

Regularization focus on minimizing testing error (generalization error)

- Ideal model is often large but regularized.

### Parameter Norm Penalties

$$\text{Form: } \tilde{J}(\theta, x, y) = J(\theta, x, y) + \boxed{\lambda \cdot \text{norm}(\theta)}$$

parameter norm penalties

- $L_2$  regularization

$$\text{Form: } \tilde{J}(w, x, y) = J(w, x, y) + \frac{1}{2} \lambda w^T \cdot w$$

$$\text{Discussion 1: } \nabla_w \tilde{J}(w, x, y) = \lambda w + \nabla_w J(w, x, y)$$

$$\text{Updating gradients } w \leftarrow w - \epsilon (\lambda w + \nabla_w J(w, x, y))$$

$$w \leftarrow \boxed{(1-\epsilon\lambda)w - \epsilon \nabla_w J(w, x, y)}$$

let  $w$  to be a smaller value.

$$\text{Discussion 2: } \hat{J}(\theta) = J(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*)$$

$$\nabla_w \hat{J}(\theta) = H(w - w^*) \xrightarrow{\text{adding normalization}}$$

$$\lambda \hat{w} + H(\hat{w} - w^*) = 0$$

$$\hat{w} = (H + \lambda I)^{-1} H w^* \xrightarrow{\substack{\text{Decomposition} \\ H}} \hat{w} = Q(1 + \lambda I)^{-1} \Lambda Q^T w^*$$

$L_2$  Regularization affects little on the larger value of  $H$  eigenvalues, which means affect the one who affects the result least.

Discussion 3: Linear Regression with  $L_2$  regularization means adding  $\lambda I$  to co-variance matrix  $XX^T$

- $L_1$  Regularization

$$\text{Form: } \tilde{J}(w, x, y) = J(w, x, y) + \lambda \cdot \|w\|_1$$

$$\nabla_w \tilde{J}(w, x, y) = \nabla_w J(w, x, y) + \lambda \cdot \text{sgn}(w)$$

$$\text{Feature selection: } \hat{J}(w, x, y) = J(w^*, x, y) + \sum_i \left[ \frac{1}{2} H_{ii} (w_i - w_i^*)^2 + \lambda |w_i| \right]$$

$$\text{Then: } w_i = \text{sign}(w_i^*) \max \left\{ |w_i^*| - \frac{\lambda}{\|x_i\|}, 0 \right\}$$

The one who is not significantly affecting the J.  
will be eliminated.

$L_1$  regularization is equal to log-prior when Bayesian + Laplace Approx.

### Norm Penalties as Constrain Optimization

Construct Langrange function for constrain optimization of Norm Penalties (Explicit constrain, instead of enforcing with penalties)

$$L(\theta, \alpha, x, y) = J(\theta, x, y) + \alpha(R(\theta) - k) \Rightarrow \text{Explicit Constrain.}$$

- Regularization can also be seen as solution to underdetermined issues.

- Dataset Augmentation: subjective
  - Application: Preprocessing (Rotate, scale)
  - General: Algorithm (Adding noise)

- Noise Robustness
  - Adding noise to input/weight: Adding additional regularization
  - Adding noise to output/target: Hard to soft.

- Semi-Supervised Learning: Learning engaged with P(X)

- Multi-task learning: Learning relative tasks using shared weights.

### Early stopping

Here we regard training steps as a hyper-parameter.

Algorithms:

- Training the model with a stopping criterion

After  $j$  steps of validation error goes up, we stop training  
and we save the optimal training steps  $j^*$

- Retrain the model with optimal  $i^*$ 
  - 1° Training from scratch with all training data, and train exactly  $i^*$  steps.
  - 2° Continue training with all data, until the validation error falls down below a threshold.

Early stopping can be seen as a  $L_2$  regularizer:

- $L_2$  regularizer result:

$$Q^\top \tilde{w} = (\lambda + dI)^{-1} \lambda Q^\top w^*$$

- Early stopping when  $w^{(0)} = 0$ :

$$w^{(T)} = w^{(T-1)} - \epsilon \nabla_w J(w^{(T-1)})$$

$$Q^\top (w^{(T)} - w^*) = (I - \epsilon \lambda) Q^\top (w^{(T-1)} - w^*)$$

$$\downarrow w^{(0)} = 0$$

$$Q^\top w^{(T)} = [I - (I - \epsilon \lambda)^T] Q^\top w^*$$

Build up  
connections

$\Rightarrow$  If we Taylor expand  $J$ , and set  $w^{(0)} = 0$ , we see in some specific value of  $d$  and  $T$ , we can equal a  $\tilde{w}$  got by early stopping equal to a  $\tilde{w}$  got by  $L_2$  regularization.

● Parameter tying and sharing: CNN a good example.

● Sparse Representations: Forces some part of  $h$  to be zero

● Bagging and other ensemble methods:

$$E[(\frac{1}{k} \sum \epsilon_i)^2] = \frac{1}{k} V + \frac{k-1}{k} C$$

Model independent Covariance between model

Dropout

Dropout can be seen as bagging with partly shared weights

Algorithms:

Training: No matter forward or backward, we randomly abandon units with a probability  $p_{\text{abn}}$ , by multiplying units with binary variable  $\mu \in \{0, 1\}$ .

Testing / Inference:

$$w \leftarrow w \cdot p_{\text{chi}} \quad \text{or} \quad h \leftarrow \frac{1}{p_{\text{chi}}} \cdot h \quad \xrightarrow{\text{when training}} Z_{\text{phi}}(D) \leftarrow$$

This is for alignment with the expectation of input in training.

Notice:

- Non-linear model cannot prove weight scaling inference but linear model can.
- Dropout needs more complicated model and training steps.
- May be not effective when data is less or with semi-supervised learning.
- Linear Regression: Dropout = weight decay
- Dropout in stochasticity has no guarantee to be optimal.
- $\mu$  sometimes can be any form (continuous  $N \sim (1, 1)$ )
- Dropout used in Hidden Units is more common  
 $\Rightarrow$  more stable/general noise: Deconstructing features.

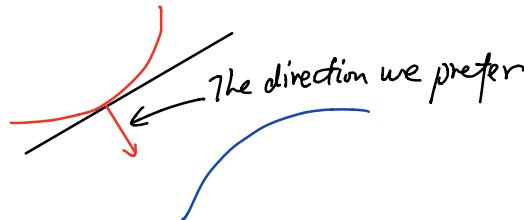
● Adversarial Training: Local Invariant

- Semi-supervised adversarial training

● Tangent Distance, Tangent Prop and Manifold Tangent Classifier

$$\pi_2(f) = \sum_i ((\nabla_x f(x))^T v_i)^2 \quad (i \text{ for each class})$$

we prefer the update direction out of the manifold of a class.



$\Rightarrow$  { Data Augmentation: Explicit manifold  
ReLU: Linear units cannot always shrink.  
Auto-encoder: A way of finding tangent manifold

- $\Rightarrow$  ① Using auto-encoder to find tangent manifold  
② Using tangent regularization to train a model (tangent prop)