



# ELASTIC NET «THE ORACLE»

## POOR MAN'S INTRODUCTION TO PENALIZED REGRESSION

Sinan İyisoğlu

# Plan

- OLS Regression
- Ridge Regression
- Lasso
- Bridge Regression
- Elastic Net
- Uygulama

# $L_q$ Norm Ailesi

- n boyutlu bir  $x = (x_1, x_2, \dots, x_n)$  vektörü için

$$L_1 \text{ normu } \|x\|_1 = |x_1| + |x_2| + \dots + |x_n|,$$

$$L_2 \text{ normu } \|x\|_2 = (x_1^2 + x_2^2 + \dots + x_n^2)^{1/2},$$

$$L_q \text{ normu } \|x\|_q = (x_1^q + x_2^q + \dots + x_n^q)^{1/q}$$

şeklinde tanımlanır.

# OLS Regression

En küçük kareler regresyonunda amacımız  $\mathbf{y}$  ve  $\mathbf{X}$  verildiğinde

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

eşitliğini sağlayan  $\boldsymbol{\beta}$  katsayılarını bulmaktır. Bunu yapmak için hata kareler toplamını minimize ederiz.  $\mathbf{X}$  full rank ise tek çözüm vardır.

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$$

$\|\cdot\|_2$  Euclid normu ( $L_2$  norm)

# OLS Regression

Başka bir ifadeyle  $p$  değişken ve  $n$  gözlem olduğunda

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2$$

Tam çoklu bağlantı olduğunda  $\mathbf{X}$  full rank olmadığından tek çözüm yoktur.

$p \gg n$  olduğunda yine  $\mathbf{X}$  full rank olmayabilir.

# OLS Regression

- OLS tahminleri genellikle yetersizdir.
- OLS tahminlerinde bias azdır fakat varyans fazladır.
- Bazı katsayılar küçültülerek ya da sıfırlanarak tahmin geçerliliği (prediction accuracy) arttırılabilir.
- Alt küme seçimi yapmak da bir yöntemdir. Fakat bu yöntemde sabit modeller elde etmek zordur.



# Ridge Regression- $L_2$ Penalty - 1970

- OLS deki parametre kestirimleri yansızdı. Bazen kestirimlerin yanlı olması daha iyi tahmin değerleri elde etmeyi sağlar.
- OLS üzerine bir  $L_2$  penalty konmasıyla Ridge regression elde edilir.

$$\hat{\beta} = \arg \min_{\beta} \|y - X\beta\|_2^2 + \|\lambda\beta\|_2^2 \text{ ya da}$$

$$\hat{\beta} = \arg \min_{\beta} \|y - X\beta\|_2^2 \text{ s. t. } \|\lambda\beta\|_2^2 < t$$



# Ridge Regression

- Ridge tahminleri daha iyidir.
- Çoklu bağlantı sorunu çözülür.
- Katsayılar küçülmüştür. (Shrinkage)
- MSE, OLS dekine göre daha küçük olabilir.
- Katsayılar değişkenlerin birimlerine duyarlı olduğundan standartlaştırma yapılır.

Fakat

- Tüm değişkenler modelde kalmıştır.



# The Lasso - $L_1$ Penalty - 1996

- Least Absolute Shrinkage and Selection Operator
- Shrinkage + Subset Selection

$$\hat{\beta} = \arg \min_{\beta} (\|y - X\beta\|_2^2 + \|\lambda\beta\|_1)$$

$\|\cdot\|_1$  Mutlak değer ( $L_1$  norm)



# Lasso

- Bazı katsayılar sıfırlanmıştır. (Sparse solution)
- $L_1$  normu olduğundan türevlenebilirlik kaybolmuştur ve sonucun açık formülü yoktur.
- Bu yüzden algoritmalar geliştirilmiştir.

1996 Quadratic Programming

2003 LARS

2008 Coordinate Descent

# Lasso

- $p > n$  durumunda Lasso en çok  $n$  tane değişken seçer.
- İkili korelasyonu yüksek bir grup değişken içerisinde sadece birini alır, hangisini seçtiğini önemsemez.
- $n > p$  olduğunda değişkenler arasında yüksek korelasyon varsa Ridge regresyonun Lassoya göre daha iyi tahmin performansı olduğu gözlenmiştir.

# Gen Seçimi

- Gen seçimi durumunda  $p \gg n$  dir.
- Bazı genlerin yüksek korelasyonlu gruplar halinde aynı biyolojik yolu paylaştığı düşünülür.
- İdeal gen seçiminde önemsiz genler silinmeli, gruptan bir gen seçilmişse diğerleri de modele girmelidir.
- Lasso bu durum için ideal değildir.

# Bridge Regression - $L_q$ Penalty 1993

- Ridge ve Lasso'nun genelleştirilmiş hali

$$\hat{\beta} = \arg \min_{\beta} (\|y - X\beta\|_2^2 + \|\lambda\beta\|_q^q)$$

$\|\cdot\|_q$   $L_q$  norm

- $1 < q < 2$  için Elastic Net'e benzer.
- Ridge de olduğu gibi değişken seçimi yapmaz.

# Elastic Net - $L_1$ & $L_2$ Penalty - 2004

- $L_1$  kısmı değişken seçimi yapar
- $L_2$  kısmı
  - (i) Seçilen değişkenlerdeki n adet sınırını kaldırır
  - (ii) Gruplama etkisini güçlendirir



$$\hat{\beta} = \arg \min_{\beta} (\|y - X\beta\|_2^2 + \|\lambda_1 \beta\|_1 + \|\lambda_2 \beta\|_2^2)$$



# Naive Elastic Net

$$\alpha = \frac{\lambda_2}{(\lambda_1 + \lambda_2)} \quad \text{seçilirse}$$

$$\hat{\beta} = \arg \min_{\beta} (\|y - X\beta\|_2^2) \quad s. t. (1 - \alpha)\|\beta\|_1 + \alpha\|\beta\|_2^2 \leq t$$

- $\alpha=1$  durumunda Ridge regression,  $\alpha=0$  olma durumunda Lasso olur.
- Elastic Net, Ridge ve Lasso'nun konveks bir birleşimidir.

# Elastic Net

- Naive Elastic Net bir Lasso problemi olarak yazılabilir.

$$\hat{\beta}^* = \mathop{\text{arg min}}_{\beta^*} (\|y^* - X^* \beta^*\|_2^2 + \frac{\lambda_1}{\sqrt{1 + \lambda_2}} \|\beta^*\|_1)$$

- Elastic Net katsayıları ile Naive Elastic Net katsayıları arasındaki ilişki aşağıdaki gibidir.

$$\hat{\beta}(enet) = (1 + \lambda_2) \hat{\beta}(naive\ enet)$$

- Elastic Net, Naive Elastic Net'e göre daha performanslıdır.

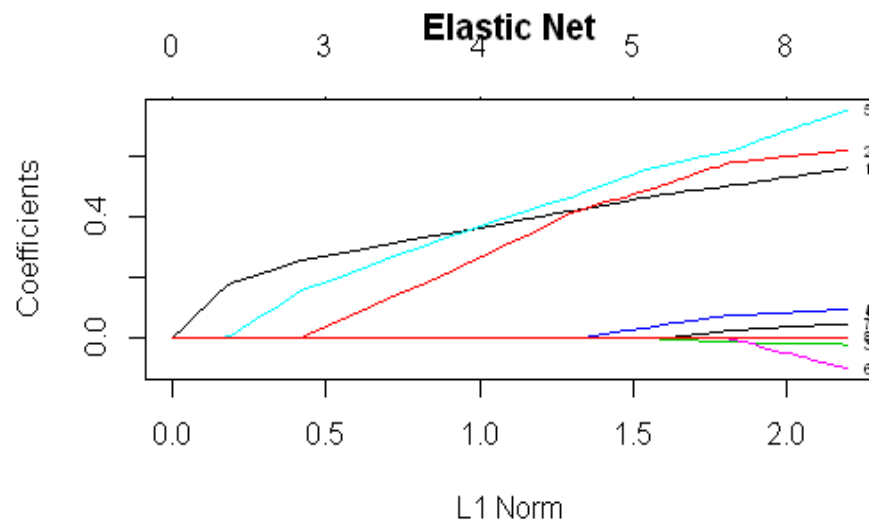
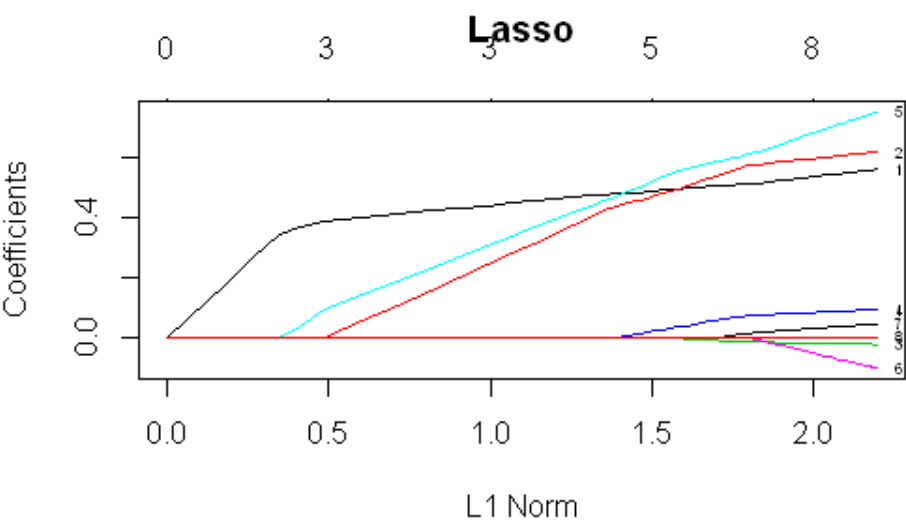
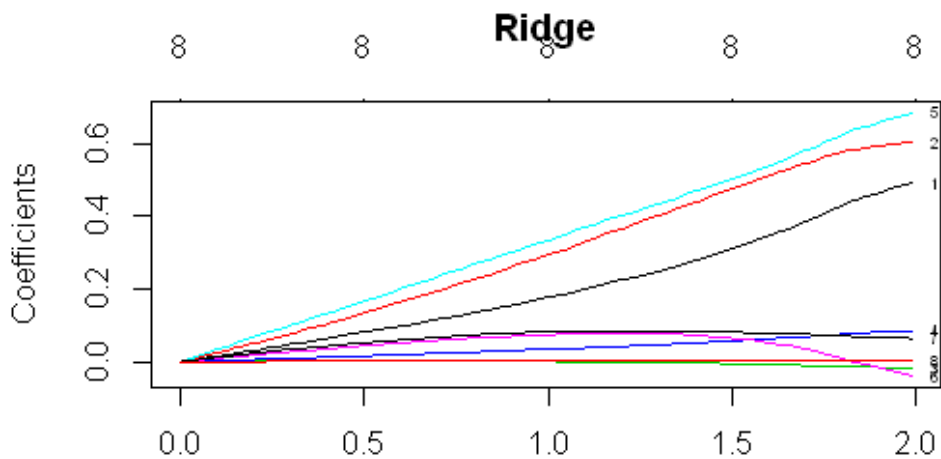
# Elastic Net

- Elastic Net iyi tahmin performansına sahip, grup etkisi içeren bir model üretir.
- Elastic Net alt küme seçimi yapar (Sparse solution).
- LARS ve Coordinate descent algoritmaları ile hesaplanabilir.
- Lassodan daha performanslıdır.
- Model için iki parametrenin seçilmesi gereklidir.

# Uygulama

- Prostat kanseri verisi
- Bir bağımlı değişken logpsa ve 8 adet bağımsız değişken
- R glmnet paketi
- Elastic net için  $\alpha=0.5$  seçildi.

# Coefficient Paths



# Katsayılar

Ridge	Lasso	ENet
(Intercept) 0.012	(Intercept) 0.154	(Intercept) 0.120
lcavol 0.492	lcavol 0.507	lcavol 0.497
lweight 0.604	lweight 0.546	lweight 0.566
age -0.017	age -0.008	age -0.011
lbph 0.086	lbph 0.062	lbph 0.070
svi 0.685	svi 0.590	svi 0.608
lcp -0.040	lcp .	lcp .
gleason 0.064	gleason 0.001	gleason 0.021
pgg45 0.003	pgg45 0.002	pgg45 0.002



# MSE

