

# Aprendizaje Automático

## Segundo Cuatrimestre de 2016

# Regresión Lineal

Clase dada en el pizarrón, basada fuertemente en el Capítulo 3 del libro: James, Witten, Hastie & Tibshirani, "An Introduction to Statistical Learning with Applications in R", 6th ed, Springer, 2015. <http://www-bcf.usc.edu/~gareth/ISL/>  
Esta presentación se ofrece sólo a modo de referencia, y no está completa!

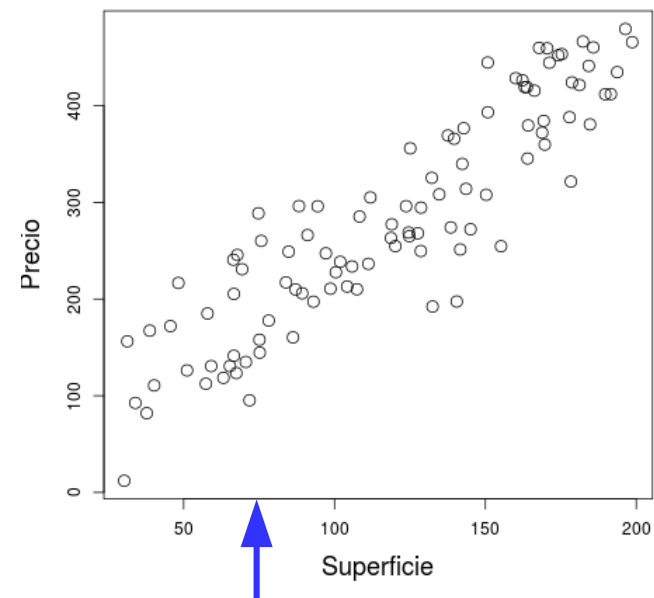


DEPARTAMENTO  
DE COMPUTACION

Facultad de Ciencias Exactas y Naturales - UBA

# Problema de Regresión

- Predecir una **cantidad**:
  - La probabilidad de que un mail sea spam.
  - El tiempo de demora de un vuelo.
  - El precio de una propiedad.



¿Precio de un departamento de 75m<sup>2</sup>?

# Regresión Lineal Simple

- Consiste en predecir una respuesta cuantitativa  $Y$  en base a una única variable predictora  $X$ .

$$Y \approx \beta_0 + \beta_1 \cdot X$$

- Ejemplo:

$$\text{Precio} \approx \beta_0 + \beta_1 \cdot \text{Superficie}$$

Ordenada al origen  
(*intercept*)

Pendiente  
(*slope*)

- $\beta_0$  y  $\beta_1$  son los coeficientes desconocidos que vamos a estimar, o ajustar en base a los datos de entrenamiento. Una vez estimados, los podemos usar para predecir:

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

Valor predicho para  $Y$   
cuando  $X=x$

Estimación de  $\beta_0$

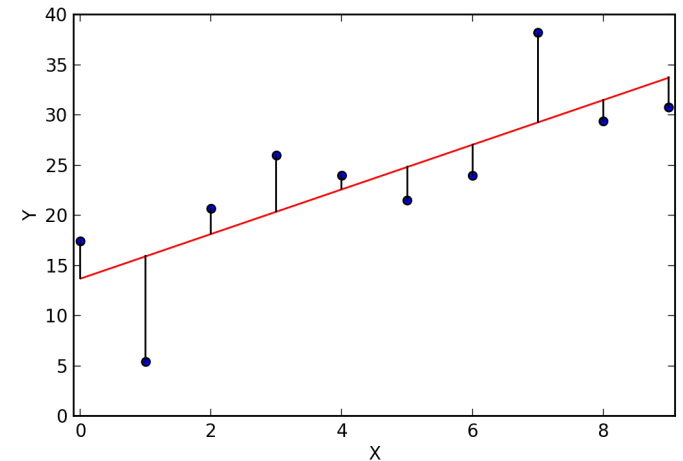
Estimación de  $\beta_1$

Nueva instancia

# Estimación de coeficientes

- Def: Residuo o error de predicción

$$e_i = y_i - \hat{y}_i$$



- Residual sum of squares:

$$\text{RSS} = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 \cdot x_i)^2$$

- Los residuos se elevan al cuadrado para sacar el signo y para que RSS sea diferenciable.
- Hay que tener cuidado con los outliers en los datos, porque RSS penaliza los residuos grandes.

# Estimación de coeficientes

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 \cdot x_i)^2$$

- Para estimar los coeficientes, buscamos minimizar RSS:

$$\frac{\delta \text{RSS}}{\delta \hat{\beta}_0} = 0$$

$$\frac{\delta \text{RSS}}{\delta \hat{\beta}_1} = 0$$

- Con un poco de análisis matemático, llegamos a estas expresiones (ejercicio 1 de la guía):

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \cdot \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\text{Cov}[X, Y]}{\text{Var}[X]}$$

- Otra opción es minimizar RSS usando la técnica iterativa de Descenso por el Gradiente. (Ejercicio de la guía.)

# Predicción de nuevos valores

- Estimados los coeficientes, los podemos usar para predecir:

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

- En nuestro ejemplo, podríamos predecir el precio de un departamento de 75m<sup>2</sup>:

$$\widehat{Precio} = \hat{\beta}_0 + \hat{\beta}_1 \cdot 75$$

# Exactitud de los coeficientes estimados

- Al correr una regresión lineal, es común reportar el error estándar de cada estimador:  $SE(\hat{\beta}_0)$  y  $SE(\hat{\beta}_1)$
- Esto es útil para estimar intervalos de confianza de las predicciones.
- Evaluar la significancia de cada estimador, mediante un test estadístico.
  - $p$ -valor bajo (típicamente,  $p < 0.05$  o  $p < 0.01$ )  $\rightarrow$  es improbable observar al azar una asociación semejante entre  $X$  e  $Y$ .
  - $p$ -valor alto  $\rightarrow$  es probable que la asociación observada sea sólo consecuencia del azar.

# Evaluación del modelo

- RSS: Variabilidad no explicada por el modelo

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- TSS (Total Sum of Squares): Variabilidad total de los datos

$$\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2$$

- $R^2$ : Proporción de la variabilidad explicada por el modelo.

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

$R^2 \rightarrow 0$  cuando el modelo explica poco de la variabilidad de los datos.

$R^2 \rightarrow 1$  cuando el modelo explica mucho de la variabilidad de los datos.



# Regresión Lineal Múltiple

- Consiste en predecir una respuesta cuantitativa  $Y$  en base a una múltiples variables predictoras  $X_1, X_2, \dots, X_p$ .

$$Y \approx \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

- RSS se define igual que para la regresión lineal simple:

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- Los coeficientes se estiman en forma análoga.
- TSS,  $R^2$  también se definen en forma similar.
- ¿Qué pasa si hay variables categóricas? (Ejercicio de la guía.)