

# **glmnet**

Kate Cook

ML4Bio software presentation

2010.02.04

# Features

- Linear regression library for R
- Makes regression models and predictions from those models
- Lasso and elastic net regression via coordinate descent (Friedman 2010)
- Very fast
  - FORTRAN-based
  - exploits sparsity in input data
- Simple to use

# Availability & installation

- `install.packages("glmnet")`
- GPL licensed
- Citation, manual etc:

<http://cran.r-project.org/web/packages/glmnet/index.html>

– Or just google “glmnet”...

# Regularization (review?)

- p features, n observations
- $y = X\beta + \varepsilon$
- Want to minimize the sum of squared errors:

$$\sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=1}^n (y_i - x_i' \hat{\beta})^2$$

- To reduce overfitting, add a penalty term
- Now we minimize:

$$\sum_{i=1}^n (y_i - x_i^T b)^2 + \lambda P_{\alpha}(b_1, \dots, b_p)$$

## Ridge, LASSO, and elastic net regularization are related

- Ridge regression, LASSO, and elastic net are part of the same family with penalty term:

$$P_{\alpha} = \sum_{j=1}^p \left[ \frac{1}{2} (1 - \alpha) b_j^2 + \alpha |b_j| \right]$$

- $\alpha = 0 \rightarrow$  ridge regression
- $\alpha = 1 \rightarrow$  LASSO
- $0 < \alpha < 1 \rightarrow$  elastic net!

# Features of LASSO and elastic net regularization

- Ridge regression shrinks correlated variables toward each other
- LASSO also does feature selection
  - if many features are correlated (eg, genes!), lasso will just pick one
- Elastic net can deal with grouped variables

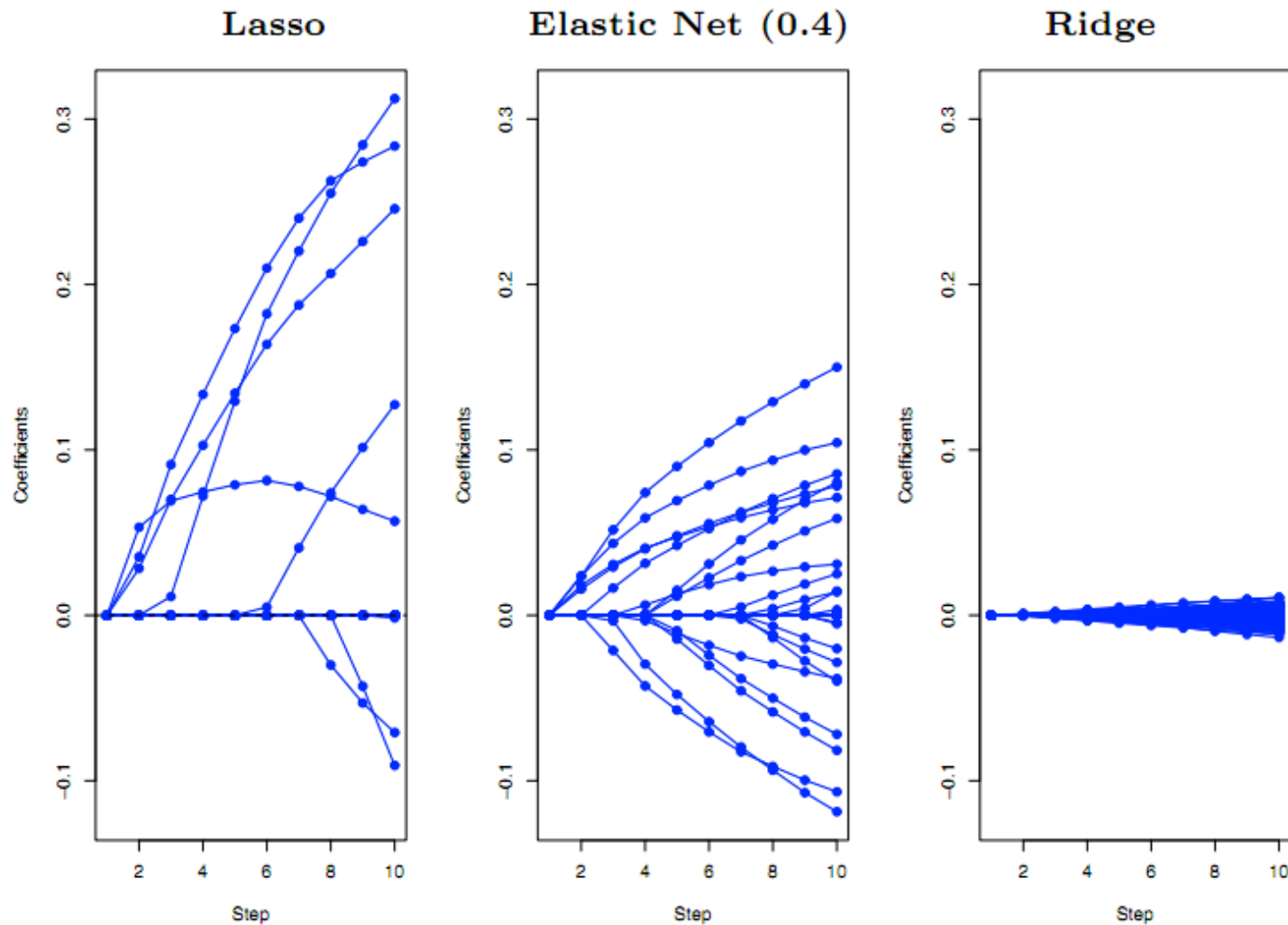
# One more detail

- Elastic net formulation above is actually the “naïve elastic net”
  - Doesn’t perform well in practice
  - Parameters are penalized twice
- How to fix it?

$$\begin{aligned}\text{Penalty} &= (1 - \alpha) |\beta|_1 + \alpha |\beta|^2 \\ &= \lambda_2 |\beta|^2 + \lambda_1 |\beta|_1 \quad \text{where} \quad \alpha = \frac{\lambda_2}{\lambda_1 + \lambda_2}\end{aligned}$$

$$\hat{\beta}(\text{elastic net}) = (1 + \lambda_2) \hat{\beta}(\text{naive elastic net}).$$

# Graphically...



Hastie, <http://www-stat.stanford.edu/~hastie/TALKS/glmnet.pdf>



# Example – intro

- Determining RNA sequence features predictive of binding to an RNA-binding protein
- Apply LASSO regression to model binding
- Use cross-validation to select the best  $\lambda$

$$\sum_{i=1}^n (y_i - x_i^T b)^2 + \lambda P_{\alpha}(b_1, \dots, b_p)$$

- Train model on first 10,000 points, test on last ~5,000

# Example - data

Data frame loaded from text file

```
> head(data_ML4bio_2)
```

|   | ratio       | E_2_4 | E_2_5 | E_3_4 | E_3_5 | E_4_4 | E_4_5 |
|---|-------------|-------|-------|-------|-------|-------|-------|
| 1 | 1.42614954  | -2.68 | -2.29 | -3.43 | -4.53 | -3.65 | -5.96 |
| 2 | 1.27598858  | -2.60 | -2.18 | -3.56 | -2.25 | -3.71 | -2.99 |
| 3 | 0.57953823  | -4.41 | -1.07 | -6.48 | -2.09 | -8.05 | -2.23 |
| 4 | -0.22087992 | -1.92 | -1.76 | -2.02 | -1.96 | -2.21 | -2.18 |
| 5 | -0.06225514 | -1.92 | -1.76 | -2.02 | -1.96 | -2.21 | -2.18 |
| 6 | 0.18564163  | -1.40 | -2.56 | -1.58 | -5.81 | -2.11 | -7.88 |

...

Affinity measurement (y)

Features (x)

- Note: features must be numeric (use dummy variables for categorical data)

# Example – fitting the model

- Syntax:

```
Fit <- cv.glmnet(X, y, ...)
```

```
> fit1_mmodel.cv<-cv.glmnet(mmodel_ML4bio[1:10000,],  
data_ML4bio$ratio[1:10000], alpha=1)
```

  
Y vector

  
LASSO

  
X matrix

# Example – making predictions

- Syntax:

```
Pred <- predict(Fit, newX)
```

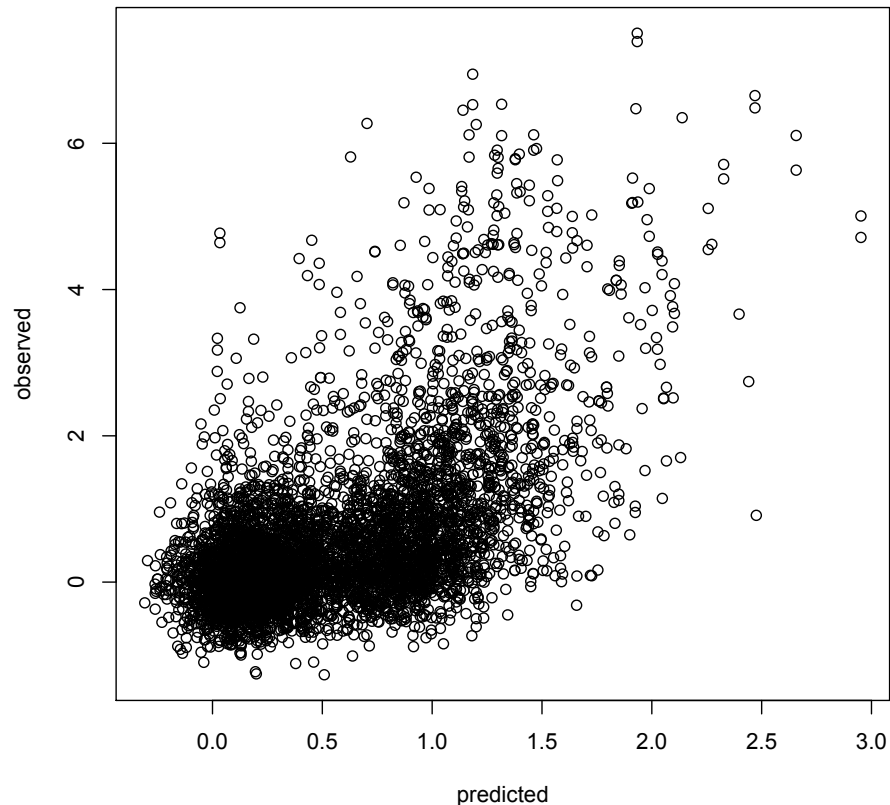
```
> pred_fit1<-predict(fit1_mmodel.cv, mmodel_ML4bio  
[10001:15490,])
```

```
> cor(pred_fit1, data_ML4bio$ratio[10001:15490])  
[ ,1]
```

```
1 0.5660715
```

# Example – comparing predicted to observed

```
> plot(pred_fit1, data_ML4bio$ratio  
[10001:15490], xlab="predicted", ylab="observed")
```



# Stuff that helped me understand how this works

- <http://www-stat.stanford.edu/~hastie/TALKS/glmnet.pdf>
  - Theory behind LARS and coordinate descent, speed trials, biological examples
- Friedman, Hastie & Tibshirani, *Regularization Paths for Generalized Linear Models via Coordinate Descent*, J Stat Soft, 2010
- Zou and Hastie, *Regularization and Variable Selection via the Elastic Net*, J Royal Stat Soc B, 2005

