Lecture 14: Regularization CS 109B, STAT 121B, AC 209B, CSE 109B

Mark Glickman and Pavlos Protopapas



Lecture 3 Regularization Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error

Outline

- Norm Penalties
- Early Stopping
- Data Augmentation
- Bagging
- Dropout

Norm Penalties



– MAP estimation with Laplacian prior

L₂ Regularization



Norm Penalties as Constraints

$$\min_{\Omega(\theta) \leq K} J(\theta; X, y)$$

- Useful if *K* is known in advance
- Optimization:
 - Construct Lagrangian and apply gradient descent
 - Projected gradient descent

Early Stopping



Early Stopping ≈ Weight Decay



Goodfellow et al. (2016)

Sparse Representations

 Weight decay on activations instead of Output of hidden layer parameters 0 $\begin{bmatrix} -14\\1\\19\\2\\23\end{bmatrix} = \begin{bmatrix} 3 & -1 & 2 & -5 & 4 & 1\\4 & 2 & -3 & -1 & 1 & 3\\-1 & 5 & 4 & 2 & -3 & -2\\3 & 1 & 2 & -3 & 0 & -3\\-5 & 4 & -2 & 2 & -5 & -1 \end{bmatrix} \begin{bmatrix} 0\\2\\0\\0\\-3\\0\end{bmatrix}$ $oldsymbol{B} \in \mathbb{R}^{m imes n}$ $oldsymbol{y} \in \mathbb{R}^m$ $oldsymbol{h} \in \mathbb{R}^n$ Weights in output layer

 $J(\theta; X, y) + \alpha \Omega(h)$

Data Augmentation



deeplearningbook.org

Noise Robustness

- Random perturbation of network weights
 - Gaussian noise: Equivalent to minimizing loss with regularization term $\mathbf{E}[\|\nabla_{W}y(x)\|]$
 - Encourages smooth function: small perturbation in weights leads to small changes in output
- Injecting noise in output labels
 - Better convergence: prevents pursuit of hard probabilities

Bagging

Original dataset





deeplearningbook.org

Dropout

 h_1

 x_1

Train all sub-networks obtained by removing non-output units from base network



Goodfellow et al. (2016)

Dropout: Stochastic GD

- For each new example/mini-batch:
 - Randomly sample a binary mask μ independently, where μ_i indicates if input/hidden node *i* is included
 - Multiply output of node *i* with μ_i , and perform gradient update
- Typically, an input node is included with prob.0.8, hidden node with prob. 0.5

Dropout: Weight Scaling

• During prediction time use all units, but scale weights with probability of inclusion



• Approximates the following inference rule:

$$ilde{p}_{ ext{ensemble}}(y \mid oldsymbol{x}) = \sqrt[2^d]{\prod_{oldsymbol{\mu}} p(y \mid oldsymbol{x}, oldsymbol{\mu})}}_{ ext{Cristina Scheau (2016)}}$$

Adversarial Examples



Training on adversarial examples is mostly intended to improve security, but can sometimes provide generic regularization.

Multi-task Learning



Goodfellow et al. (2016)