

Regularization in Deep Neural Networks

Paul Stey

March 3, 2017

Table of Contents

- 1 Introduction
 - Regularization
 - Penalized Regression

- 2 Broader View
 - Non-Obvious Examples

What is Regularization?

- ① Method of calibrating the complexity of our model
- ② Very general (i.e., applicable to variety of models)
- ③ Penalized regression
 - 1 Ridge regression
 - 2 Lasso
 - 3 Elastic net
- ④ XGBoost improvement on AdaBoost

L_2 Regularization Regression

- 1 Ridge regression, also called Tikhonov regularization (Tikhonov, 1963)
- 2 Penalize the L_2 norm
- 3 Constrains the Euclidian distance of β

Given the linear model

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$$

the ridge penalty constrains

$$\|\beta\|_2 = \sqrt{\beta_1^2 + \beta_2^2 + \dots + \beta_p^2}$$

L_2 Regularization cont.

Recall the least squares solution is obtained by

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y.$$

The ridge penalty gives us

$$\hat{\beta}^{ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T y$$

where λ is a penalty term, and \mathbf{I} is the $p \times p$ identity matrix.

L_2 Regularization cont.

- ① Penalty λ is in the interval $[0, \infty)$
- ② Shrinks β_j values towards 0 (and each other)
- ③ Benefits:
 - 1 Reduce model variance (i.e., more generalizable)
 - 2 Gain stability in estimates involving correlated data
 - 3 Computationally efficient
- ④ Drawbacks:
 - 1 Sacrifice unbiasedness of maximum likelihood estimate

L_1 Regularization

The lasso

- ① Similar to ridge regression
- ② Penalize L_1 norm of β
- ③ Constrains taxicab distance spanned by vector of regression coefficients

$$\|\beta\|_1 = |\beta_1| + |\beta_2| + \dots + |\beta_p|$$

- ④ As λ increases, $\beta_j \rightarrow 0$
 - 1 β_j for less important predictors shrink faster
 - 2 Can be used for “variable selection”

Elastic Net

- 1 Penalize *both* the L_1 and L_2 norms of β
- 2 Advantage of combining strengths of both approaches

Caveat: Note that we often use ridge or lasso for very specific reason (i.e., ridge for accuracy, lasso for variable selection). One might suggest we give away this clarity of purpose when using the elastic net.

Regularization More Broadly

Recall that regularization is merely a method of controlling model complexity.

The standard shrinkage method examples make this very transparent with the use of explicit penalty terms.

Broader View cont.

Consider less obvious examples:

- 1 Number of boosting iterations
- 2 Learning rate for gradient descent
- 3 Bagging proportion
- 4 Proportion in training vs. test set