ALAN BADILLO SALAS

# MÓDULO II

OCTUBRE 2023

# CLASIFICACIÓN

# REGRESIÓN

# REGRESIÓN

Formula

$$R^2 = 1 - \frac{RSS}{TSS}$$

$R^2$ = coefficient of determination
$RSS$ = sum of squares of residuals
$TSS$ = total sum of squares

$$\text{RSS} = \Sigma\left(y_i - \widehat{y_i}\right)^2$$

Where: $y_i$ is the actual value and, $\widehat{y_i}$ is the predicted value.

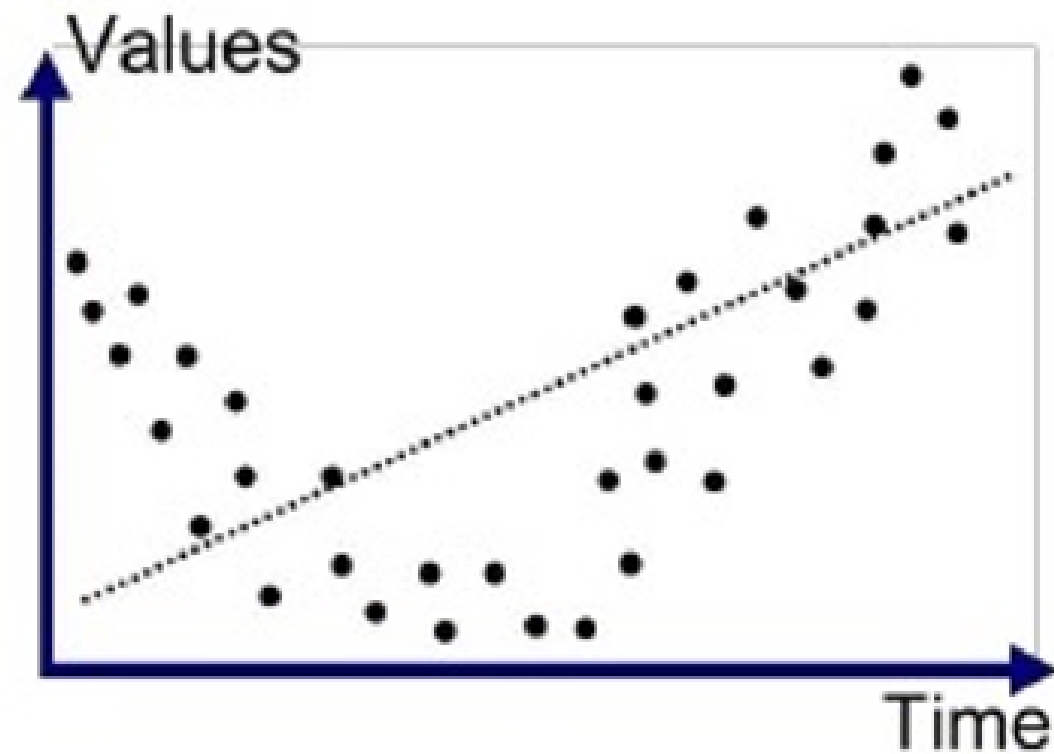$$\text{TSS} = \Sigma\left(y_i - \overline{y}\right)^2$$

Where: $y_i$ is the actual value and $\overline{y}$ is the mean value of the variable/feature
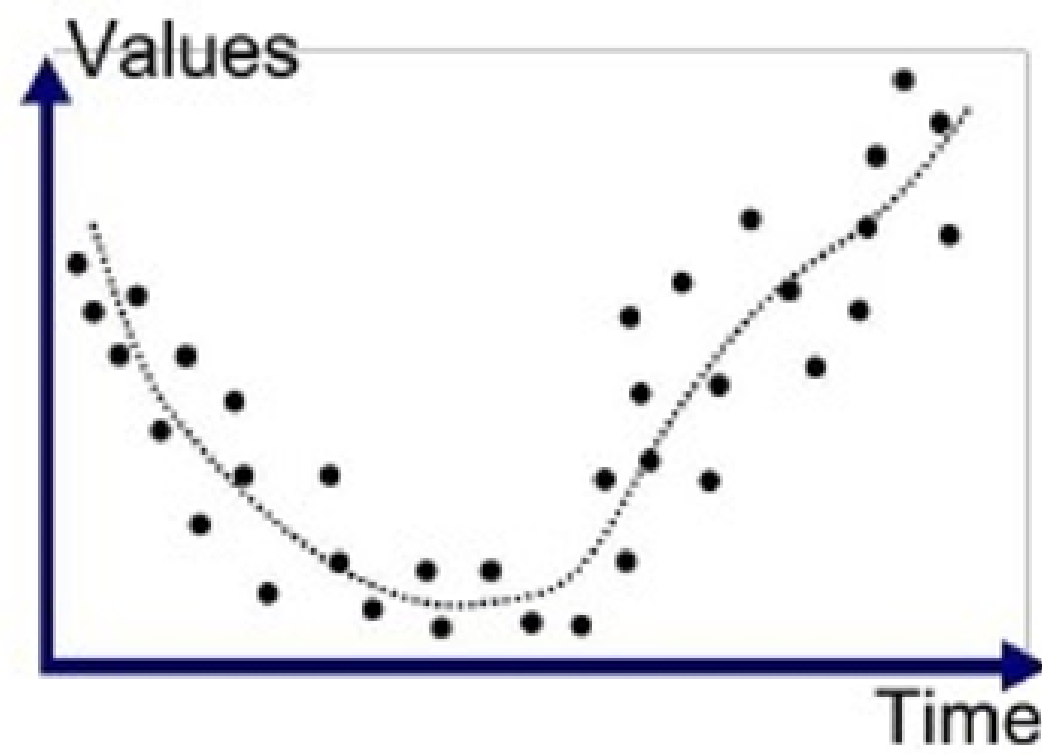
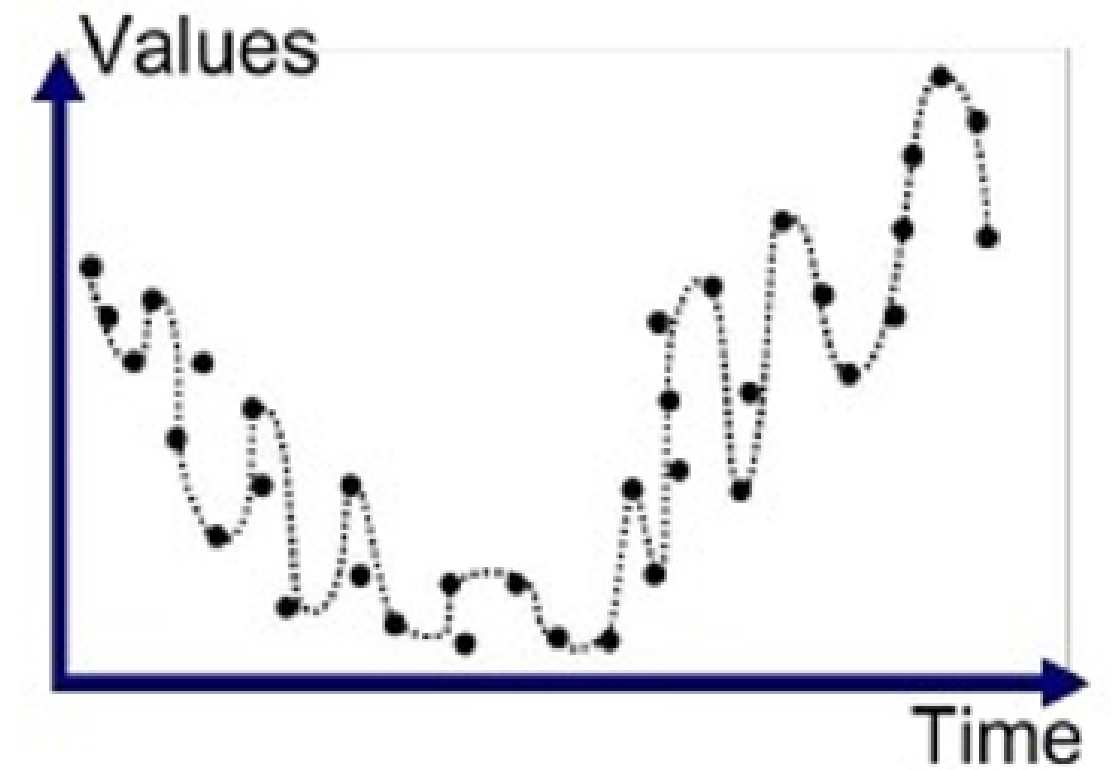# GENERALIZACIÓN, SOBREAJUSTE Y SUBAJUSTE

# GENERALIZACIÓN, SOBREAJUSTE Y SUBAJUSTE



Underfitted

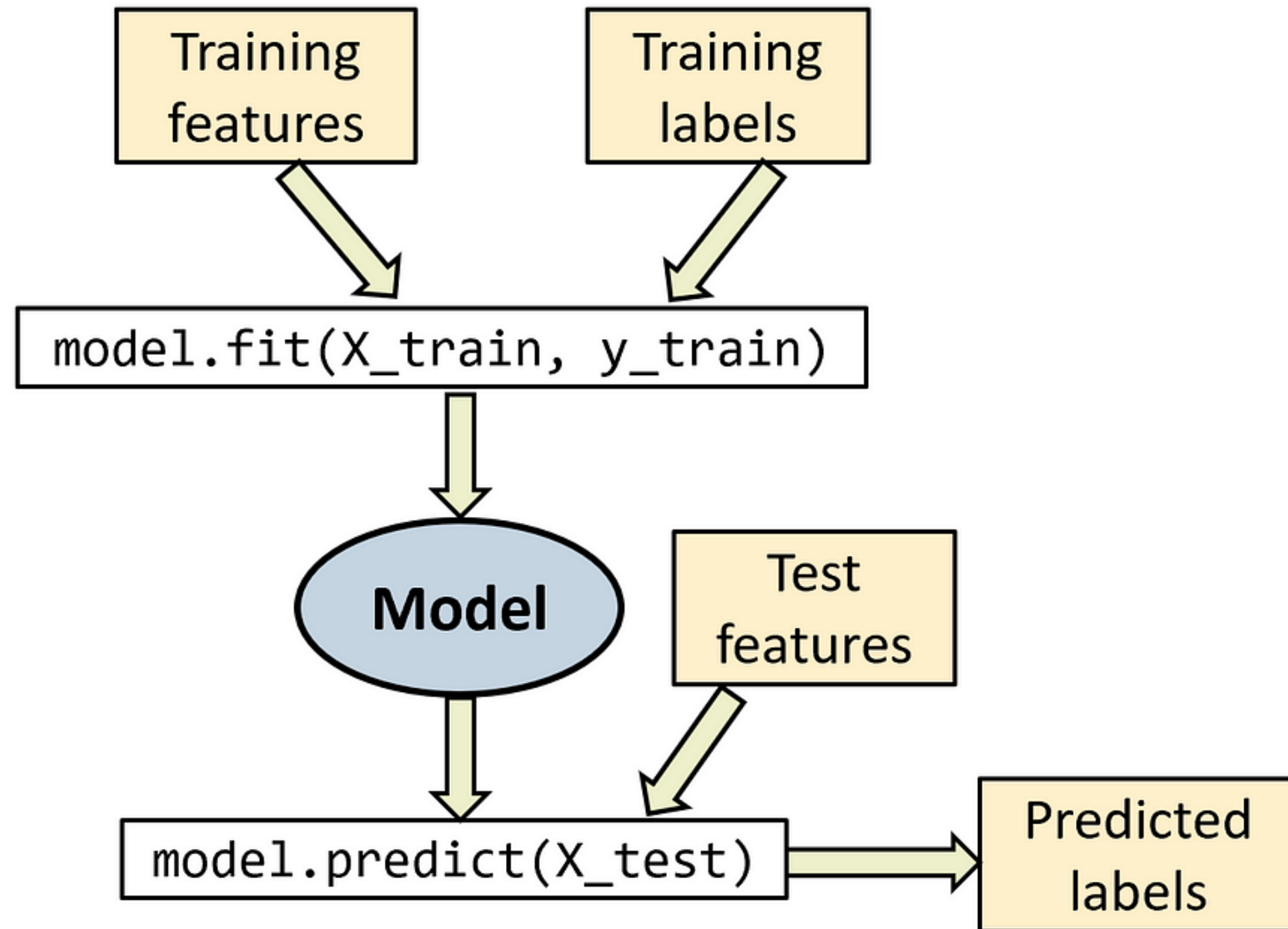Good Fit/Robust

Overfitted

# Algoritmo de Machine Learning Supervisado

# Algoritmo de Machine Learning Supervisado

```python
from sklearn import linear_model as lm

X = iris[["petal_length"]]
y = iris["petal_width"]

# Fit the linear model
model = lm.LinearRegression()
results = model.fit(X, y)

# Print the coefficients
print model.intercept_, model.coef_
```

```
-0.363075521319 [ 0.41575542]
```

# Estimaciones de Incertidumbre
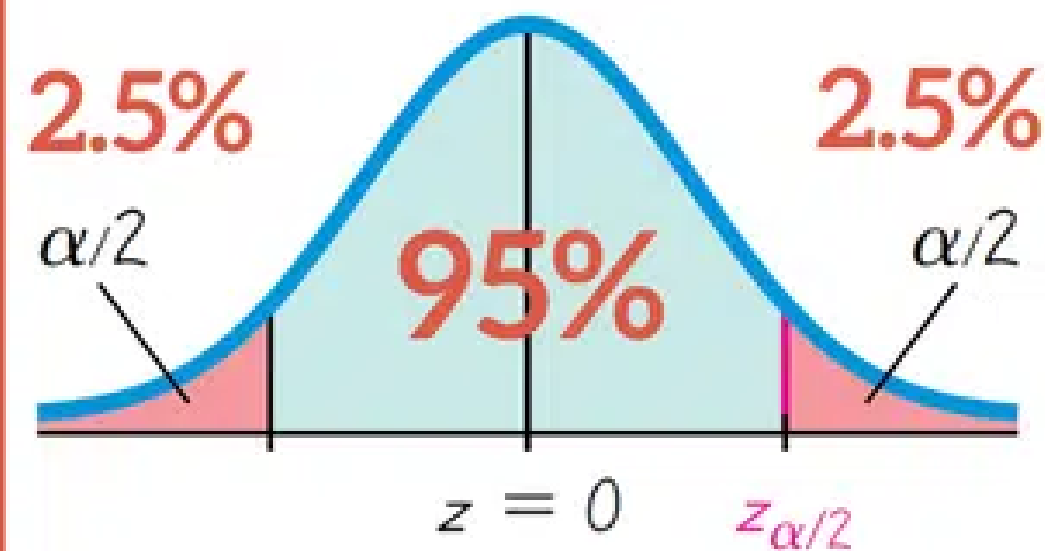


INTERVALO DE CONFIANZA DE LA MEDIA
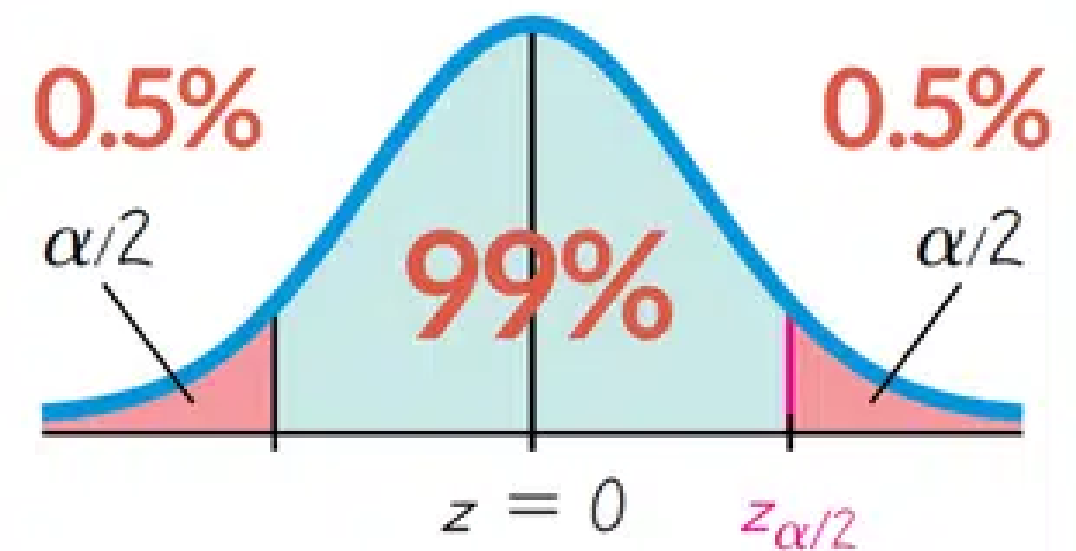95%
Nidel de significación alpha
$\alpha=100\%-95\% = 5\%$
2.5%    2.5%
$\alpha/2$    95%    $\alpha/2$
$z = 0$    $z_{\alpha/2}$
1.96

INTERVALO DE CONFIANZA DE LA MEDIA
99%
Nidel de significación alpha
$\alpha=100\%-99\% = 1\%$
0.5%    0.5%
$\alpha/2$    99%    $\alpha/2$
$z = 0$    $z_{\alpha/2}$
2.57