

# CSE4088 Introduction to Machine Learning

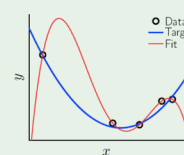
## Regularization

Slides are adopted from lecture notes of Yaser Abu-Mostafa

## Review of last lecture

### • Overfitting

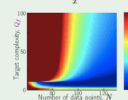
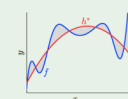
Fitting the data more than is warranted



VC allows it; doesn't predict it

Fitting the noise, stochastic/deterministic

### • Deterministic noise



## Outline

- Regularization – informal
- Regularization – formal
- Weight decay
- Choosing a regularizer

## Two approaches to regularization

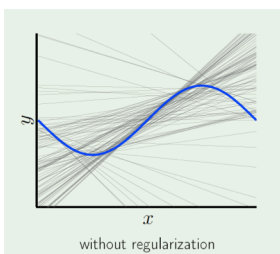
### Mathematical:

Ill-posed problems in function approximation

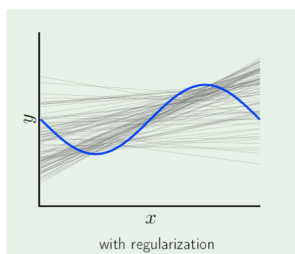
### Heuristic:

Handicapping the minimization of  $E_{in}$

## A familiar example

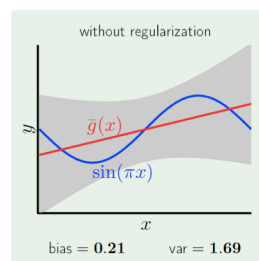


without regularization

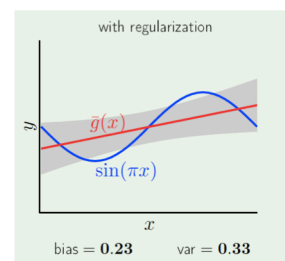


with regularization

## And the winner is ...



without regularization  
bias = 0.21 var = 1.69



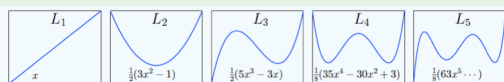
with regularization  
bias = 0.23 var = 0.33

## The polynomial model

$\mathcal{H}_Q$ : polynomials of order  $Q$  linear regression in  $\mathcal{Z}$  space

$$\mathbf{z} = \begin{bmatrix} 1 \\ L_1(x) \\ \vdots \\ L_Q(x) \end{bmatrix} \quad \mathcal{H}_Q = \left\{ \sum_{q=0}^Q w_q L_q(x) \right\}$$

Legendre polynomials:



## Unconstrained solution

Given  $(x_1, y_1), \dots, (x_N, y_N) \longrightarrow (z_1, y_1), \dots, (z_N, y_N)$

$$\text{Minimize } E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{z}_n - y_n)^2$$

$$\text{Minimize } \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^T (\mathbf{Z}\mathbf{w} - \mathbf{y})$$

$$\mathbf{w}_{\text{lin}} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y}$$

## Constraining the weights

Hard constraint:  $\mathcal{H}_2$  is constrained version of  $\mathcal{H}_{10}$  with  $w_q = 0$  for  $q > 2$

Softer version:  $\sum_{q=0}^Q w_q^2 \leq C$  "soft-order" constraint

$$\text{Minimize } \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^T (\mathbf{Z}\mathbf{w} - \mathbf{y})$$

$$\text{subject to: } \mathbf{w}^T \mathbf{w} \leq C$$

Solution:  $\mathbf{w}_{\text{reg}}$  instead of  $\mathbf{w}_{\text{lin}}$

Solving for  $\mathbf{w}_{\text{reg}}$

$$\text{Minimize } E_{\text{in}}(\mathbf{w}) = \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^T (\mathbf{Z}\mathbf{w} - \mathbf{y})$$

$$\text{subject to: } \mathbf{w}^T \mathbf{w} \leq C$$

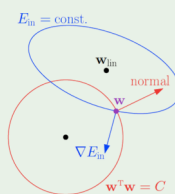
$$\nabla E_{\text{in}}(\mathbf{w}_{\text{reg}}) \propto -\mathbf{w}_{\text{reg}}$$

$$= -2 \frac{\lambda}{N} \mathbf{w}_{\text{reg}}$$

$$\nabla E_{\text{in}}(\mathbf{w}_{\text{reg}}) + 2 \frac{\lambda}{N} \mathbf{w}_{\text{reg}} = \mathbf{0}$$

$$\text{Minimize } E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$$

$$C \uparrow \lambda \downarrow$$



## Augmented error

$$\text{Minimizing } E_{\text{aug}}(\mathbf{w}) = E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$$

$$= \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^T (\mathbf{Z}\mathbf{w} - \mathbf{y}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w} \quad \text{unconditionally}$$

— solves —

$$\text{Minimizing } E_{\text{in}}(\mathbf{w}) = \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^T (\mathbf{Z}\mathbf{w} - \mathbf{y})$$

$$\text{subject to: } \mathbf{w}^T \mathbf{w} \leq C$$

← VC formulation

Since some hypothesis are not allowed, we expect the VC dimension to decrease and hence generalization behavior will improve.

## The solution

$$\text{Minimize } E_{\text{aug}}(\mathbf{w}) = E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$$

$$= \frac{1}{N} \left( (\mathbf{Z}\mathbf{w} - \mathbf{y})^T (\mathbf{Z}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w} \right)$$

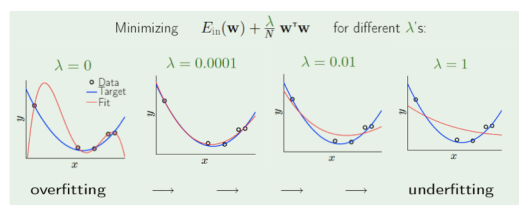
$$\nabla E_{\text{aug}}(\mathbf{w}) = \mathbf{0} \implies \mathbf{Z}^T (\mathbf{Z}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w} = \mathbf{0}$$

$$\mathbf{w}_{\text{reg}} = (\mathbf{Z}^T \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^T \mathbf{y} \quad (\text{with regularization})$$

$$\text{as opposed to } \mathbf{w}_{\text{lin}} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y} \quad (\text{without regularization})$$

If  $\lambda$  is large, then we will have  $\frac{1}{\lambda}$  on the right hand side and  $\mathbf{w}$  will be pushed towards zero  $\rightarrow$  strong regularization.

## The result



## Weight 'decay'

Minimizing  $E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$  is called weight decay. Why?

Gradient descent:

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \nabla E_{\text{in}}(\mathbf{w}(t)) - 2\eta \frac{\lambda}{N} \mathbf{w}(t)$$

$$= \mathbf{w}(t) \left(1 - 2\eta \frac{\lambda}{N}\right) - \eta \nabla E_{\text{in}}(\mathbf{w}(t))$$

Applies in neural networks:

$$\mathbf{w}^T \mathbf{w} = \sum_{l=1}^L \sum_{i=0}^{d^{(l-1)}} \sum_{j=1}^{d^{(l)}} \left(w_{ij}^{(l)}\right)^2$$

## Variations of weight decay

Emphasis of certain weights:

$$\sum_{q=0}^Q \gamma_q w_q^2$$

Examples:

- $\gamma_q = 2^q \Rightarrow$  low-order fit
- $\gamma_q = 2^{-q} \Rightarrow$  high-order fit

Neural networks: different layers get different  $\gamma$ 's

Tikhonov regularizer:  $\mathbf{w}^T \mathbf{\Gamma} \mathbf{w}$

## Even weight growth!

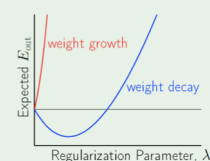
We 'constrain' the weights to be large - bad!

Practical rule:

stochastic noise is 'high-frequency'

deterministic noise is also non-smooth

$\Rightarrow$  constrain learning towards smoother hypotheses



## General form of augmented error

Calling the regularizer  $\Omega = \Omega(h)$ , we minimize

$$E_{\text{aug}}(h) = E_{\text{in}}(h) + \frac{\lambda}{N} \Omega(h)$$

Rings a bell?  $\downarrow \downarrow$

$$E_{\text{out}}(h) \leq E_{\text{in}}(h) + \Omega(\mathcal{H})$$

$E_{\text{aug}}$  is better than  $E_{\text{in}}$  as a proxy for  $E_{\text{out}}$

## Outline

- Regularization – informal
- Regularization – formal
- Weight decay
- Choosing a regularizer

The perfect regularizer  $\Omega$ 

Constraint in the 'direction' of the target function (going in circles ☺)

Guiding principle:

Direction of **smoother** or "simpler"

Chose a bad  $\Omega$ ?

We still have  $\lambda$ !

## Neural-network regularizers

Weight decay: From linear to logical

Weight elimination:

Fewer weights  $\implies$  smaller VC dimension

Soft weight elimination:

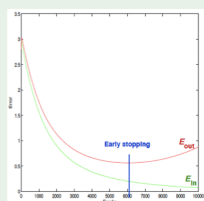
$$\Omega(\mathbf{w}) = \sum_{i,j,l} \frac{(w_{ij}^{(l)})^2}{\beta^2 + (w_{ij}^{(l)})^2}$$



## Early stopping as a regularizer

Regularization through the optimizer!

When to stop? **validation**

The optimal  $\lambda$ 