

# ANÁLISIS INTELIGENTE DE DATOS:

## TAREA 2

Felipe Valdivia  
Anibal Pérez

# INTRODUCCIÓN

Se utilizarán los métodos de regresión lineal estudiados en clases para analizar:

- prostate-cancer, datos de pacientes luego de aplicada prostatectomía.
- Datos sobre Utilidades en estreno de Películas.

# DESCRIPCIÓN DATASET PROSTATE-CANCER

- 97 registros.
- 9 Columnas.
- 8 serán predictores.
- 1 de respuesta (lpsa).
- Datos de Train y Test.

# REGRESION LINEAL ORDINARIA (LSS)

- MSE en datos de entrenamiento (67)

$$\text{MSE} = 0.439$$

- MSE en datos de prueba (30)

$$\text{MSE} = 0.521$$

- MSE en Cross-Validation K = 5

$$\text{MSE} = 0.957$$

- MSE en Cross-Validation K = 10

$$\text{MSE} = 0.757$$

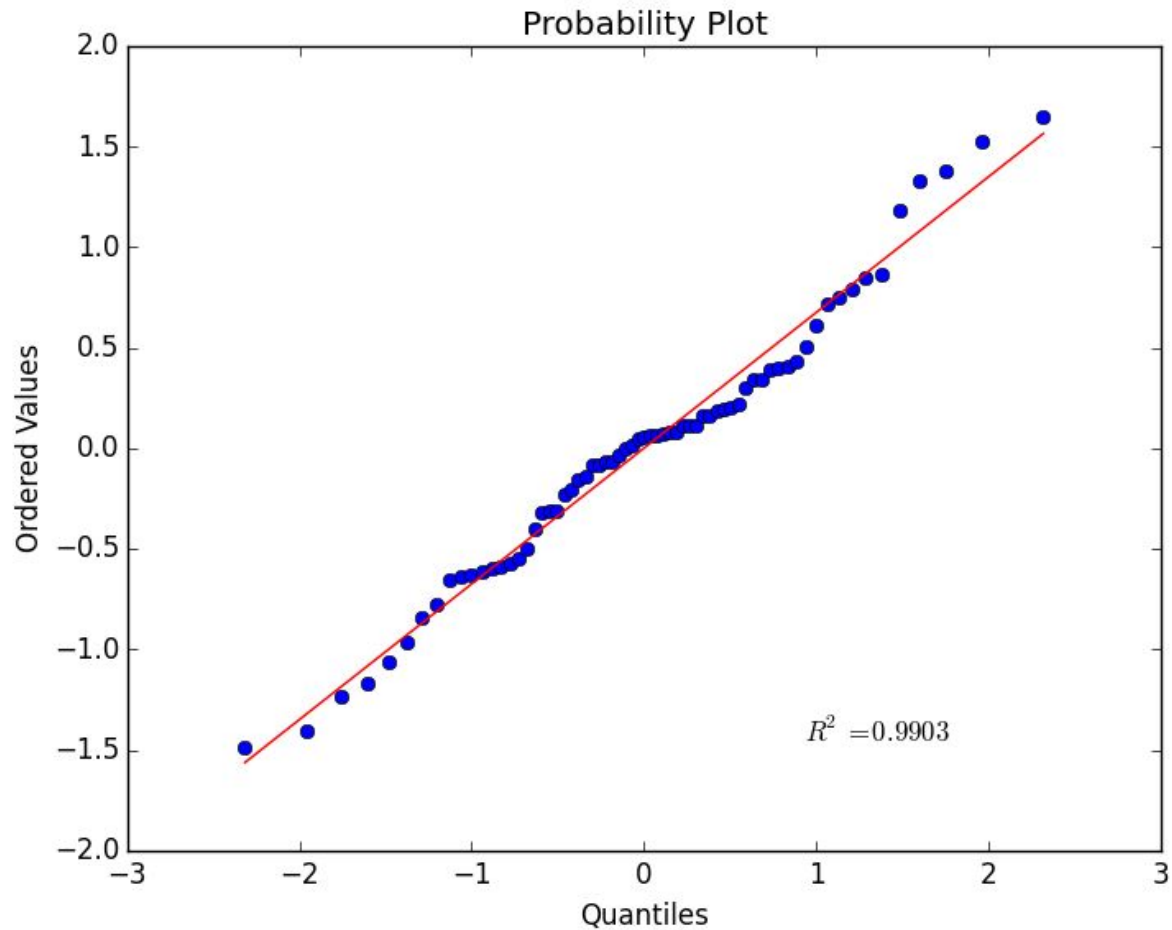
# REGRESION LINEAL ORDINARIA (LSS)

- t-Student con 58  $gl(67-8-1)$ :
- $a=0.05 \rightarrow 1.6716$
- `lcavol`
- `lweight`
- `svi`
- `age`, `gleason`, `pgg45` no tienen significancia.

TABLA DE PESOS Y Z-SCORE (TRAIN SET)

Atributo	Coefficiente b	Z-Score
Intercept	2.4649	27.3592
lcavol	0.6760	5.3198
lweight	0.2616	2.7269
age	-0.1407	-1.3838
lbph	0.2090	2.0380
svi	0.3036	2.4478
lcp	-0.2870	-1.8507
gleason	-0.0211	-0.1454
pgg45	0.2655	1.7227

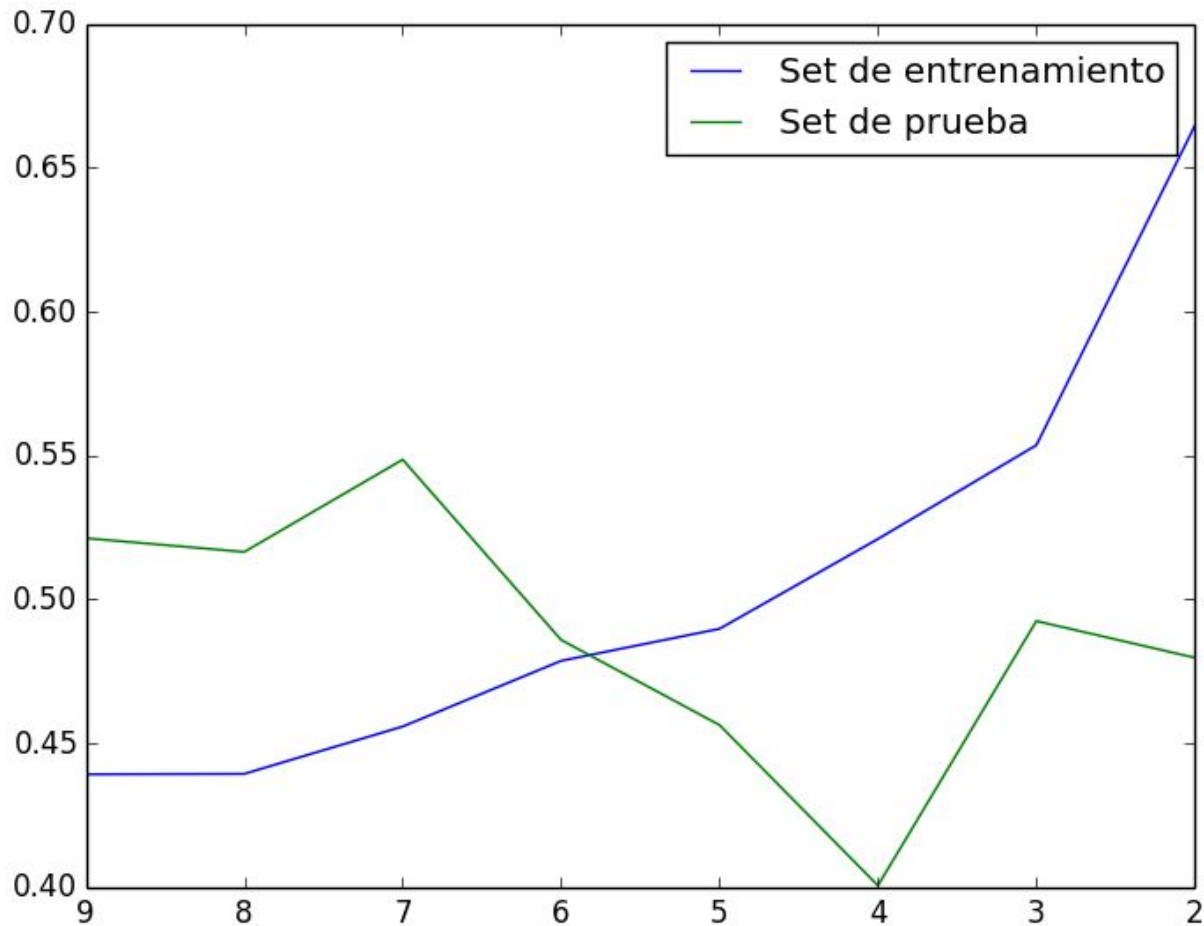
# NORMALIDAD DE RESIDUOS



# SELECCIÓN DE ATRIBUTOS

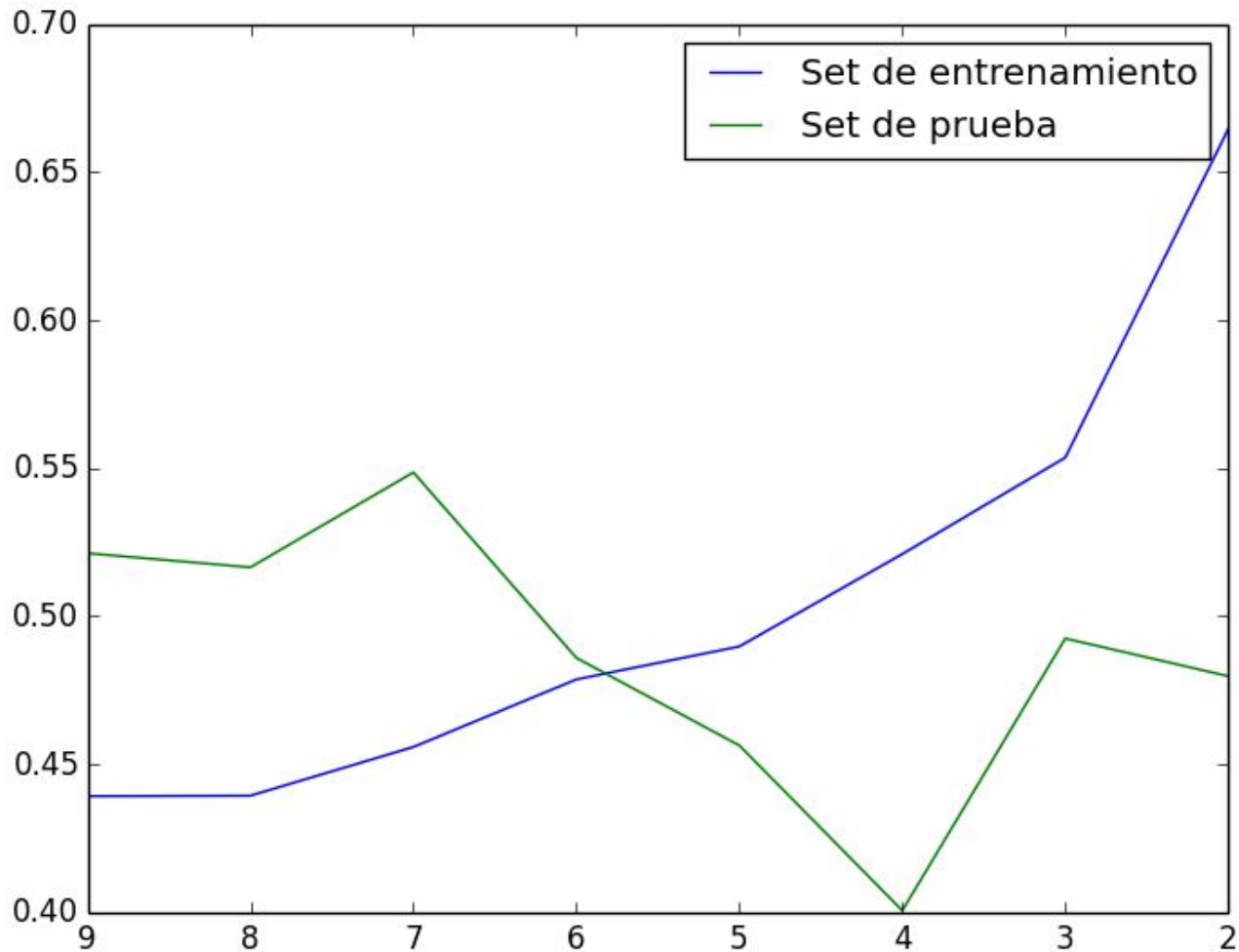
Se pide realizar Backward Step-wise y Forward Step-wise para selección de atributos, la primera consiste en partir con todas las variables y eliminar una a una aquella que tiene el menor peso, determinado por el Z-score, mientras que la segunda, busca partir desde la menor cantidad de parámetros hasta llegar a la mayor, agregando en cada iteración aquel más importante. La técnica se aplica al mismo dataset entregado para la pregunta 1 y se obtienen los siguientes resultados

# ERROR EN FUNCIÓN DEL NÚMERO DE PARÁMETROS(FSS)

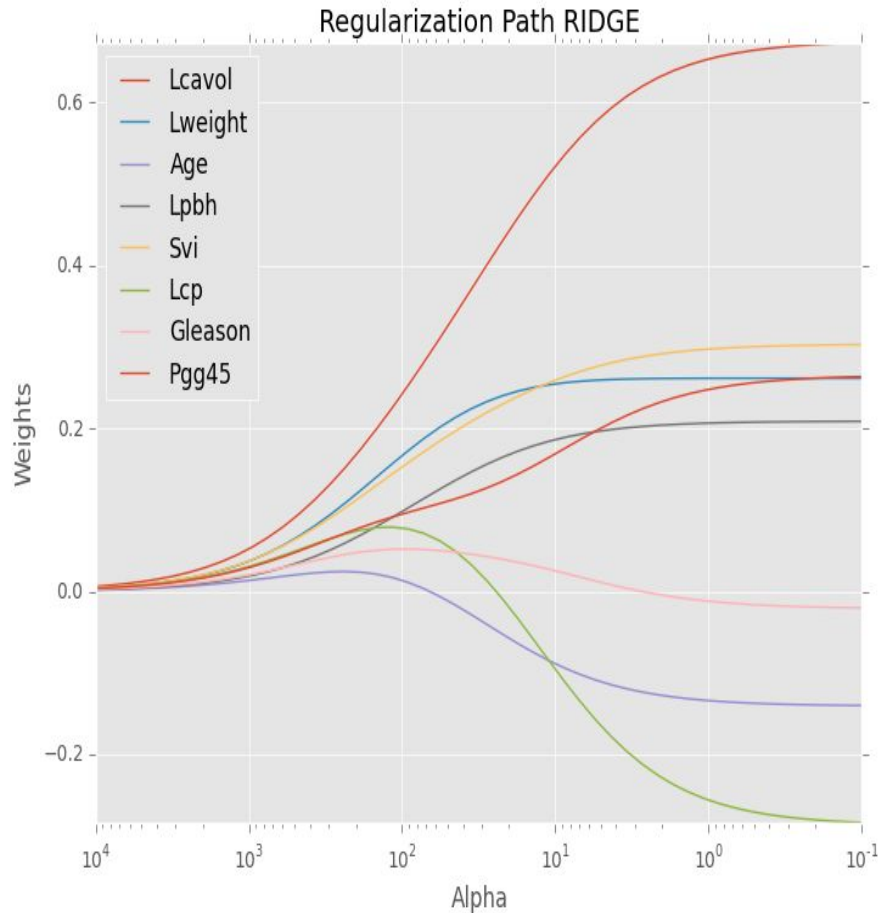




# ERROR EN FUNCIÓN DEL NÚMERO DE PARÁMETROS(BSS)

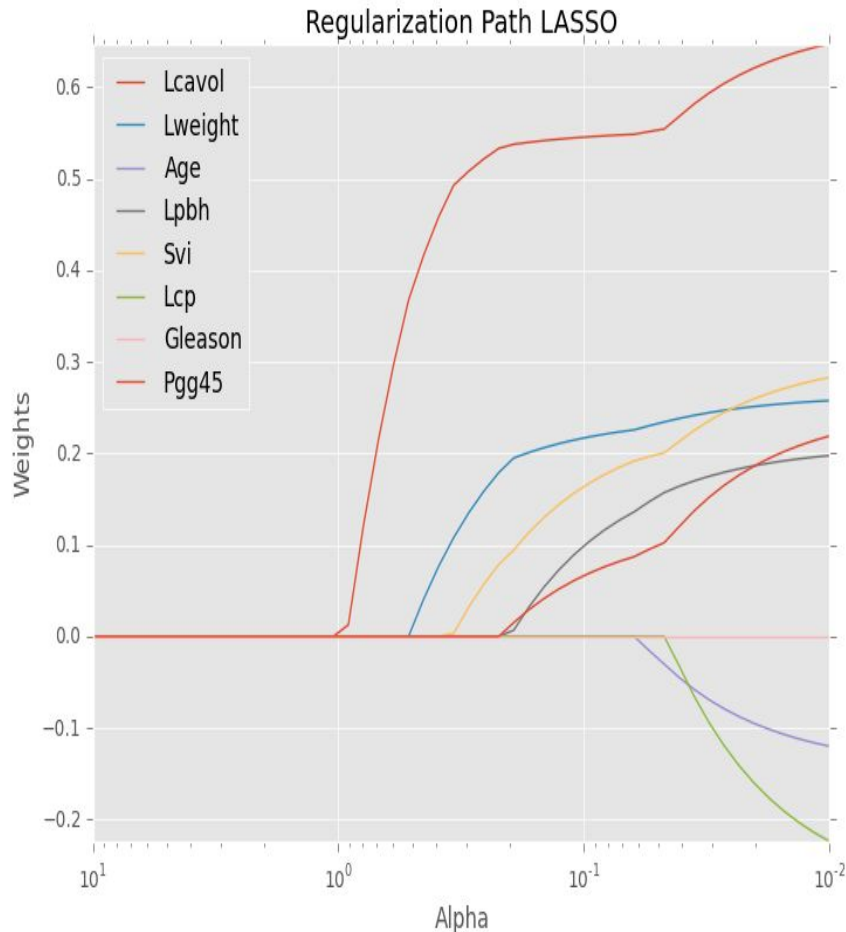


# COEFICIENTES DE RIDGE EN FUNCION DE PARAM DE REG



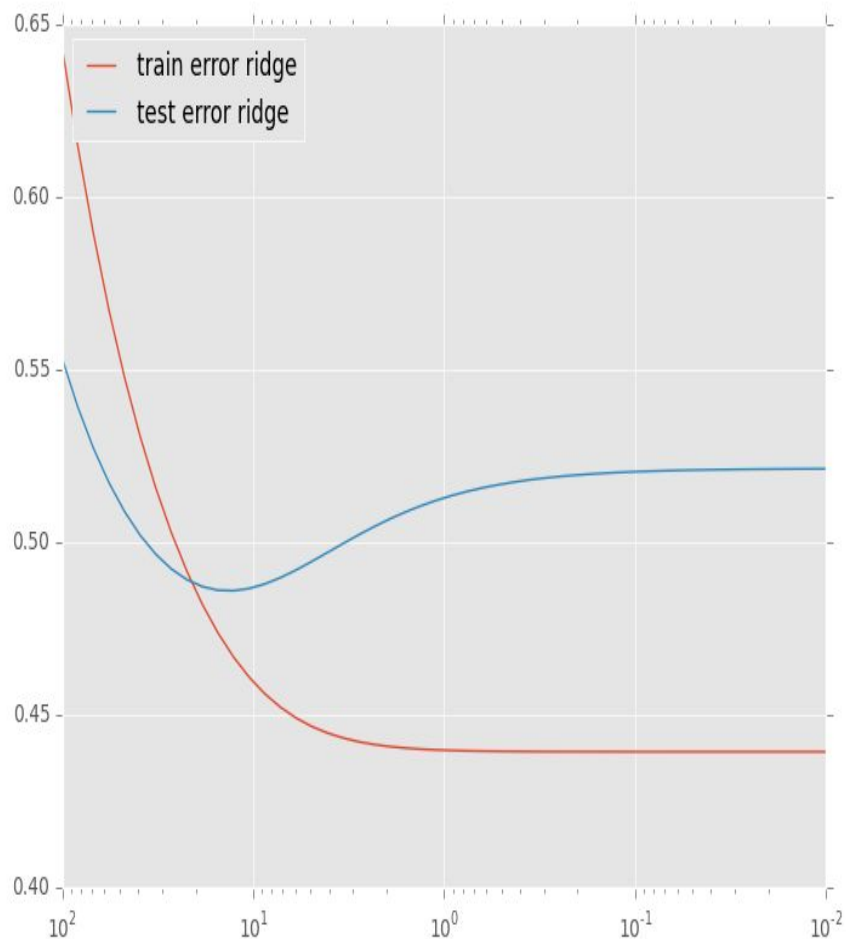
- Tiende a LSS con alpha pequeño.
- Alpha grande sobre-regulariza.
- Luego de  $10^3$  no hay significancia, pero nunca es 0.
- Poco efecto, hasta 10.

# COEFICIENTES DE LASSO EN FUNCION DE PARAM DE REG



- Lcavol, Lweight, svi siguen con mayor peso ( $\alpha < 0.5$ )
- $0.1 < \alpha < 1$  se hacen lweight, svi = 0.
- $\alpha > 1$  desaparece el modelo.
- $\alpha > 0.1$ , lcp y age desaparecen.
- Alpha grande sobre-regulariza.

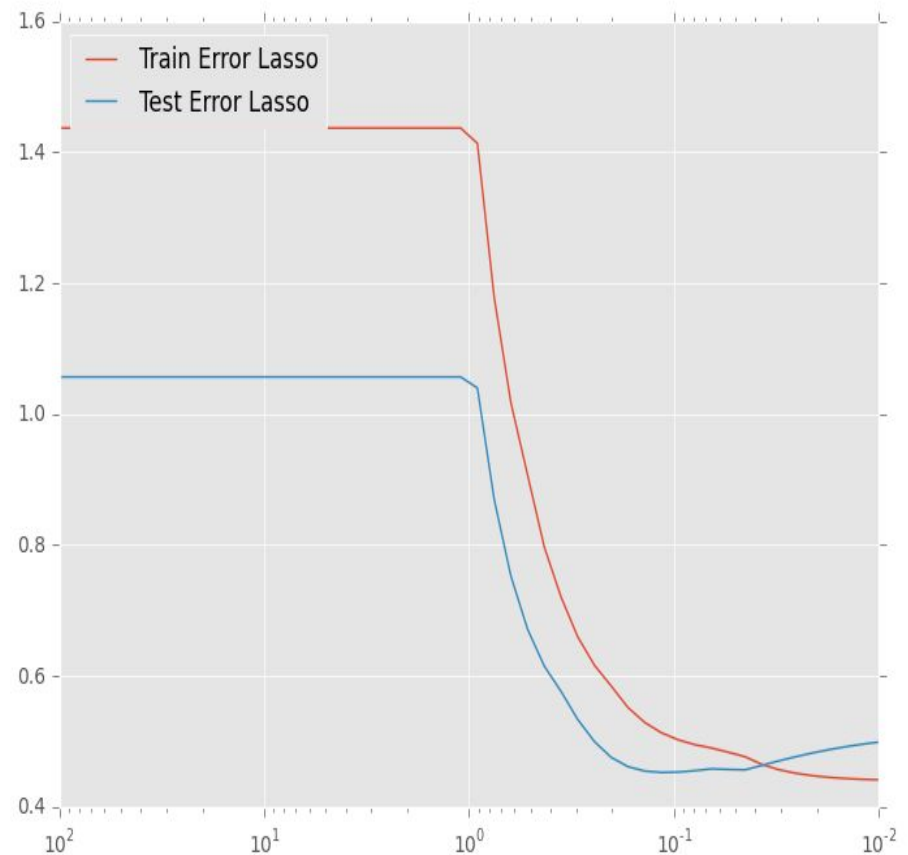
# COMPARACION ERROR TRAIN Y TEST, RIDGE REGRESSION



- $a > 1$  errores constantes
  - train  $\sim 0.53$
  - test  $\sim 0.44$
- error de train menor al de test ( $a < 30$ )
- error de test menor al de train ( $a > 30$ )
- Train siempre baja al disminuir  $\alpha$ .

# COMPARACION ERROR TRAIN Y TEST, LASSO REGRESSION

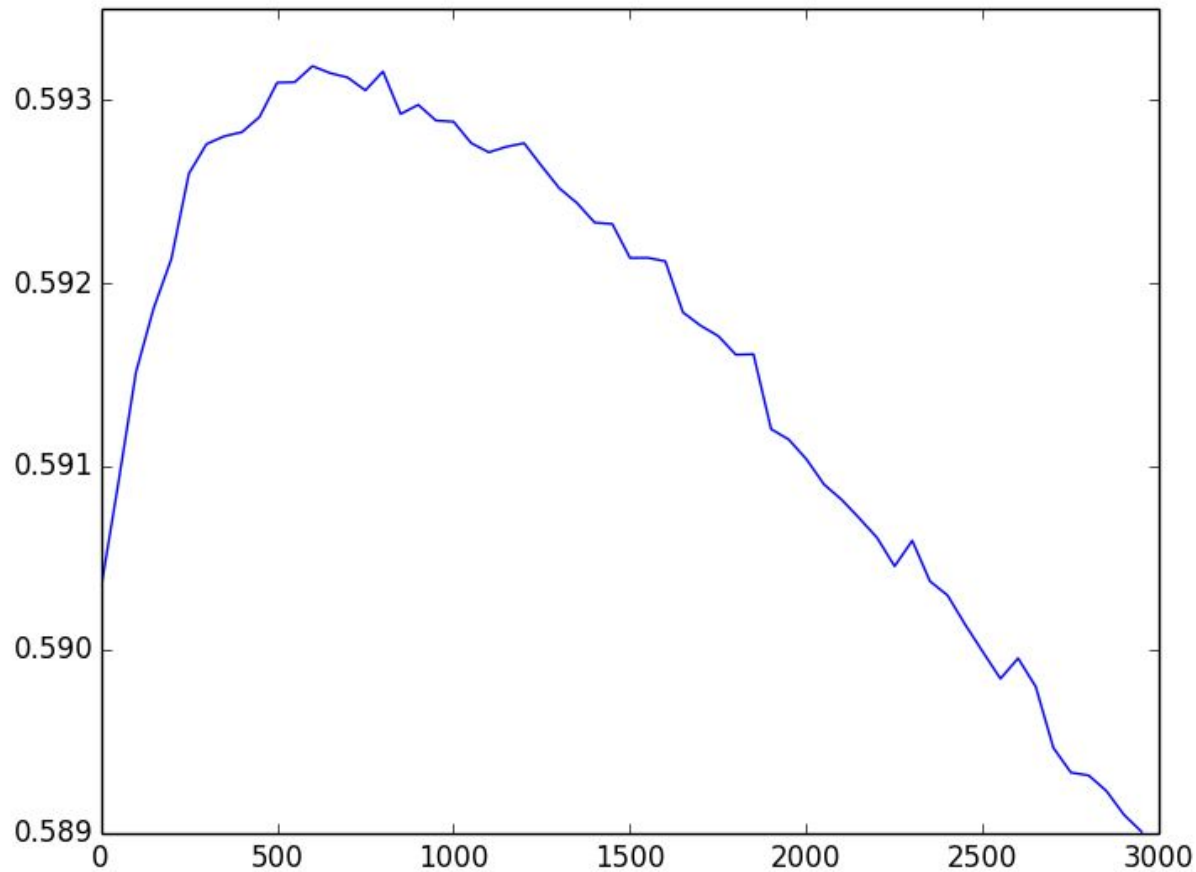
- Baja el MSE a medida que baja  $\alpha$ .
- $\alpha < 1$  errores disminuyen rápidamente.
- Error en  $\alpha > 1$ :
  - train  $\sim 1.42$
  - test  $\sim 1.03$
- Test más bajo hasta  $\alpha > 0.05$ .

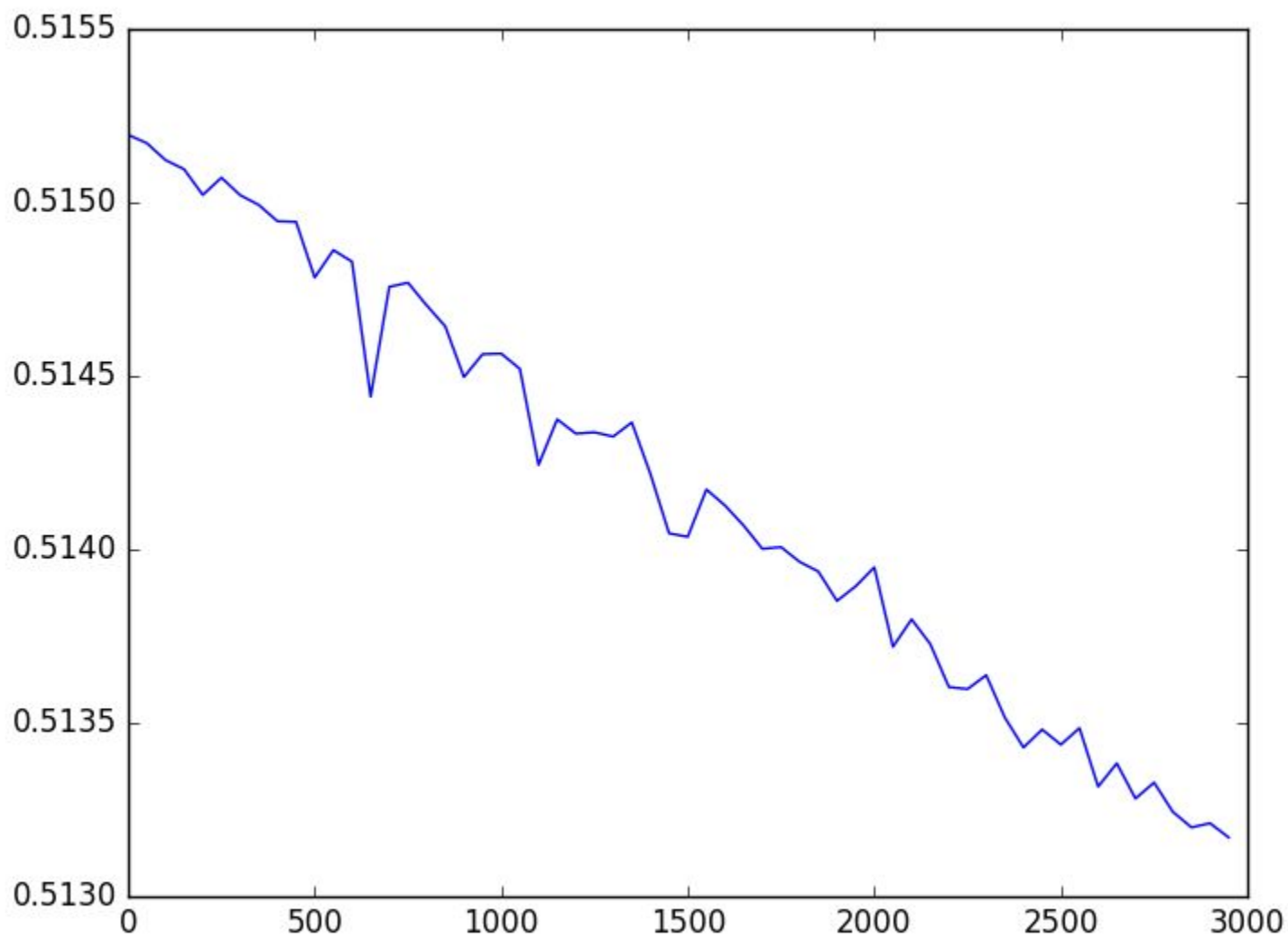


# PREDICCIÓN DE UTILIDADES DE UNA PELÍCULA

En esta parte, se intenta encontrar un predictor para las utilidades de una película, a través de ciertos atributos como lugar de origen, presupuesto, número de puntos de proyección, etc. Para esto, se probaron 3 métodos, el primero una regresión lineal sin regularizar, el segundo una regresión lineal por el método Lasso, y por último, una regresión lineal por el método Ridge.

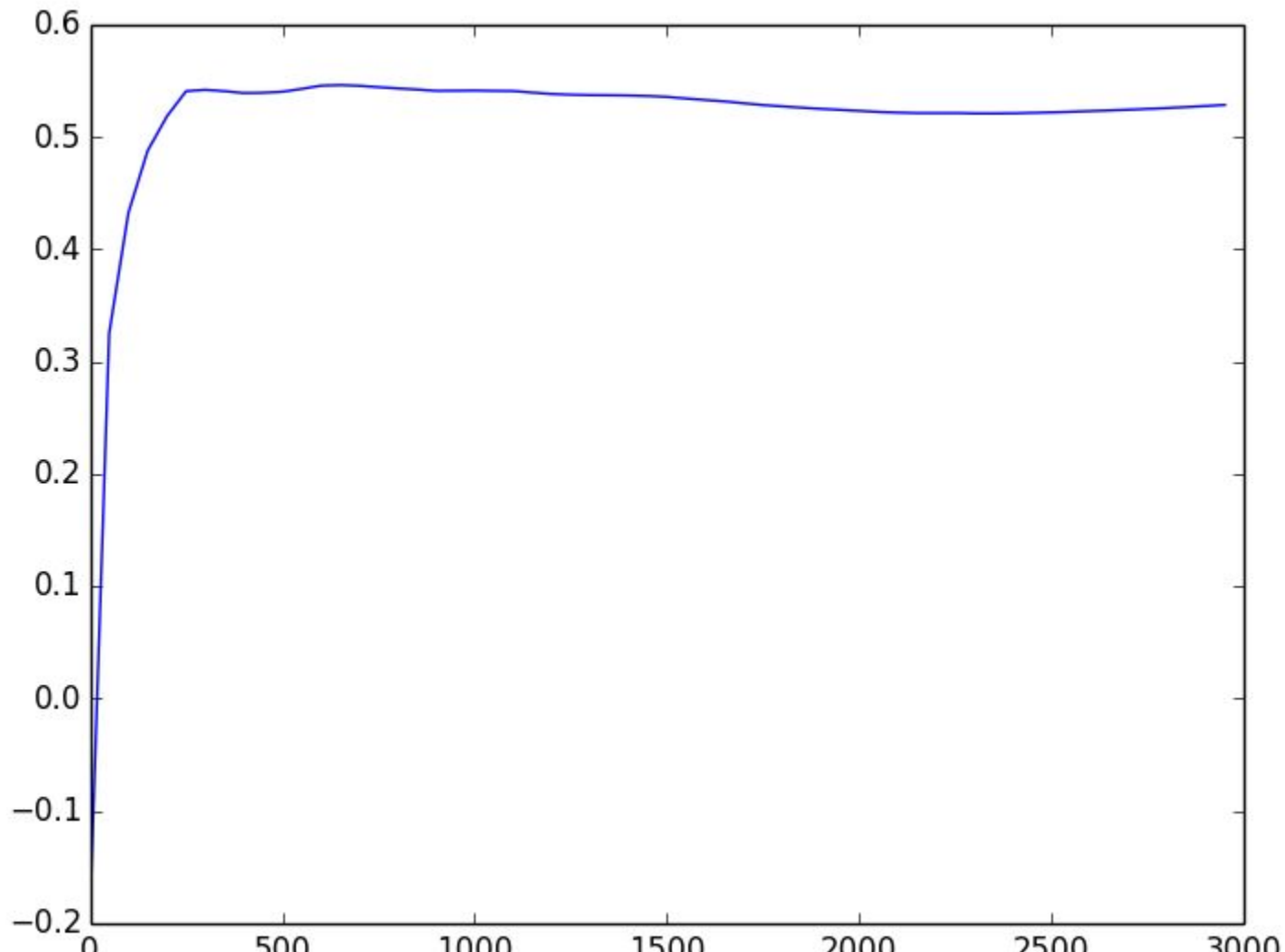
# COEFICIENTE DE DETERMINACIÓN EN FUNCIÓN DE ALFA (RIDGE)







# COEFICIENTE DE DETERMINACIÓN EN FUNCIÓN DE ALFA (LASSO)



GRACIAS POR SU ATENCIÓN!