# Parameter Norm Penalties

Sargur N. Srihari

srihari@cedar.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
6. Parameter tying and parameter sharing
7. Sparse representations
8. Bagging and other ensemble methods
9. Dropout
10. Adversarial training
11. Tangent methods

# Topics in Parameter Norm Penalties

1. Overview (limiting model capacity)
2. $L^2$ parameter regularization
3. $L^1$ regularization

# Limiting Model Capacity

- ## Regularization has been used for decades prior to advent of deep learning

- ## Linear- and logistic-regression allow simple, straightforward and effective regularization strategies

  - Adding a parameter norm penalty $\Omega(\theta)$ to the objective function $J$ :

  $$\tilde{J}(\boldsymbol{\theta}; X, y) = J(\boldsymbol{\theta}; X, y) + \alpha\Omega(\boldsymbol{\theta})$$

    - where $\alpha\varepsilon[0,\theta)$ is a hyperparameter that weight the relative contribution of the norm penalty term $\Omega$

      - Setting $\alpha$ to $0$ results in no regularization. Larger values correspond to more regularization

4

# Norm Penalty

- When our training algorithm minimizes the regularized objective function

$$\tilde{J}(\boldsymbol{\theta}; X, y) = J(\boldsymbol{\theta}; X, y) + \alpha\Omega(\boldsymbol{\theta})$$

  - it will decrease both the the original objective $J$ on the training data and some measure of the size of the parameters $\theta$

- Different choices of the parameter norm $\Omega$ can result in different solutions preferred
  - We discuss effects of various norms

# No penalty for biases

- Norm penalty $\Omega$ penalizes only weights at each layer and leaves biases unregularized
  - Biases require less data to fit than weights
  - Each weight specifies how variables interact
    - Fitting weights requires observing both variables in a variety of conditions

- Each bias controls only a single variable
  - We do not induce too much variance by leaving biases unregularized

- $w$ indicates all weights affected by norm penalty

- $\theta$ denotes both $w$ and biases

6

# Different or Same $\alpha$s for layers?

- Sometimes it is desirable to use a separate penalty with a different $\alpha$ for each layer

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha \Omega(\theta)$$

- Because it can be expensive to search for the correct value of multiple hyperparameters, it is still reasonable to use same weight decay at all layers to reduce search space

# $L^2$ parameter Regularization

- Simplest and most common kind
- Called *Weight decay*
- Drives weights closer to the origin
  - by adding a regularization term to the objectve function

$$\Omega(\theta) = \frac{1}{2} \| w \|_2^2$$

- In other communities also known as *ridge regression* or *Tikhonov regularization*

# Gradient of Regularized Objective

- Objective function (with no bias parameter)

$$\tilde{J}(w; X, y) = \frac{\alpha}{2} w^T w + J(w; X, y)$$

- Corresponding parameter gradient

$$\nabla_w \tilde{J}(w; X, y) = \alpha w + \nabla_w J(w; X, y)$$

- To perform single gradient step, perform update:

$$w \leftarrow w - \varepsilon \left( \alpha w + \nabla_w J(w; X, y) \right)$$

- Written another way, the update is

$$w \leftarrow (1 - \varepsilon \alpha) w - \varepsilon \nabla_w J(w; X, y)$$

  - We have modified learning rule to shrink $w$ by constant factor $1 - \varepsilon \alpha$ at each step

9

# To study effect on entire training

- Make quadratic approximation to the objective function in the neighborhood of minimal unregularized cost $\boldsymbol{w}^* = \arg\min_{\boldsymbol{w}} J(\boldsymbol{w})$

- The approximation is given by

$$J(\boldsymbol{w}^*) + \tfrac{1}{2}(\boldsymbol{w}\text{-}\boldsymbol{w}^*)^T H(\boldsymbol{w}\text{-}\boldsymbol{w}^*)$$

- Where $H$ is the Hessian matrix of $J$ wrt $\boldsymbol{w}$ evaluated at $\boldsymbol{w}^*$

# Illustration of $L^2$ regularization
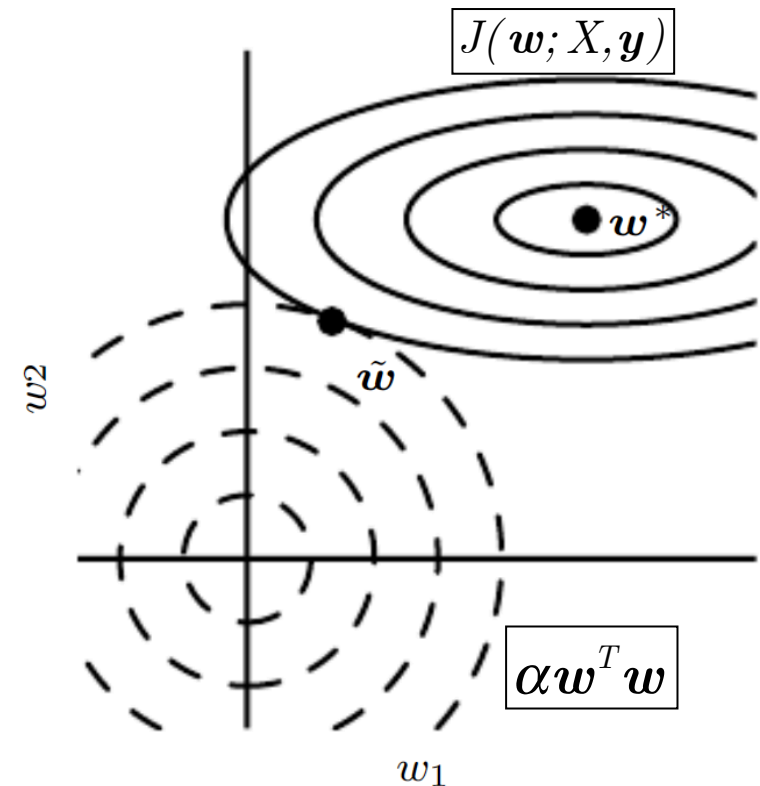
Effect on value of optimal $\boldsymbol{w}$

Solid ellipses:

  contours of equal value of unregularized objective $J$

Dotted circles:

  contours of equal value of $L^2$ regularizer

At point $\boldsymbol{w}$ competing objectives reach equilibrium

$\boxed{J(\boldsymbol{w}; X, \boldsymbol{y})}$

$w_2$

$\tilde{w}$

$\boxed{\alpha \boldsymbol{w}^T \boldsymbol{w}}$

$w_1$

$\boldsymbol{w}^*$

| | |
|---|---|
| Along $w_1$, eigen value of Hessian of $J$ is small. $J$ does not increase much when moving horizontally away from $\boldsymbol{w}*$. Because $J$ does not have a strong preference along this direction, the regularizer has a strong effect on this axis. The regularizer pulls $w_1$ close to $0$. | Along $w_2$, $J$ is very sensitive to movements away from $\boldsymbol{w}*$. The corresponding eigenvalue is large, indicating high curvature. As a result, weight decay affects the position of $w_2$ relatively little |

11

# $L^1$ Regularization

- While $L^2$ weight decay is the most common form of weight decay there are other ways to penalize the size of model parameters

- $L^1$ regularization is defined as

$$\Omega(\boldsymbol{\theta}) = \left\|\boldsymbol{w}\right\|_1 = \sum_i \left|w_i\right|_1$$

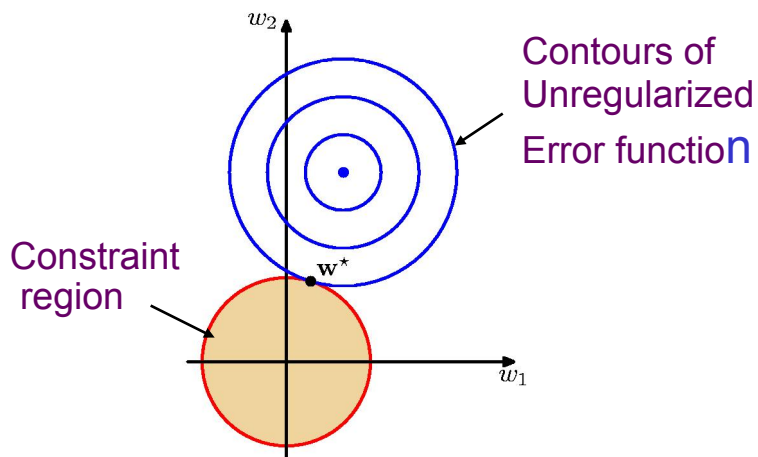  – which is the sum of the absolute values of the individual parameters

# Sparsity and Feature Selection

• The sparsity property induced by $L^1$ regularization has been used extensively as a feature selection mechanism

– Feature selection simplifies an ML problem by choosing subset of available features

• LASSO (Least Absolute Shrinkage and Selection Operator) integrates an $L^1$ penalty with a linear model and least squares cost function

– The $L^1$ penalty causes a subset of the weights to become zero, suggesting that those features can be discarded

# Sparsity with Lasso constraint

- With $q=1$ and $\lambda$ is sufficiently large, some of the coefficients $w_j$ are driven to zero
- Leads to a sparse model
  - where corresponding basis functions play no role
- Origin of sparsity is illustrated here:

Quadratic solution where $w_1{}^*$ and $w_0{}^*$ are nonzero

Minimization with Lasso Regularizer A sparse solution with $w_1{}^*=0$



Contours of Unregularized Error function

Constraint region

14