

Deep Learning

6.1a Regularization for Deep Learning

Dr. Konda Reddy Mopuri
kmopuri@iittp.ac.in
Dept. of CSE, IIT Tirupati

Deep Regularization

- ① Most of the regularization techniques for deep learning are based on regularizing estimators

Deep Regularization

- ① Most of the regularization techniques for deep learning are based on regularizing estimators
- ② Trade increased bias for decreased variance

Deep Regularization

- ① An overly complex model family need not include the target function

Deep Regularization

- ① An overly complex model family need not include the target function
- ② In practice we almost never have access to the true data generating process, and which is almost certainly outside the model family

Deep Regularization

- ① Most often the best-fitting model is a large model that has been appropriately regularized

Deep Regularization

- Parameter Norm penalties (l_2, l_1 , etc.)
- Dataset Augmentation
- Noise Robustness
- Semi-Supervised Learning
- Multi-Task Learning (Parameter sharing)
- Sparse Representation
- Dropout
- etc.

Parameter Norm Penalties

- ① For neural networks, typically only the weights of the affine transformations are regularized leaving the biases unregularized

Parameter Norm Penalties

- ① For neural networks, typically only the weights of the affine transformations are regularized leaving the biases unregularized
- ② Bias controls only a single variable as opposed to weight which connects two

Parameter Norm Penalties

- ① For neural networks, typically only the weights of the affine transformations are regularized leaving the biases unregularized
- ② Bias controls only a single variable as opposed to weight which connects two
- ③ Regularizing biases induces underfitting

Parameter Norm Penalties

① L_2 parameter regularization: $\tilde{\mathcal{J}} = \frac{\alpha}{2} w^T w + \mathcal{J}(w; X, y)$

Parameter Norm Penalties

- ① L_2 parameter regularization: $\tilde{\mathcal{J}} = \frac{\alpha}{2} w^T w + \mathcal{J}(w; X, y)$
- ② L_1 regularization: $\tilde{\mathcal{J}} = \alpha |w|_1 + \mathcal{J}(w; X, y)$

Parameter Norm Penalties

- ① L_2 parameter regularization: $\tilde{\mathcal{J}} = \frac{\alpha}{2}w^T w + \mathcal{J}(w; X, y)$
- ② L_1 regularization: $\tilde{\mathcal{J}} = \alpha|w|_1 + \mathcal{J}(w; X, y)$
- ③ Norm penalties induce different desired behaviors based on the exact penalty imposed

Dataset Augmentation

- ① Bestway to make ML model generalize better is to train with more data

Dataset Augmentation

- ① Bestway to make ML model generalize better is to train with more data
- ② In practice training data is limited

Dataset Augmentation

- ① Bestway to make ML model generalize better is to train with more data
- ② In practice training data is limited
- ③ Create fake data and add it to the training data, called Dataset augmentation

Dataset Augmentation

- ① Easier for classification

Dataset Augmentation

- ① Has been particularly effective for specific classification problems such as object recognition

Dataset Augmentation

- ① Has been particularly effective for specific classification problems such as object recognition
- ② Operations such as translation by few pixels, rotating slightly, adding mild noise, etc. greatly improve generalization

Dataset Augmentation

- ① Has been particularly effective for specific classification problems such as object recognition
- ② Operations such as translation by few pixels, rotating slightly, adding mild noise, etc. greatly improve generalization
- ③ Hand-designed augmentations in some domains can result in dramatic improvements
- ④ Should restrict to label preserving transformations

Multi-Task Learning

- ① Improves generalization by collecting samples arising out of multiple tasks

Multi-Task Learning

- ① Improves generalization by collecting samples arising out of multiple tasks
- ② Similar to additional data samples, multi-task samples also put more pressure on the parameters of the shared layers to be better

Multi-Task Learning

- ① Improves generalization by collecting samples arising out of multiple tasks
- ② Similar to additional data samples, multi-task samples also put more pressure on the parameters of the shared layers to be better

