

Noise Robustness

Sargur N. Srihari
srihari@buffalo.edu

Regularization Strategies

1. Parameter Norm Penalties
2. Norm Penalties as Constrained Optimization
3. Regularization and Under-constrained Problems
4. Data Set Augmentation
5. Noise Robustness
6. Semi-supervised learning
7. Multi-task learning
8. Early Stopping
9. Parameter tying and parameter sharing
10. Sparse representations
11. Bagging and other ensemble methods
12. Dropout
13. Adversarial training
14. Tangent methods

Topics in Noise Robustness

- Noise applied at inputs
- Noise applied to weights
- Injecting noise at the output targets

Noise injection is powerful

- Noise applied to inputs is a data augmentation
 - For some models addition of noise with infinitesimal variance at the input is equivalent to imposing a penalty on the norm of the weights, e.g., $\lambda \mathbf{w}^T \mathbf{w}$
- Noise applied to hidden units
 - Noise injection can be much more powerful than simply shrinking the parameters
 - Noise applied to hidden units is so important that it merits its own separate discussion
 - Dropout is the main development of this approach

Adding Noise to Weights

- This technique primarily used with RNNs
- This can be interpreted as a stochastic implementation of Bayesian inference over the weights
 - Bayesian learning considers model weights to be uncertain and representable via a probability distribution $p(\mathbf{w})$ that reflects that uncertainty,
 - Adding noise to weights is a practical, stochastic way to reflect this uncertainty

Noise at weights encourages stability

- Noise applied to weights is equivalent to traditional regularization, encouraging stability
- This can be seen in a regression setting
 - Train $\hat{y}(\mathbf{x})$ to map \mathbf{x} to a scalar using least squares between model prediction $\hat{y}(\mathbf{x})$ and true values y :

$$J = E_{p(\mathbf{x}, y)} \left[\left(\hat{y}(\mathbf{x}) - y \right)^2 \right]$$

- Training set: m labeled samples $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$
 - We perturb each input with $\epsilon_W \sim N(\epsilon; 0, \eta \mathbf{I})$
 - For small η , this is equivalent to a regularization term $\eta E_{p(\mathbf{x}, y)} \left[\left\| \nabla_W \hat{y}(\mathbf{x}) \right\|^2 \right]$
 - It encourages parameters to regions where small perturbations of weights have small influence on output

Injecting Noise at Output Targets

- Most datasets have some mistakes in y labels
 - Harmful to maximize $\log p(y|\mathbf{x})$ when y is a mistake
- To prevent it we explicitly model noise on labels
 - Ex: we assume training set label y is correct with probability $1-\epsilon$, and otherwise any of the other labels may be correct
 - This can be incorporated into the cost function
 - Ex: *Local Smoothing* regularizes a model based on a softmax with k output values by replacing the hard 0 and 1 classification targets with targets of $\epsilon/(k-1)$ and $1-\epsilon$ respectively