

# TTIC 31230, Fundamentals of Deep Learning

David McAllester, Autumn 2020

Early Stopping meets Shrinkage

$L_1$  Regularization and Sparsity

Ensembles

## Shrinkage meets Early Stopping

Early stopping can limit  $||\Phi||$ .

But early stopping more directly limits  $||\Phi - \Phi_{\text{init}}||$ .

It seems better to take the prior on  $\Phi$  to be

$$p(\Phi) \propto \exp \left( -\frac{||\Phi - \Phi_{\text{init}}||^2}{2\sigma^2} \right)$$

giving

$$\Phi_{t+1} = \Phi_t - \eta \hat{g} - \gamma(\Phi_t - \Phi_{\text{init}})$$

## $L_1$ Regularization

$$p(\Phi) \propto e^{-\lambda ||\Phi||_1} \quad ||\Phi||_1 = \sum_i |\Phi_i|$$

$$\Phi^* = \operatorname{argmax}_{\Phi} p(\Phi) \prod_i P_{\Phi}(y_i|x_i)$$

$$\Phi^* = \operatorname{argmin}_{\Phi} \left( \sum_i -\ln P_{\Phi}(y_i|x_i) \right) + \lambda ||\Phi||_1$$

$$\Phi^* = \operatorname{argmin}_{\Phi} \hat{\mathcal{L}}(\Phi) + \frac{\lambda}{N_{\text{Train}}} ||\Phi||_1$$

## $L_1$ Regularization

$$\Phi^* = \underset{\Phi}{\operatorname{argmin}} \quad \hat{\mathcal{L}}(\Phi) + \frac{\lambda}{N_{\text{Train}}} \|\Phi\|_1$$

$$\Phi_i \leftarrow \eta \left( \hat{g}_i + \frac{\lambda}{N_{\text{Train}}} \operatorname{sign}(\Phi_i) \right)$$

$$\eta = (1 - \mu) B \eta_0$$

## Sparsity

$$\Phi_i \leftarrow \eta \left( \hat{g}_i + \frac{\lambda}{N_{\text{Train}}} \text{sign}(\Phi_i) \right)$$

For  $\Phi^*$  the gradient of the objective, and hence the average update, must be zero:

$$\Phi_i^* = 0 \quad \text{if } |g_i| < \lambda/N_{\text{Train}}$$

$$g_i = -(\lambda/N_{\text{Train}})\text{sign}(\Phi_i) \quad \text{otherwise}$$

But in practice  $\Phi_i$  will never be exactly zero.

## Ensembles

Train several models  $\text{Ens} = (\Phi_1, \dots, \Phi_K)$  from different initializations and/or under different meta parameters.

We define the ensemble model by

$$P_{\text{Ens}}(y|x) = \frac{1}{K} \sum_k P_{\Phi_k}(y|x) = E_k P_k(y|x)$$

Ensemble models almost always perform better than any single model.

## Ensembles Under Cross Entropy Loss

$$\begin{aligned}\mathcal{L}(P_{\text{Ens}}) &= E_{\langle x, y \rangle \sim \text{Pop}} - \ln P_{\text{Ens}}(y|x) \\ &= E_{\langle x, y \rangle \sim \text{Pop}} - \ln E_k P_k(y|x) \\ &\leq E_{\langle x, y \rangle \sim \text{Pop}} E_k - \ln P_k(y|x) \\ &= E_k \mathcal{L}(P_k)\end{aligned}$$

## Ensembles Under Cross Entropy Loss

It is important to note that

$$-\ln E_k P_k(y|x) \leq E_k - \ln P_k(y|x)$$

for each individual pair  $\langle x, y \rangle$ .

This may explain why in practice an ensemble model is typically better than any single component model.



**END**