

# Regularization in Deep Neural Networks

Paul Stey

Isabel Restrepo

March 3, 2017

# Table of Contents

- 1 Introduction
  - Regularization
  - Penalized Regression
- 2 Norm Penalties
  - Norm Penalties in Neural Networks
- 3 Broader View
  - Non-Obvious Examples
- 4 More Regularization Methods
  - Drop Out
  - Data Augmentation
  - Adversarial Training
- 5 Conclusion
  - Summary
  - Other Methods

# What is Regularization?

Method of calibrating the complexity of our model

# What is Regularization?

- ① Very general (i.e., applicable to variety of models)
- ② Penalized regression
  - 1 Ridge regression
  - 2 Lasso
  - 3 Elastic net
- ③ 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  $\beta$

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

# $L_2$ Regularization cont.

- ❶ Penalty  $\lambda \in [0, \infty)$
- ❷ Shrinks  $\beta_j$  values towards 0 (and each other)
- ❸ Benefits:
  - 1 Reduce model variance (i.e., more generalizable)
  - 2 Computationally efficient
  - 3 Gain stability in estimates involving correlated data

[\[example in R\]](#)

- ❹ Drawbacks:
  - 1 Sacrifice unbiasedness of maximum likelihood estimate

# $L_1$ Regularization

## The lasso

- 1 Similar to ridge regression
- 2 Penalize  $L_1$  norm of  $\beta$
- 3 Constrains Manhattan distance spanned by vector of regression coefficients

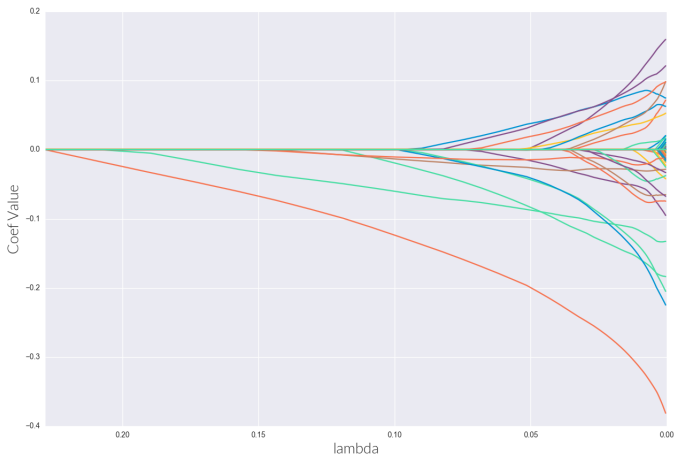
$$\|\beta\|_1 = |\beta_1| + |\beta_2| + \dots + |\beta_p|$$

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



# Drop Out cont.

## Coefficient Path



# Elastic Net

- 1 Penalize *both* the  $L_1$  and  $L_2$  norms of  $\beta$
- 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.

# Terminology

Many papers and texts use the term “weight decay” to refer norm penalization methods in neural networks

# $L_2$ Norm Penalty

- ① Same basic idea as ridge regression
- ② Cost function gets penalty term  $\lambda$  to constrain weights

$$C = C_0 + \frac{\lambda}{2} \sum_i w_i^2$$

where  $C_0$  is the original cost.

- ③ Stabilizes weights for highly correlated features

# $L_1$ Norm Penalty

- ① Same idea as the lasso
- ② Add penalty to cost function, thereby shrinking weights

$$C = C_0 + \lambda \sum_i |w_i|$$

- ③ Some weights will be shrunk to 0
- ④ Produces a more sparse model

# 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

# Broader View cont.

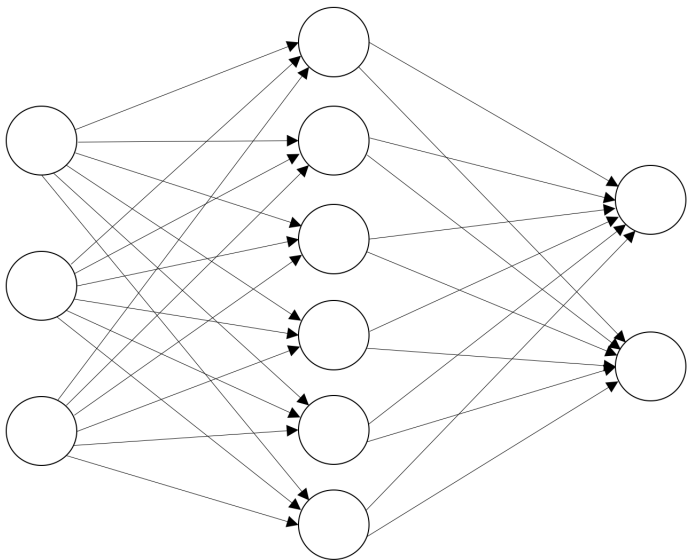
The essential feature of regularization is that it involves using some mechanism to control the complexity of our model. This makes our model more generalizable.



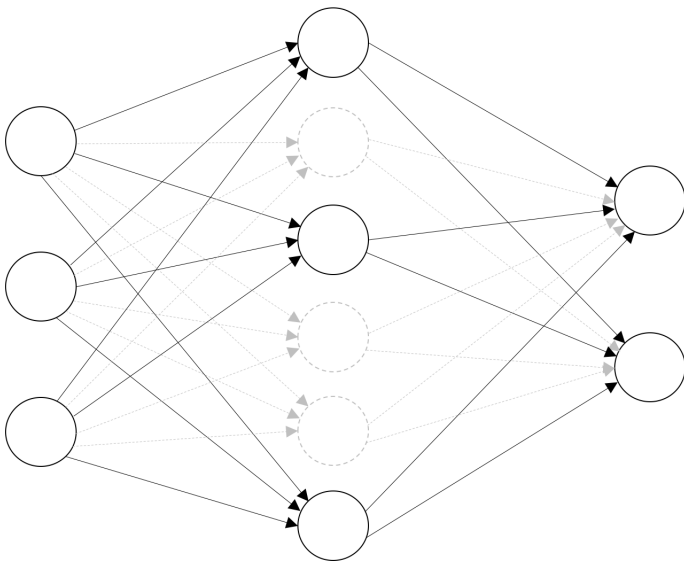
# Drop Out

- 1 Hinton *et al.* 2012
- 2 Subset of hidden and input units are randomly omitted
- 3 Reduces number of weights and biases being estimated
- 4 Thereby reducing complexity of model

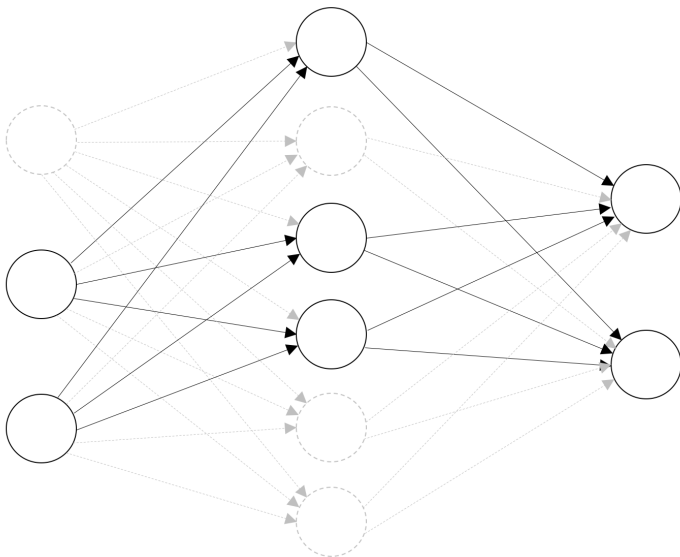
# Drop Out cont.



# Drop Out cont.



# Drop Out cont.



# Drop Out in Practice

- ① Units omitted according to meta-parameters we specify
  - 1 Input unit probability of inclusion frequently 0.8
  - 2 Hidden unit 0.5
- ② In most cases, omission means multiplying activation by 0

# Compare with Bagging and Random Forests

Consider similarities with bagging and feature subset selection in random forests.

These methods have a regularizing effect on model by using randomization to artificially restrict the parameters.

# Data Augmentation

- ➊ Adding [synthetic] data
- ➋ Or perturbing existing data
- ➌ Or both!

# Data Augmentation Examples

- 1 Translating images by a few pixels
- 2 Injecting noise in speech recognition data



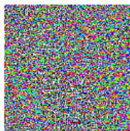
# Adversarial Training

- 1 Szegedy *et al.* 2014
- 2 Similar in a sense to data augmentation
- 3 Generate new observations in such a way as to maximize the chance they are misclassified
- 4 These in-domain examples that have been altered are referred to as adversarial examples
- 5 Related to “rubbish class” examples (LeCun *et al.*, 1998)

# Adversarial Training Example



+ .007 ×



=



$x$

$y$  = “panda”  
w/ 57.7%  
confidence

$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”  
w/ 8.2%  
confidence

$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”  
w/ 99.3 %  
confidence

Goodfellow, Bengio, & Courville, 2016

# Methods we Did Not Discuss

- 1 Early stopping
- 2 Bagging (Breiman, 1994)
- 3 Parameter sharing (Laserre *et al.*, 2006)

## *TensorFlow examples*