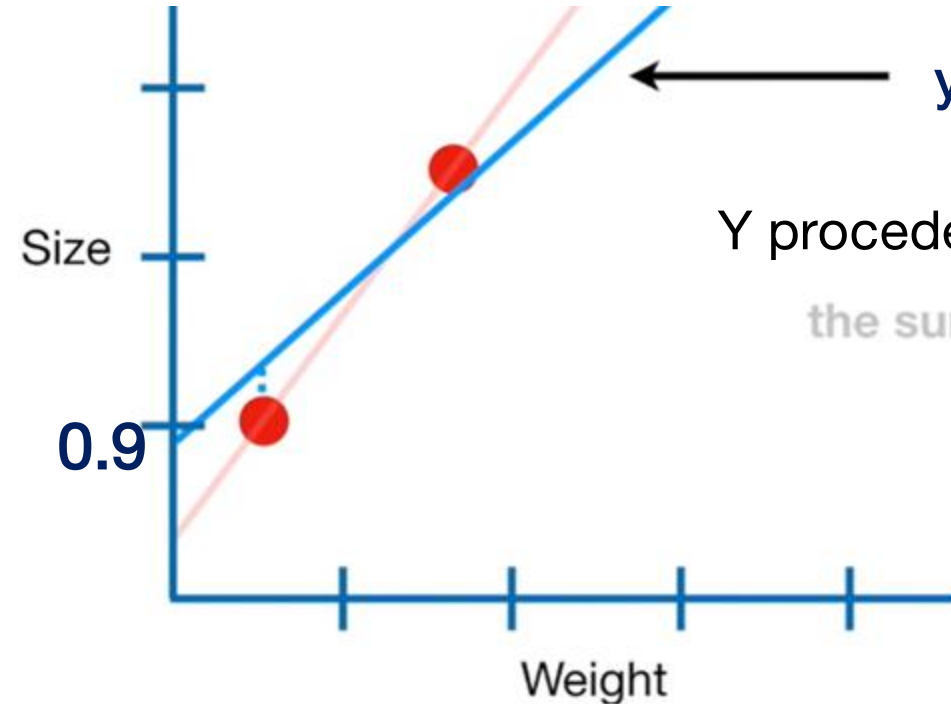
$$y = 0.9 + 0.8X$$

Y procedemos a calcular un nuevo error total:

the sum of the squared residuals

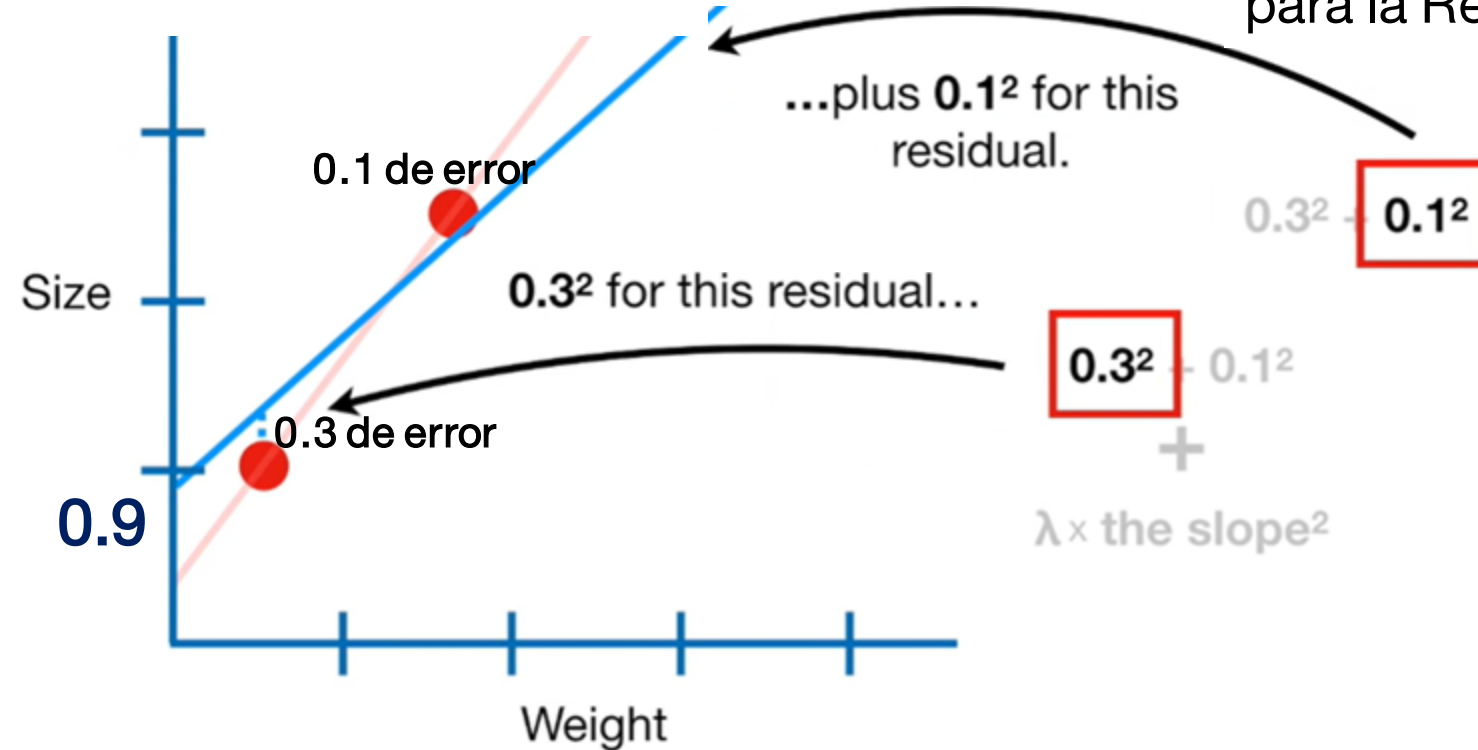
+

$\lambda \times \text{the slope}^2$

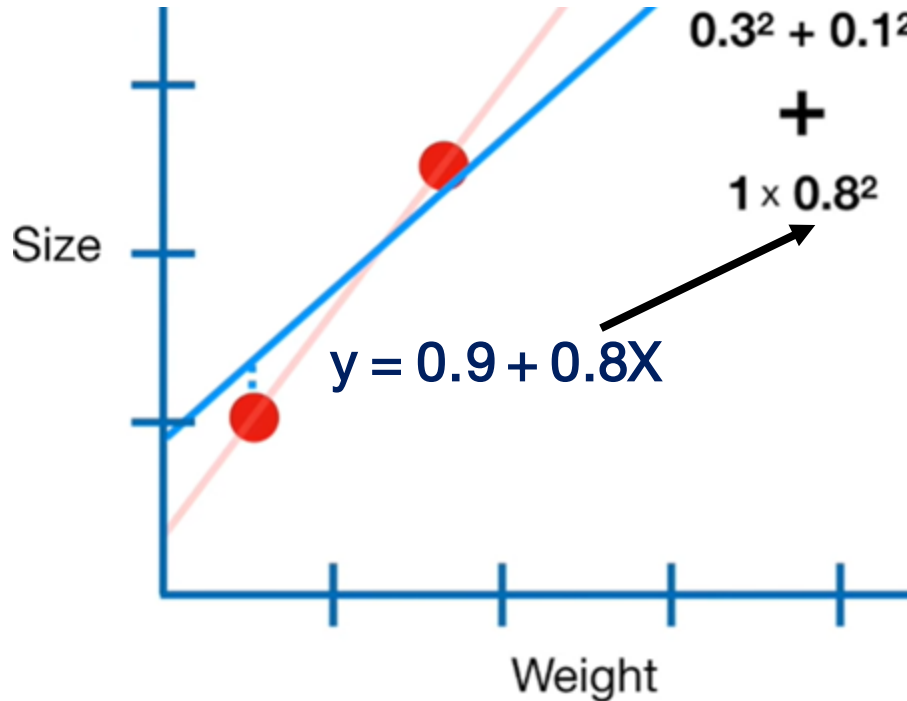


# REGRESIÓN RIDGE

Calculamos **nuevos errores**  
para la Regresión Ridge



# REGRESIÓN RIDGE



$$0.3^2 + 0.1^2$$

+

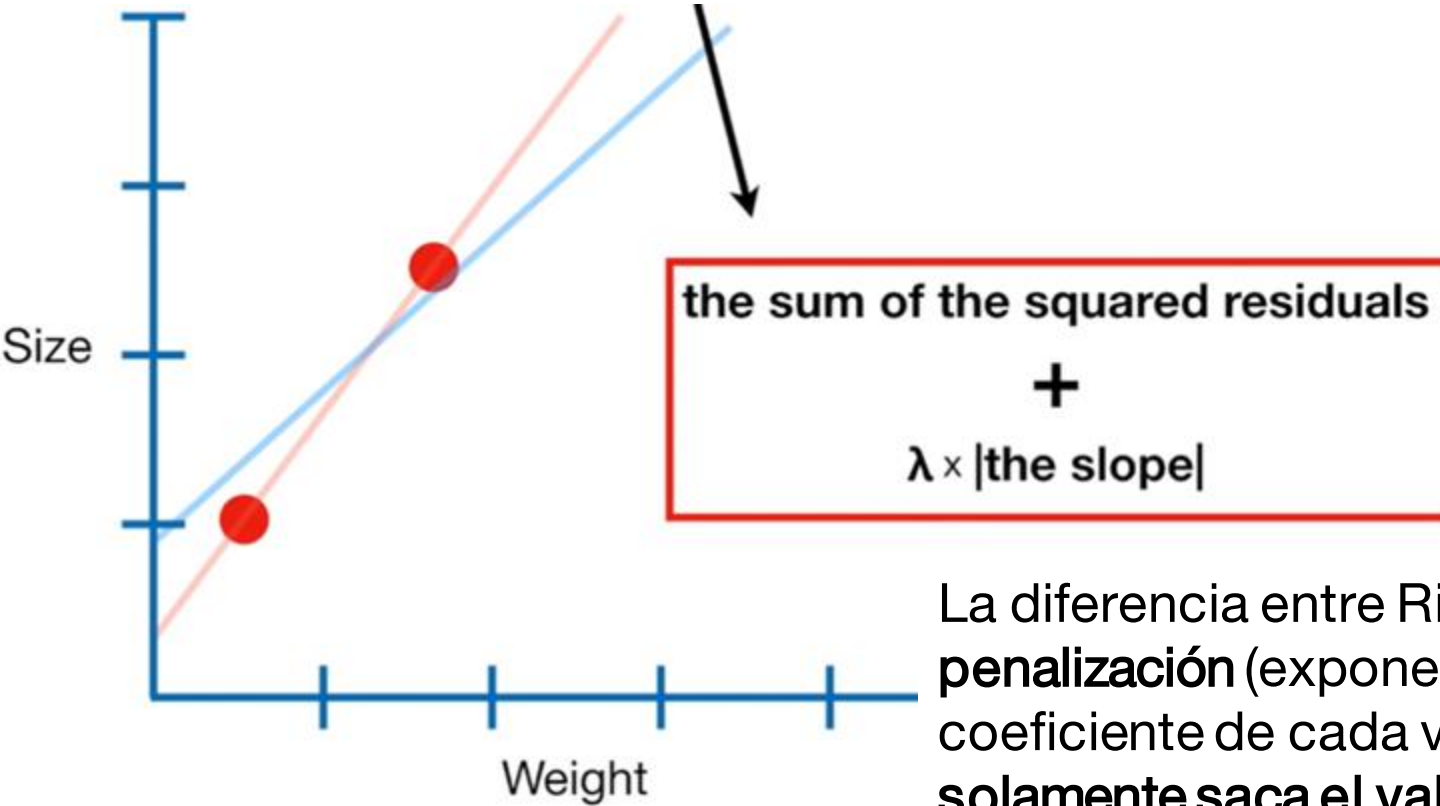
$$1 \times 0.8^2$$

$$= 0.09 + 0.01 + 0.64$$

$$= 0.74$$

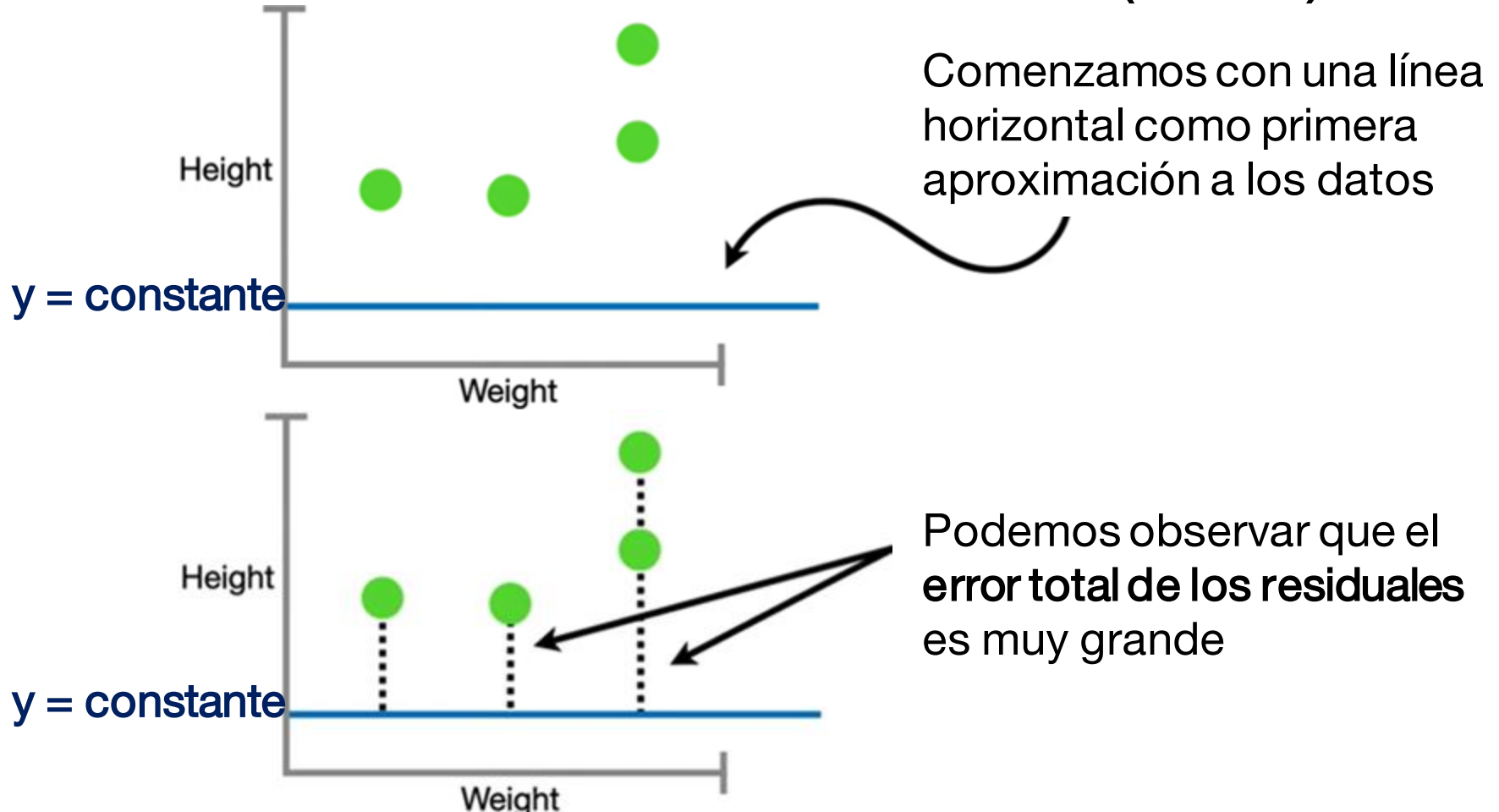
Esta nueva recta obtuvo un **menor error total**. De esta manera podemos seguir creando nuevas rectas hasta encontrar el **mínimo error total**.

# REGRESIÓN LASSO

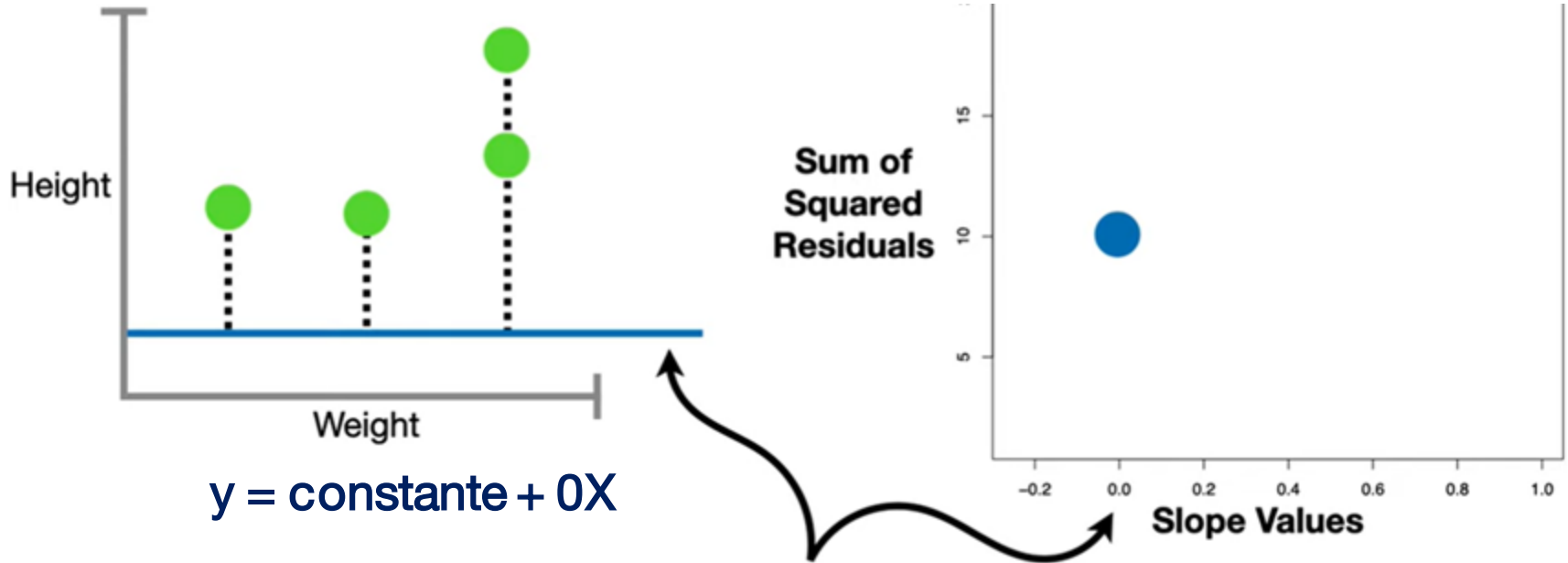


La diferencia entre Ridge y Lasso es la **penalización** (exponente de **slope** o coeficiente de cada variable). Lasso **solamente saca el valor absoluto**, es decir **penaliza menos**

# OPTIMIZACION DEL COEFICIENTE (SLOPE)

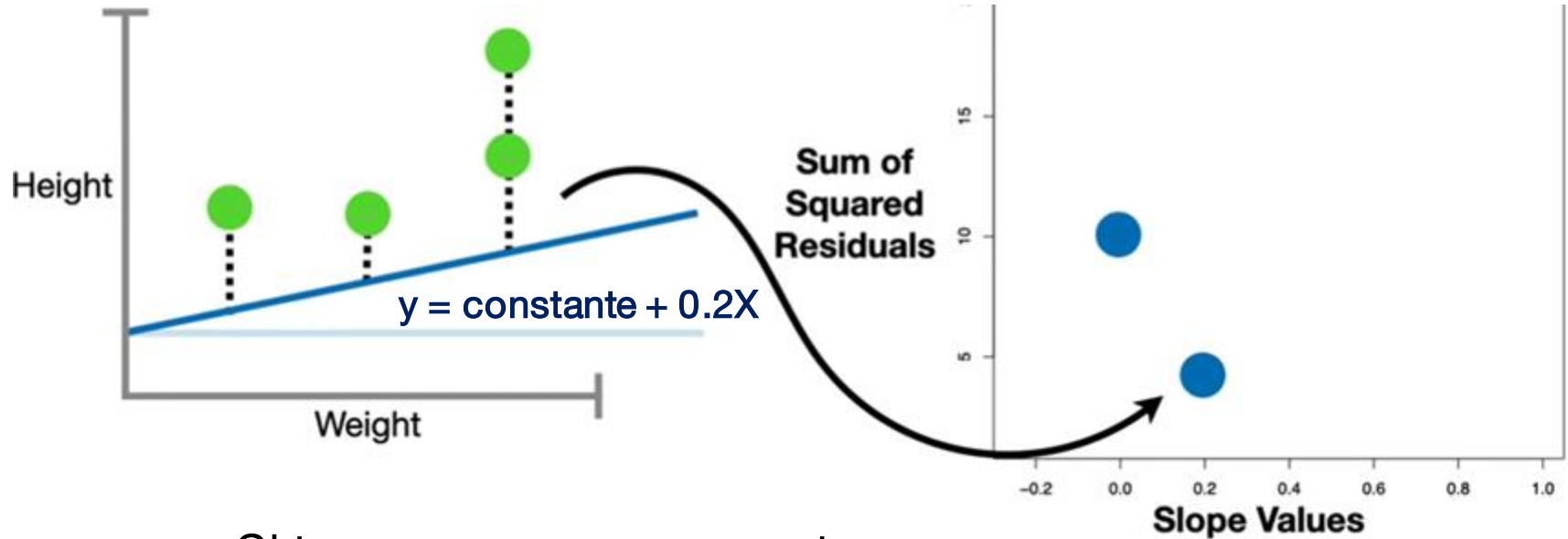


# OPTIMIZACION DEL COEFICIENTE (SLOPE)



Cuando se traza una línea recta el **coeficiente** que acompaña a **Weight = 0**

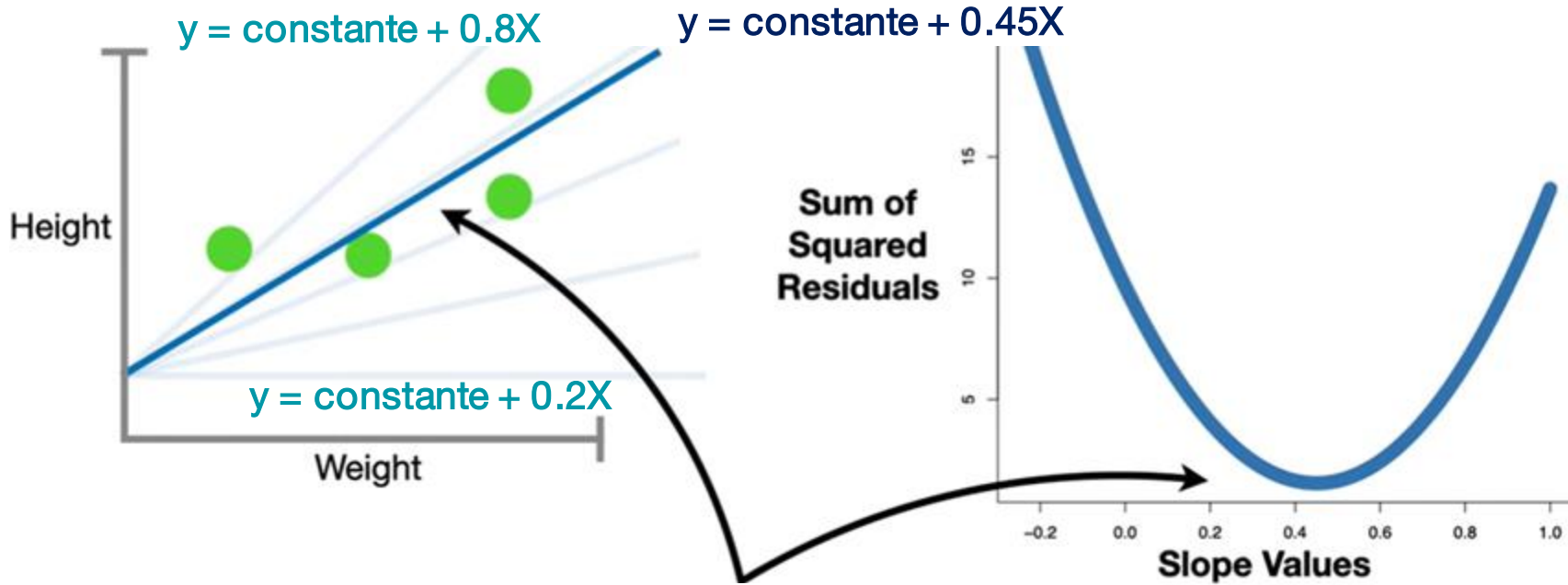
# OPTIMIZACION DEL COEFICIENTE (SLOPE)



Si trazamos una nueva recta con un **coeficiente diferente de 0 (0.2)** generamos un **menor error total de residuales**



# OPTIMIZACION DEL COEFICIENTE (SLOPE)



El punto óptimo con el **mínimo error total** de los **residuales** ocurre cuando el **coeficiente es 0.45**

$$\text{Size} = \text{y-intercept} + \text{slope}_1 \times \text{Weight} + \text{diet difference} \times \text{High Fat Diet} \\ + \text{slope}_2 \times \text{Age} + \text{slope}_3 \times \text{Size of Father}$$

Cuando **todas** las variables predictoras son **significativas** es correcto **penalizar los coeficientes de forma significativa también**.  
Para esta opción se recomienda la **Regresión Ridge**.

Cuando **no estamos seguros** si todas las variables predictoras son significativas es correcto **penalizar menos los coeficientes**.  
Para esta opción se recomienda la **Regresión Lasso**.



$$\text{Size} = \text{y-intercept} + \text{slope} \times \text{Weight} + \text{diet difference} \times \text{High Fat Diet}$$

~~$$+ \text{astrological offset} \times \text{Sign} + \text{airspeed scalar} \times \text{Airspeed of Swallow}$$~~

# MÉTRICAS DE REGRESIÓN

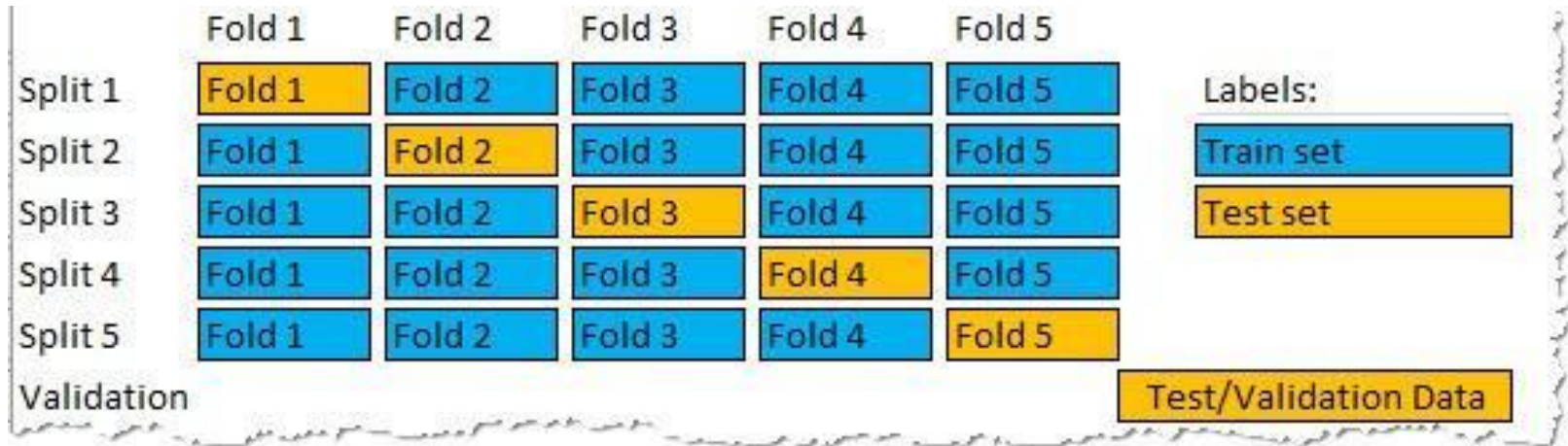
$$MSE = \frac{1}{m} \sum_{i=1}^m (y_i - t_i)^2$$

$$RMSE = \sqrt{\frac{1}{m} \sum_{i=1}^m (y_i - t_i)^2}$$

$$MAE = \frac{1}{m} \sum_{i=1}^m |y_i - t_i|$$

$$MSLE = \frac{1}{m} \sum_{i=1}^m (\log(y_i + 1) - \log(t_i + 1))^2$$

# CROSS-VALIDATION: K-FOLDS



- Mayor información en la siguiente ruta URL:

<https://towardsdatascience.com/complete-guide-to-pythons-cross-validation-with-examples-a9676b5cac12>