

# Advanced Machine Learning, 13 Aug 2019

Sourish Das, Madhavan Mukund

## Deep Learning

- └ Acyclic NN (feed forward)
  - └ Loss functions, Optimization Strategies, Regularization
- └ Convolutional NN
- └ Recurrent NN, LSTM

Reinforcement- learning

↳ Alpha Go

Probabilistic Graphical Models

Assignments ~ 35%

Mid Sem Exam ~ 25%

Final Exam ~ 40%

# Motivation

ML - Builds a model of a "concept"  
from training examples

Concept?

Data Space  $D$   $\langle x_1, x_2, \dots, x_n \rangle$

Concept  $C \subseteq D$

Build a model - construct a description  
of  $C$  from training data

## Regression

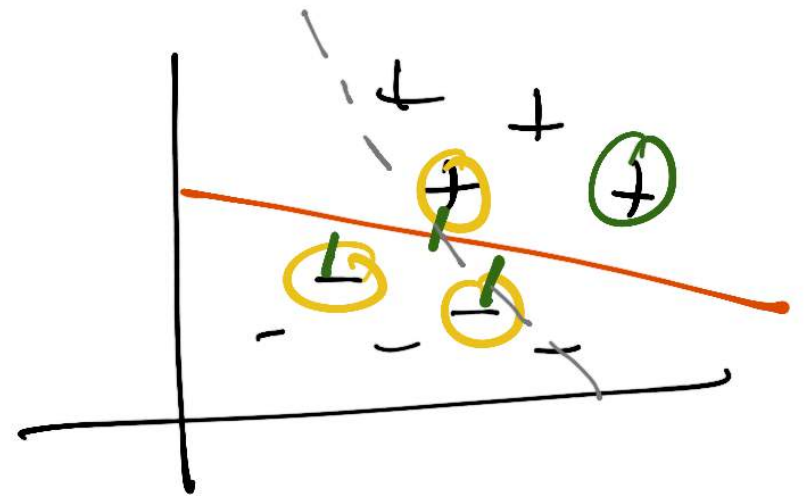
$$y = f(x_1, \dots, x_n) \\ = \sum w_i x_i + b$$

Optimize  $w_i$ 's

Cost/Loss associated  
with  $w_i$ 's

What distinguishes ML from optimization?

## SVM



Optimize min  
distance to  
boundary point

Challenge is generalization

How well the model performs on  
unseen data

Prime face - impossible

Unseen data is "similar" to training data

Probabilistic setting

- Distribution over  $D$  -  $p$

- Training & unknown samples are i.i.d w.r.t  $p$



$D$  - all data points

$C \subseteq D$  concept

$p$  - distribution over  $D$

Space of possible models

$H$  - hypothesis space

e.g. regression or SVM

- all possible hyperplanes

# PAC Learning (Leslie Valiant)

Probabilistically Approximately Correct

Fix  $p, D, \epsilon, H$

Connect error on training data

(optimization while constructing  $h \in H$ )

to potential error on unknown (general)  
data

Larger training data  $\Rightarrow$  Better  $h$

Given  $p, D, C, \mathcal{H}$

fix parameters  $\delta, \epsilon$

Can compute an  $n$  (as a function of  $\delta, \epsilon$ )

s.t. that any  $h$  built from a

training sample of size  $n$

with probability  $1 - \delta$

if  $h$  is  $\epsilon$ -away from training data,

$h$  is also  $\epsilon$ -away from unseen data



Expressiveness of  $H$  - Capacity of model

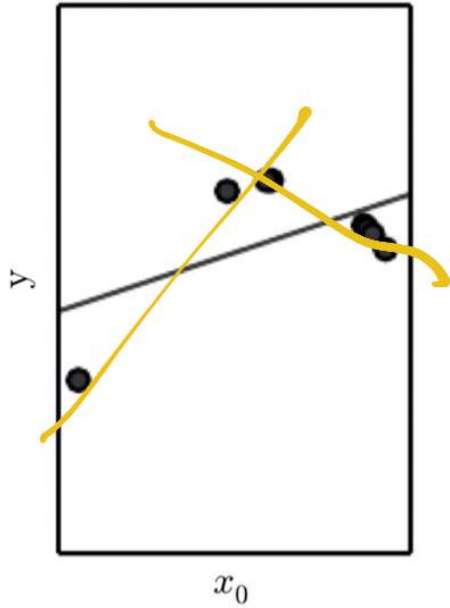
Statistical ML  $\rightarrow$  Vapnik-Chervonenkis (VC)  
Dimension

Regression

$$w_1x + b \rightarrow w_1x + w_2x^2 + b$$

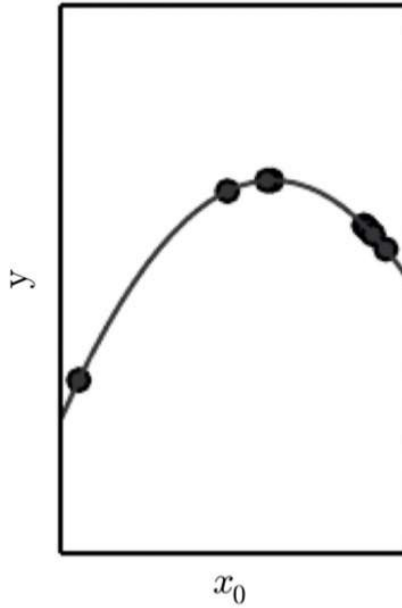
$$\dots w_1x \dots + w_kx^k + b$$

Underfitting

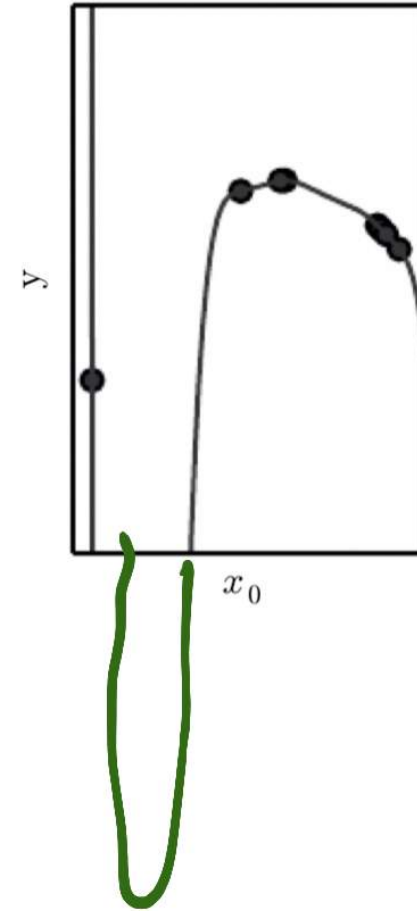


Low  
Capacity

Appropriate capacity



Overfitting



Match expressiveness of  $\mathcal{H}$  to  
complexity of  $\mathcal{C}$

Occam's Razor - Simplest model that fits  
is the best

Penalize complex models in  $\mathcal{H}$

Polynomial:  $\sum_{i=0}^k w_i x^i$   $W = \langle w_0, \dots, w_k \rangle$

Add a cost proportional to  $w_i$ 's

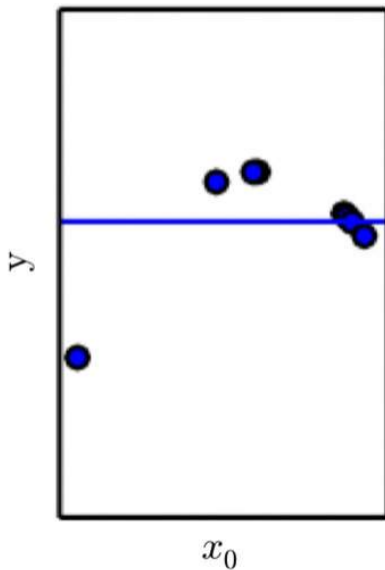
Typical  $W^T W$

Optimize: Original loss +  $\lambda \cdot$  Complexity

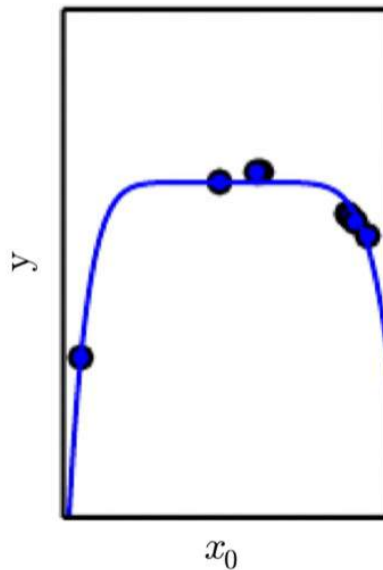
# Regularization

Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.

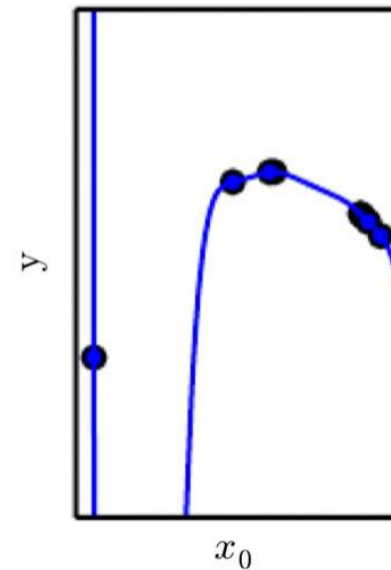
Underfitting  
(Excessive  $\lambda$ )

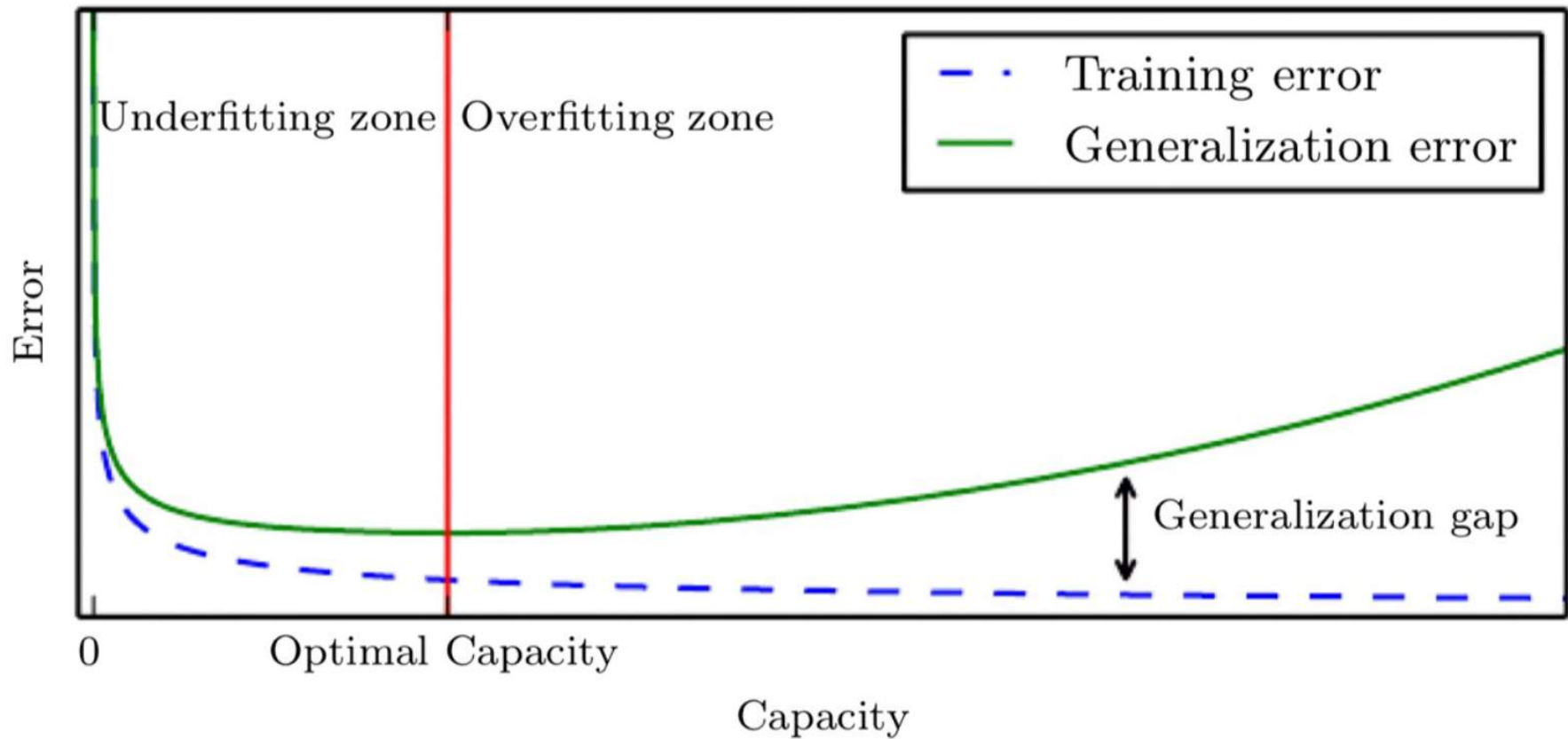


Appropriate weight decay  
(Medium  $\lambda$ )



Overfitting  
( $\lambda \rightarrow 0$ )





Regularization tries to push the hypothesis towards the red line from the right



# Hyperparameters

ML estimates parameters of a given  $h$

Degree of polynomial?

Don't want ML to "discover" this

Freeze a priori - Hyperparameter

Another Example - regularization coefficient  
 $\lambda$

# Bias & Variance

Bias - Expressiveness / Ability to describe  $C$

Variance - Sensitivity to training data

## Statistically

Across all training data =  $E(C)$

Real  $C$

"Difference"  $E(C) - C = \text{bias}$

# Variance

Statistical variance between  $E(c)$   
and a given model  $h$

---

## Neural Networks

Networks of linear separators

SVM - non linear data via geometric  
transforms - "Kernel trick"

Network of perceptrons has (arbitrarily)  
high capacity

Applying regularization &  
optimization (variations on  
gradient descent)