

Introducción a la Inteligencia Artificial

Clase ~~4~~ 3



Índice

1. Notas clase anterior
2. Análisis de la regresión lineal (R^2)
3. Descomposición Bias-Variance
4. Práctica

Analizamos el caso simple: (1 regresor, 1 var. dep)

partimos de $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

con esto podemos definir los **residuos**

$$r_i = y_i - \hat{y}_i \quad \forall i \in [1, \dots, N]$$

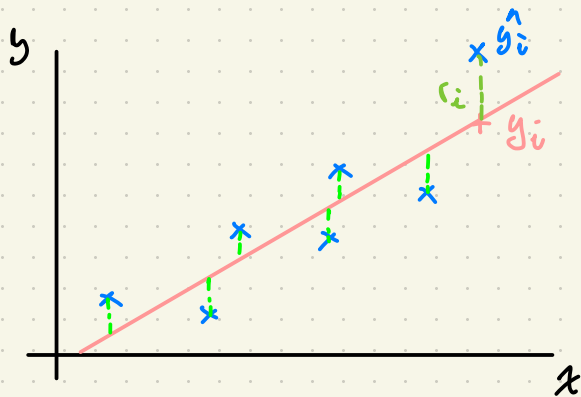
bondad del ajuste $\propto \sum_i r_i^2$

$$r_i \sim N(0, \sigma_r)$$

$$(\hat{\beta}_0, \hat{\beta}_1) = \arg \min_{\beta} \sum_i \underbrace{(y_i - \beta_0 - \beta_1 x_i)^2}_{r_i^2}$$

$$\begin{cases} \partial_{\beta_0} f = -2 \sum (y_i - \beta_0 - \beta_1 x_i) = 0 \rightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \\ \partial_{\beta_1} f = -2 \sum (y_i - \beta_0 - \beta_1 x_i) \cdot x_i = 0 \rightarrow \textcircled{1} \end{cases}$$

⊕ Aquí \bar{y} y \bar{x} representan los promedios.



$$\bar{r} = \bar{y} - \hat{\bar{y}} = \bar{y} - \bar{x} \hat{\beta} = \bar{y} - \bar{H} \bar{y}$$

$$\bar{r} = \underbrace{(\mathbb{I} - H)}_{\substack{\text{simétrica,} \\ \text{idempotente,} \\ \text{tr} = n-p}} \bar{y}$$

$$(1) \rightarrow \sum (y_i - \beta_0 - \beta_1 x_i) x_i = 0$$

$$\sum (y_i x_i - \beta_0 x_i - \beta_1 x_i^2) = 0$$

$$\sum (y_i x_i - \bar{y} x_i - \beta_1 x_i - \beta_1 x_i^2) = 0$$

Desarrollando (Ver apunte) llegamos a:

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}}{\frac{\sum (x_i - \bar{x})^2}{n}}$$

\hat{C}_{xy}
 $\hat{\sigma}_x^2$

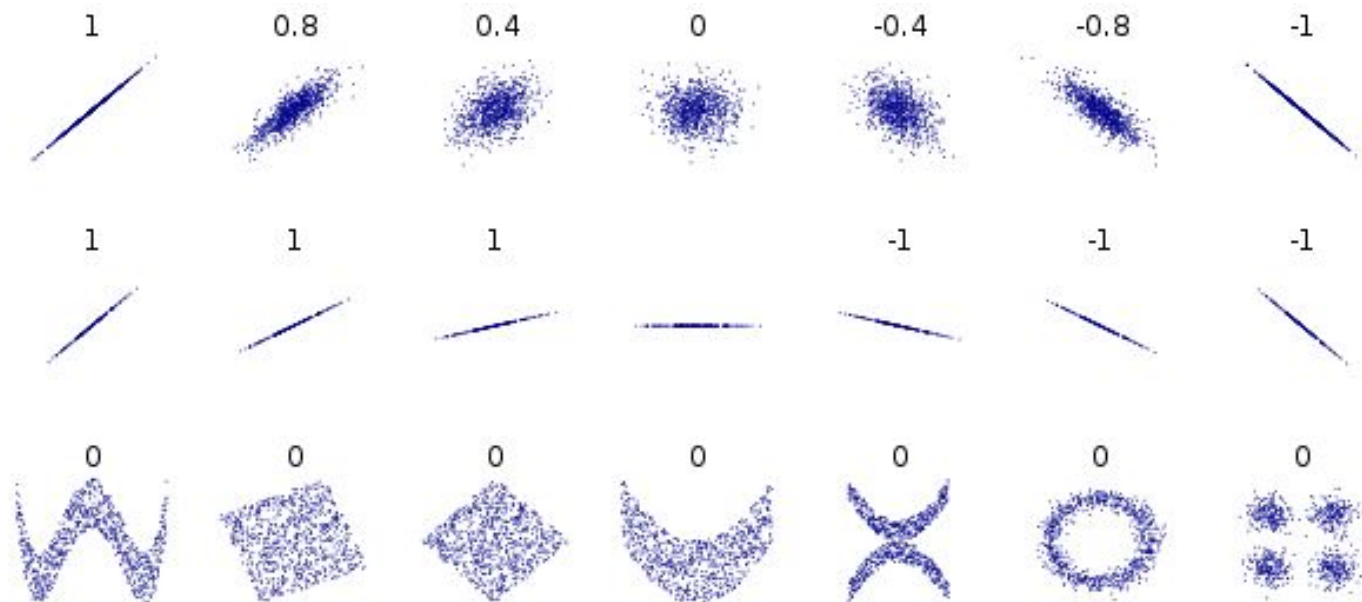
$$\begin{cases} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \end{cases}$$

$$\begin{cases} \hat{\beta}_1 = \frac{\text{covar}(x, y)}{\text{var}(x)} = \hat{\rho}_{xy} \end{cases}$$

coef. de correlación
lineal de Pearson

$$\rho_{xy} \in [-1, 1]$$

Coeficiente de correlación de Pearson (Lineal)



Analizamos los errores de la regresión:

$$y_i = \hat{y}_i + r_i$$

\bar{y} : media muestral (promedio)

el término
de los
residuos
es por
indep.

$$y_i - \bar{y} = \hat{y}_i - \bar{y} + r_i \Rightarrow \sum_{i=0}^N (y_i - \bar{y})^2 = \sum_{i=0}^N (\hat{y}_i - \bar{y})^2 + \sum_{i=0}^N r_i^2$$

(1) (2) (3)

① TSS: tasa de variabilidad total

② ESS: tasa de variabilidad explicada

③ RSS: suma de los residuos al cuadrado

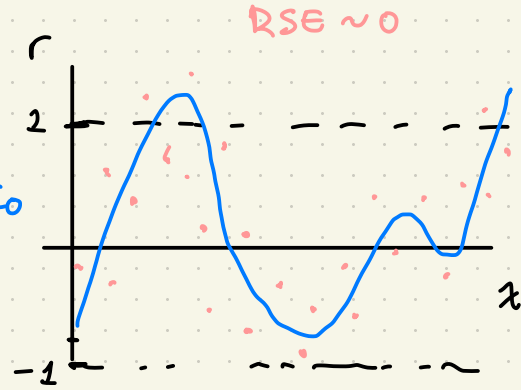
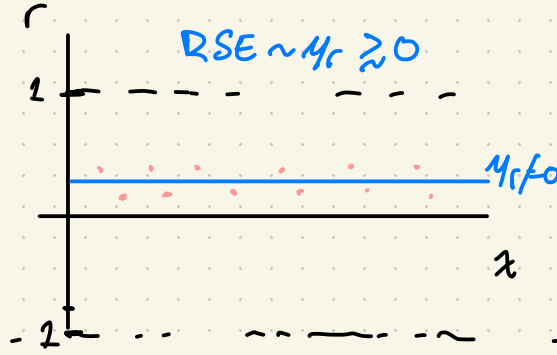
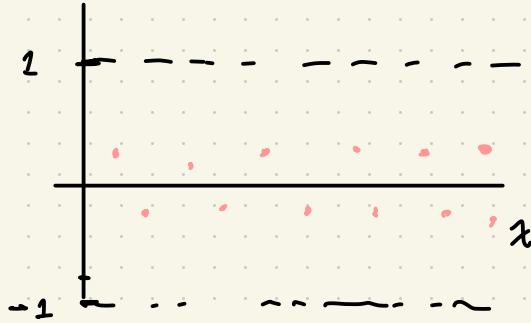
Estos valores me permiten construir métricas de bondad de ajuste para diagnosticar mi modelo.

métricas posibles

- error residual (RSE)
- R^2 (coef. de pearson)
- Est. F (Análisis "Avanzado")

• Error residual:
$$RSE = \sqrt{\frac{RSS}{N-1}} = \frac{1}{N-2} \sqrt{\sum_i (y_i - \hat{y}_i)^2}$$

r_n ← residuo
normalizado
 $RSE \sim 0$



• Si RSE es bajo ~~entonces~~ nuestro ajuste ~~es~~ bueno
puede que sea

• RSE no es sensible a la distribución y/o tendencia funcional de los residuos.

Coeficiente de determinación - “R cuadrado”

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

RSS

TSS

SS = Sum of squares

res = residuos

tot = total

+ coef de pearson R^2 :

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS} \stackrel{(1)}{=} (\rho_{xy})^2$$

① es válido únicamente bajo el régimen de reg. simple

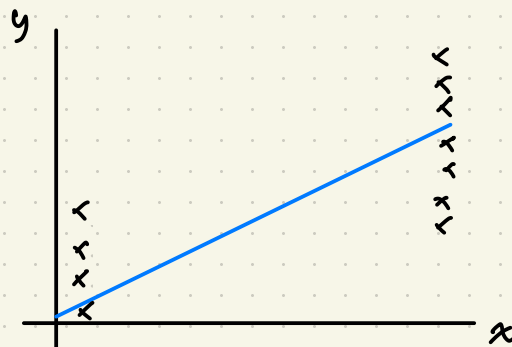
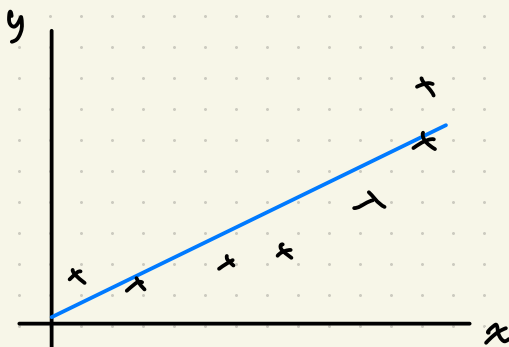
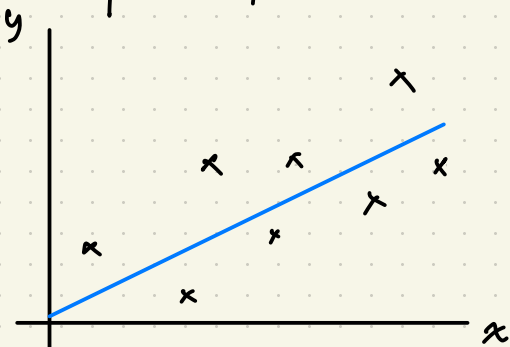
• R^2 no depende de las escalas, solo depende de las proporciones.

• $R^2 \in [0, 1] \Rightarrow R^2 \sim 1 \Rightarrow$ el modelo es "bueno"

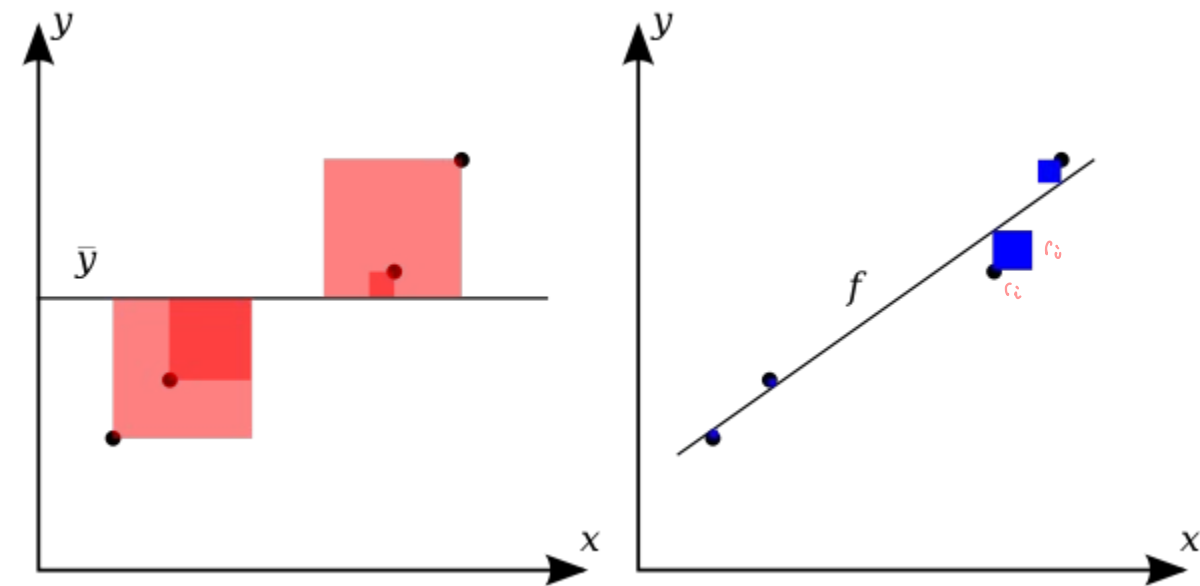
↑
si hay una
reg. lineal

$R^2 \lesssim 0 \Rightarrow$ el modelo es peor que haber aproximado con la media.

$\beta_0 \sim 3, \beta_1 \sim 0,5, R^2 \sim 0,7$



Regresión Lineal - R2



$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

ESS

$$SS_{reg} = \sum_i (f_i - \bar{y})^2$$

RSS

$$SS_{res} = \sum_i (y_i - f_i)^2 = \sum_i e_i^2$$

TSS

$$\begin{aligned} SS_{tot} &= \sum_i (y_i - \bar{y})^2 \\ &= SS_{res} + SS_{reg} \end{aligned}$$

¿Similar a σ^2 ?

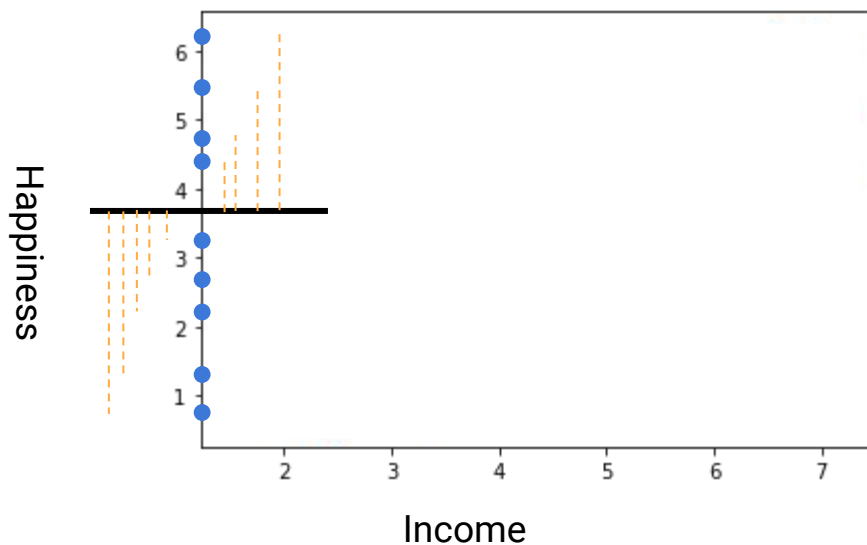
Regresión Lineal - R2

$$\begin{aligned} R^2 &= 1 - \frac{SS_{res}}{SS_{tot}} \\ &= 1 - \left(\frac{SS_{res}}{SS_{tot}} * \frac{n}{n} \right) \\ &= 1 - \frac{\sigma_{res}}{\sigma_{tot}} \end{aligned}$$

Proporción de varianza no explicada

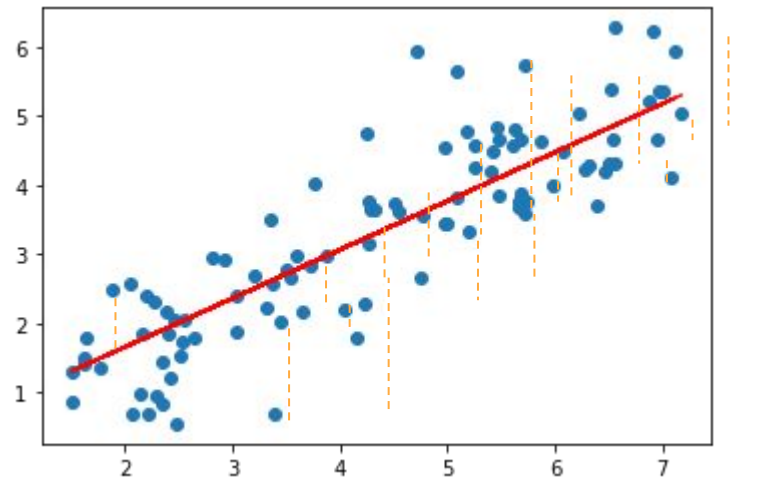
Proporción de varianza explicada

Regresión Lineal - R2



$$SS(media) = (happiness - media)^2$$

$$Variación(media) = \frac{(happiness - media)^2}{n}$$

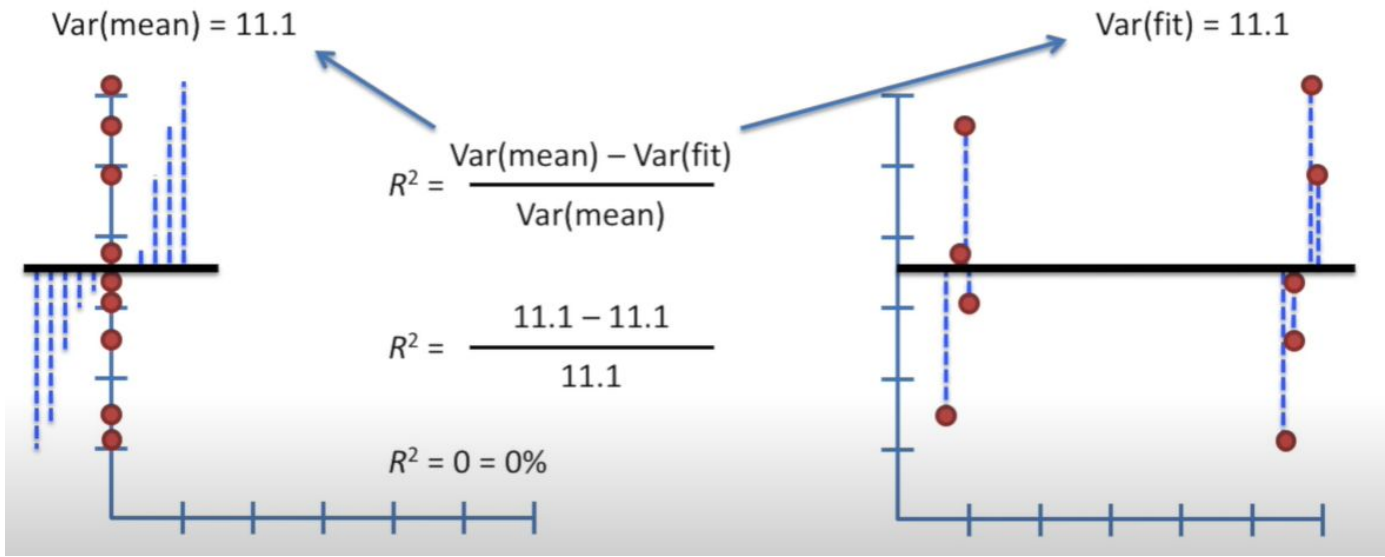


$$SS(fit) = (happiness - lr_fit)^2$$

$$Variación(fit) = \frac{(happiness - lr_fit)^2}{n}$$

Regresión Lineal - R2

$$R^2 = \frac{\text{Variación(media)} - \text{Variación(fit)}}{\text{Variación(media)}}$$



Fuente: StatQuest with Josh Starmer

Regresión Lineal - R2

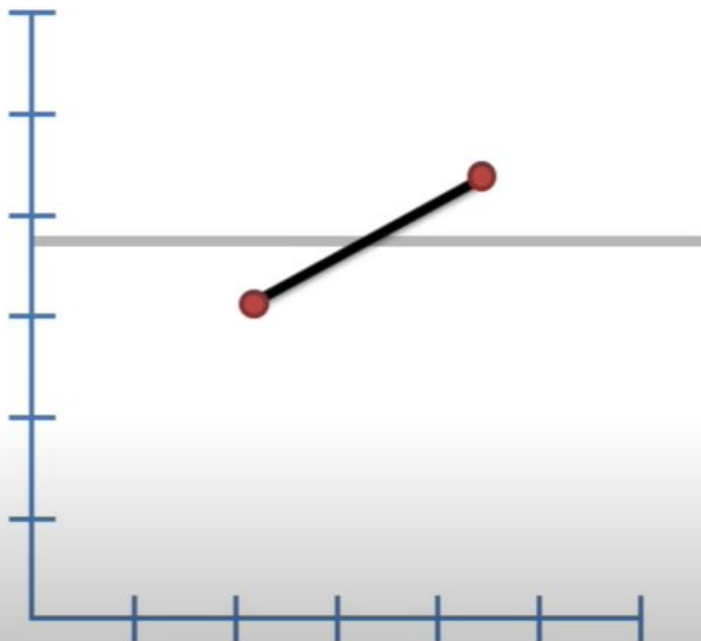
$$F = \frac{\text{Variación en happiness explicada por income}}{\text{Variación en happiness no explicada por income}}$$

$$SS(\text{mean}) = 10$$

$$SS(\text{fit}) = 0$$

$$R^2 = \frac{SS(\text{mean}) - SS(\text{fit})}{SS(\text{mean})}$$

$$= \frac{100 - 0}{100} = 100\%$$



Fuente: StatQuest with Josh Starmer

R² y el coeficiente de correlación de Pearson

$$\rho = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} \in [-1, 1]$$

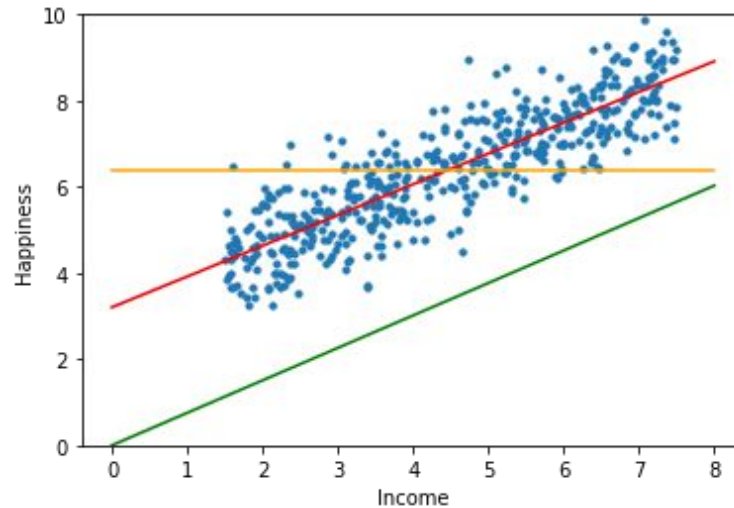
Regresión múltiple **con ordenada** → Correlación entre observación y predicción

Regresión **con ordenada** → Correlación entre variable dependiente e independiente

$$\rho^2 = R^2 \in [0, 1]$$

¿R² negativo?

$$R^2 = 1 - \frac{\sigma_{res}}{\sigma_{tot}} < 0 \Leftrightarrow \sigma_{res} > \sigma_{tot}$$



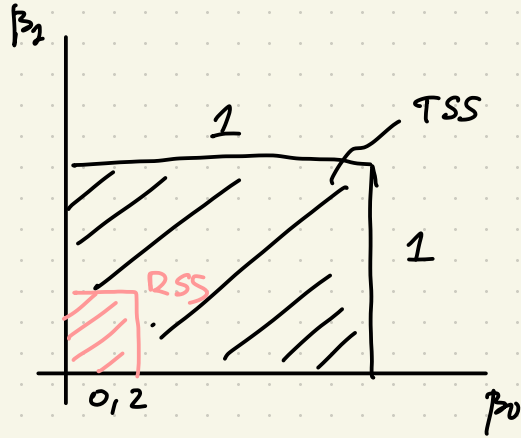
Predecir con el promedio es mejor que el modelo

R² inflation

(caso regresión multivariada)

↑ cantidad de predictores → ↑ R² → F-test para comparación válida entre modelos.

curso regresión lineal simple

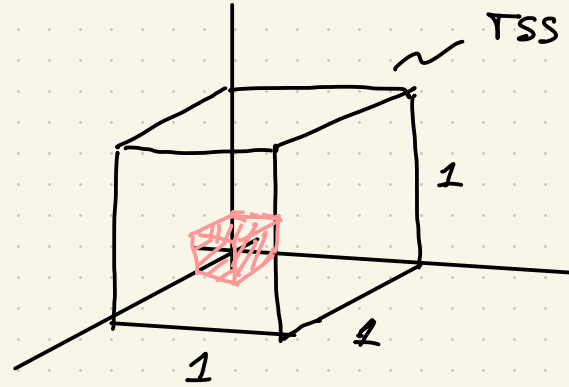


$$R^2 = 1 - \frac{RSS_2}{TSS_2}$$

$$y = \beta_0 + \beta_1 x$$

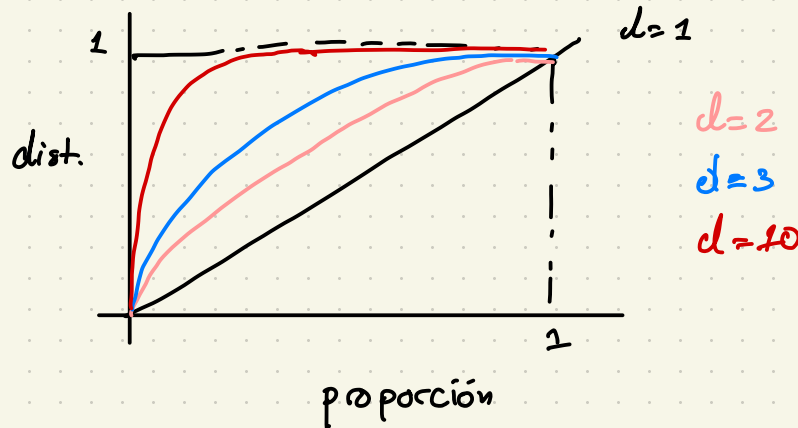
R² Inflation

regresión lineal multivariable



$$1 - \frac{RSS_3}{TSS_3}$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$



Estadístico F:

Vamos a suponer que $X \in \mathbb{R}^{n \times p}$ donde n es la cantidad de datos y p la cantidad de regresores (columnas).

Armando la tabla ANOVA de la siguiente forma.

Fuente de variación	suma de los cuadrados	grados de libertad	cuadrados medios
Explicada	$\sum (\hat{y}_i - \bar{y})^2$	$p - 1$	$S_E = \frac{1}{p - 1} \sum_i (\hat{y}_i - \bar{y})^2$
residual	$\sum r_i^2$	$n - p$	$S_R = \frac{1}{n - p} \sum_i r_i^2$
total	$\sum (y_i - \bar{y})^2$	$n - 1$	

bajo estas condiciones podemos analizar múltiples test de hipótesis como:

- $H_0: \exists j / \beta_j = 0$; $H_A: \nexists j / \beta_j = 0$
- $H_0: \beta_i = \beta_j$; $H_A: \beta_i \neq \beta_j$ - etc.

con esto planteamos nuestro test de hipótesis:

$$T_H \begin{cases} H_0 : \beta_j = 0 \quad \forall j=0 & \text{Aquí buscamos rechazar } H_0 \\ H_1 : \beta_j \neq 0 \quad j \in [1, p] & \text{con una significancia de } \alpha \end{cases}$$

$$T_H \sim \mathcal{F}_{\alpha, \beta} \text{ (distrib. F) } \rightarrow \text{buscamos } P(F \geq f_{crit})$$

Repaso

¿Qué es un test de hipótesis?

Un test de hipótesis es un procedimiento estadístico que permite tomar una decisión acerca de una afirmación o supuesto (hipótesis) con base en los datos recolectados. Existen dos tipos principales de hipótesis:

1. Hipótesis nula (H_0): Propone que no hay una relación significativa o un efecto notable entre dos fenómenos. Es la hipótesis que se asume verdadera hasta que se demuestre lo contrario.
2. Hipótesis alternativa (H_1): Sostiene que hay una relación significativa o un efecto notable entre dos fenómenos. Esta hipótesis se acepta si los datos recolectados proporcionan suficiente evidencia en contra de la hipótesis nula.

El proceso de prueba en sí consiste en recoger datos, analizarlos y luego decidir si la evidencia respalda la hipótesis nula o la alternativa. Si la evidencia va en contra de la hipótesis nula (basado en un nivel de significancia preestablecido, usualmente 0.05), "rechazamos" la hipótesis nula y aceptamos la alternativa. Por otro lado, si la evidencia no es suficientemente fuerte, "aceptamos" la hipótesis nula, que no significa necesariamente que la hipótesis nula sea verdadera, sino que no tenemos suficientes pruebas para rechazarla. Finalmente, es interesante notar que estas conclusiones están sujetas a errores. Por ejemplo, podríamos rechazar una hipótesis nula verdadera (Error tipo I) o aceptar una hipótesis nula falsa (Error tipo II).

¿Cómo calculamos el F test? → `Scipy.stats.f` or `scipy.stats.f_oneway`

1. armamos la tabla ANOVA

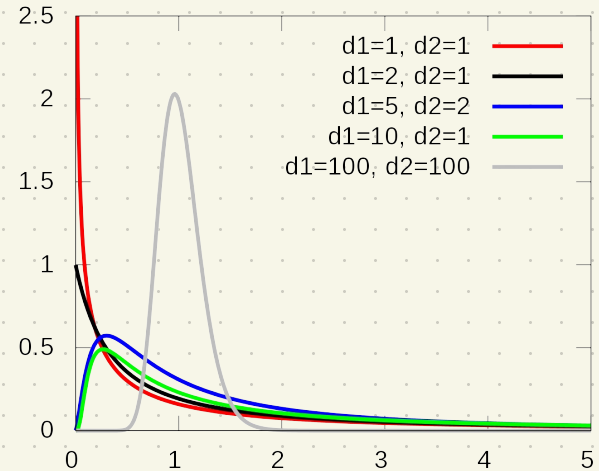
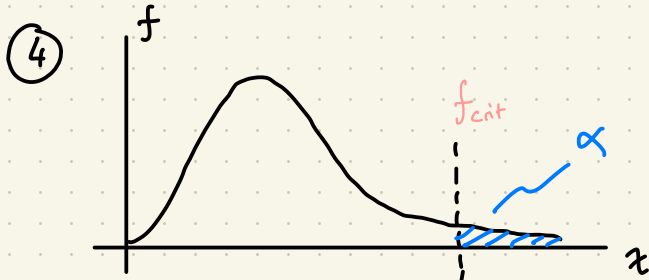
2. obtenemos el estadístico $F = \text{ESS} / S_r^2$

3. Encontramos f_{crit} (p-valor), definimos el nivel de significancia α

$$f \sim \tilde{F}_{gl_{\text{ESS}}, gl_{\text{RSS}}} \quad \text{gl: 'grados de libertad'}$$

4. buscamos $P(\tilde{F} \geq f_{\text{crit}}) = \alpha$

5. si $F \geq f_{\text{crit}}$ rechazamos H_0



Otras medidas a tener en cuenta

Mean Absolute Error

$$\text{MAE} = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n} = \frac{\sum_{i=1}^n |e_i|}{n}.$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2. \rightarrow \text{Mean Square Error}$$

$$\text{RMSD} = \sqrt{\frac{\sum_{t=1}^T (\hat{y}_t - y_t)^2}{T}}.$$

\hookrightarrow Root Mean Square Deviation

Bias-Variance Tradeoff

Cuando utilizamos el **error cuadrático medio** en un modelo de ML, podemos descomponer el mismo en términos de bias (sesgo) y variance (varianza).

1. error de estimación.

$$MSE = \underbrace{Bias(\hat{f})^2 + Var(\hat{f})}_{\text{error de estimación}} + \sigma_\epsilon^2$$

2. sesgo.

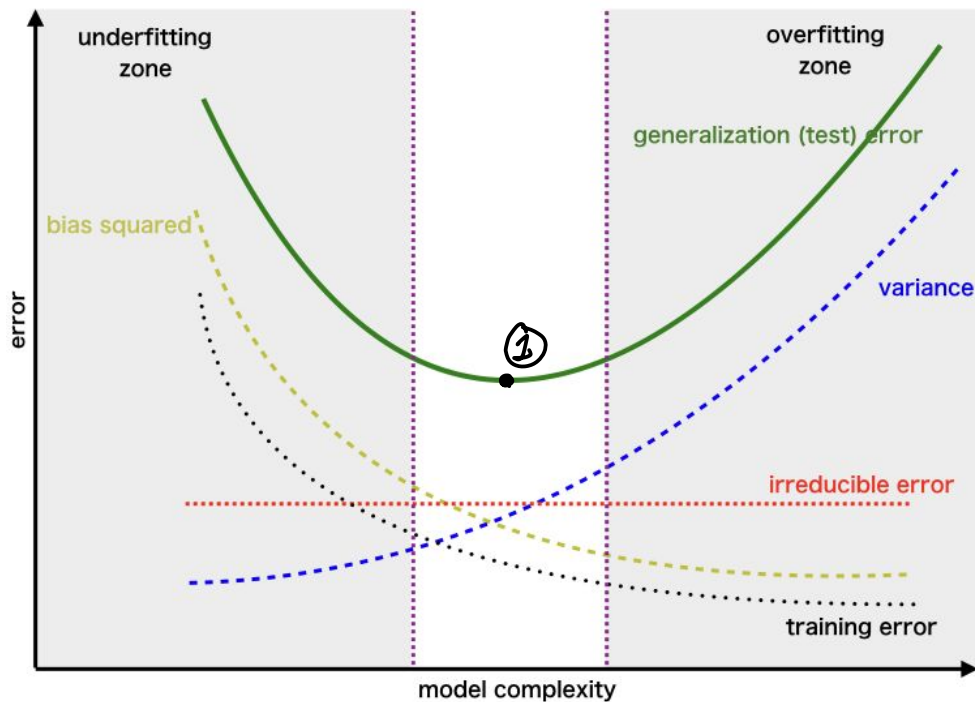
$$\begin{aligned}\hat{\beta} - \beta &= (X^t X)^{-1} X^t y - \beta \\ &= (X^t X)^{-1} X^t (X\beta + \epsilon) - \beta \\ &= (X^t X)^{-1} X^t \epsilon\end{aligned}$$

$$Bias = E[\hat{f} - f]$$

$$Var(\hat{f}) = E[(E[\hat{f}] - \hat{f})^2]$$

$$E(\hat{\beta} - \beta) = (X^t X)^{-1} X^t \underbrace{E(\epsilon)}_{=0} = 0$$

Bias-Variance Tradeoff



② este punto se obtiene por optimización (early-stop) y la implementación depende del problema y del modelo.

cota cramer-rao de Varianza de un estimador

↳ en gal # iteraciones
epochs

Nota adicional:

Existen múltiples test de interés en la regresión lineal, en los vamos a tener relacionados a los param. del modelo (los β 's), estos buscan responder:

+ ¿cuál es la validez del modelo? + ¿cuáles son variables importantes?

En gen. planteamos:

$$- H_0: \bar{R} \bar{\beta} = \bar{r} \text{ vs. } H_1: \bar{R} \bar{\beta} \neq \bar{r} \quad \text{donde } \bar{R} \in \mathbb{R}^{J \times p} \wedge \bar{r} \in \mathbb{R}^{J \times 1}$$

Algunos casos de interés:

$$H_0: \beta_1 = \beta_2 = \dots = 0 \quad (\beta_j = 0 \quad \forall j \in [1, p-1]) \quad \text{test de bondad de la regresión}$$

Aquí:

$$\bar{R} = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \in \mathbb{R}^{p-1, p} \quad \bar{r} = [0, \dots, 0]$$

Bibliografía

- The Elements of Statistical Learning | Trevor Hastie | Springer
- An Introduction to Statistical Learning | Gareth James | Springer
- Deep Learning | Ian Goodfellow | <https://www.deeplearningbook.org/>
- Stanford | CS229T/STATS231: Statistical Learning Theory | <http://web.stanford.edu/class/cs229t/>
- Mathematics for Machine Learning | Deisenroth, Faisal, Ong
- Artificial Intelligence, A Modern Approach | Stuart J. Russell, Peter Norvig