



#### The Alan Turing Institute

#### Introduction to Machine Learning

Prof. Sotirios A Tsaftaris WP3.1 Lead in UDRC 3 Chair in Machine Learning and Computer Vision Canon Medical / Royal Academy of Engineering Research Chair in Healthcare Al

Turing Fellow https://vios.science

S.Tsaftaris@ed.ac.uk



Royal Academy of Engineering

Some slides adopted from: Mike Davies, A. Katsaggelos, S. Theodoridis, A. Ng

# Disclaimers

- The following slides contain material from several sources.
- They are to be used by participants of the summer school and cannot be distributed without permission from the lecturer.
- Copyright of figures remains on the copyright holders which may not be solely of the author. Attribution has been given when possible but has not been exhaustive.
- Material shared maybe more than what we will be covered at delivery. Said material is made available for self-learning.
- Material subject to change.

# **Common questions**

Q: Are the slides going to be shared?

A: Yes, a preview version is already in the dropbox link https://www.dropbox.com/s/am2tspic0t0fy2x/Summer\_Sch ool\_Slides\_tsaftaris\_2021.pptx?dl=0

Updated version is available at the link of the chat!

- Q: Will the session be recorded?
- A: Yes and will be shared afterwards.

Q: Are you going to show us practical examples?A: Yes via web-based demos. More hands on follows later in the day.

# Our expertise

# Computer vision & Image Analysis:

- Extract information from images
- Applications: medicine, plants
- Machine learning & pattern recognition
  - Shallow or deep representation learning (the process of making sense of data)
  - Learning without lots of data
  - Combining information sources Z = f
- Distributed learning









# We do Al...



Canon Medical ——

Intern

# An example of ML in UDRC

Given data can
 I learn their distribution?



Real images (CIFAR-10)

• Learn a distribution that can generate that data?



Generated images

# Some common applications

- Email spam filtering
- Netflix/Amazon recommendations
- Google suggested queries
- The Google index itself
- Predicting stock prices
- Classifying threats in images
- etc

# Extreme (...) applications

- MIT flight
- <u>http://www.youtube.com/watch?v=aiNX-vpDhMo</u>
- Robot in the dessert
- http://www.youtube.com/watch?v=OIOtOmyySQo
- Google car
- <u>http://www.youtube.com/watch?v=cdgQpa1pUUE</u>

On purpose I leave these "old" extremes to show how quickly ML has advanced in less than 10 years.

#### Autonomous Drones



https://www.youtube.com/watch?v=2DsJLigMrS0

#### Use AI to design new chips...



https://www.nature.com/articles/s41586-021-03544-w

#### Autonomous driving



Source: Andrej Karpathy Director of AI and Autopilot Vision at Tesla, Multi-Task Learning in the Wilderness ICML 2019



#### **Protein folding**

https://www.nature.com/articles/s41586-019-1923-7.epdf?author\_access\_token=Z\_KaZKDqtKzbE7Wd5Htwl9RgN0jAjWel9jnR3Zo Tv0MCcgAwHMgRx9mvLjNQdB2TIQQaa7I420UCtGo8vYQ39gg8IFWR9mAZtvs N\_1PrccXflbc6e-tGSgazNL\_XdtQzn1PHfy21qdcxV7Pw-k3htw%3D%3D

# What is machine learning?

 Arthur Samuel [1959] (informal definition) Gives computers ability to learn without being explicitly programmed.

→ He built the very first checker's program

- Tom Mitchell [98] (more formal): A well-posed learning problem is defined as follows:
  - A computer program is set to learn from an experience E with respect to some task T and some performance measure P if its performance on T as measured by P improves with experience E.

## Machine learning the 10 year challenge

$$2009 \qquad 2019 \\ Y = BX + E \qquad Y = BX + E \\ STATISTICS \qquad MACHINE LEARNING \\ X 10 YEARS CHALLENGE$$

#### **Text Books**

• Useful texts ...







Most can be found online or libraries: e.g.
 <a href="https://github.com/jermwatt/machine\_learning\_refined">https://github.com/jermwatt/machine\_learning\_refined</a>
 <a href="https://www.deeplearningbook.org/">https://www.deeplearningbook.org/</a>





# What \*we\* do with ML...

- We build algorithms to analyze imaging data (2D, 3D, 2D+t, 3D+t)
- From a variety of domains
- Use machine learning throughout
- Some examples...

# Restoring faults in images

- Recover gaps from images of plant roots
- A similar problem is present in medical imaging as well
   →Retina fundus





Chen et al. "Adversarial Large-scale Root Gap Inpainting ", CVPRW - CVPPP 2019

### Learning to age



Xia, Chartsias, Wang, & Tsaftaris (2019) Learning to synthesise the ageing brain without longitudinal data. arXiv 1912.02620 Xia, Chartsias & Tsaftaris (2019) Consistent Brain Ageing Synthesis. MICCAI

# Doing more with less

- Build systems that learn from few data, few annotations
  - Understand the world
  - Build intuition about the world



Chartsias, Joyce, Papanastasiou, Semple, Williams, Newby, Dharmakumar & Tsaftaris (2018 and 2019) Disentangled Representation Learning in Cardiac Image Analysis. Medical Image Analysis 2019 and MICCAI 2018 Code: https://github.com/agis85/anatomy\_modality\_decomposition

# Generate high quality synthetic data





# Create **high quality** synthetic images combining data



Thermos, Liu, O'Neil, & Tsaftaris (2021) Controllable cardiac synthesis via disentangled anatomy arithmetic. To appear MICCAI 2021

# Test-time adaptation

- Encoders, decoders, discriminators capture priors
- At test time we throw them away.
- Can we reuse them and adapt networks during testing?

G. Valvano, A. Leo and S. A. Tsaftaris, Stop Throwing Away Discriminators! Re-using Adversaries for Test-Time Training, Work in progress

Dataset	Input	Pred Before TTT	After TTT	True
ACDC <sub>1.5T→3T</sub>		<b>Ç</b> •	•	•
M & Ms		<b>0</b>	•	•
ACDC		<b>`</b>	•	<ul> <li></li> <li></li></ul>
LVSC		, <b>o</b>	<b>o</b>	0
CHAOS		<ul> <li>.</li> <li>.</li> </ul>	<ul> <li>.</li> <li>.</li> </ul>	<ul> <li>.</li> <li>.</li> </ul>

# Anomaly detection





- 1. N. Dionelis, M. Yaghoobi, S. A. Tsaftaris, "Boundary of Distribution Support Generator (BDSG): Sample Generation on the Boundary," *IEEE ICIP* 2020.
- 2. N. Dionelis, M. Yaghoobi, S. A. Tsaftaris, "Tail of Distribution GAN: GAN-Based Boundary of Distribution Formation," *SSPD* 2020.
- 3. N. Dionelis, M. Yaghoobi, S. A. Tsaftaris, "Few-Shot Adaptive Detection of Objects of Concern Using Generative Models with Negative Retraining (REFGAN model)", Work in Progress, 2021.

- The major directions of learning are:
  - Supervised: Patterns whose class is known a-priori are used for training.
  - Unsupervised: The number of classes/groups is (in general) unknown and no training patterns are available.
  - Semisupervised: A mixed type of patterns is available. For some of them, their corresponding class is known and for the rest is not.

#### Supervised example



#### A classification example



Goal find a "line" to separate the two classes

#### Hmm this looks harder



Looks like separating the data needs a more "complex" line

#### Hmm this looks harder



#### **Unsupervised clustering**



#### **Unsupervised clustering**



## Semi-supervised learning



## Semi-supervised learning



# Topics that we will try to cover

- Supervised methods
  - Linear regression
  - Logistic regression
  - Support vector machines (briefly)
  - Perceptron Classifier
- Unsupervised (and dimensionality reduction)
  - PCA, kernel PCA
- Learning theory (simple view)
- Introduction; simple feed forward neural network architecture; how to train a neural network; backpropagation (briefly); introduction to convolutional neural networks.

# PART 1: INTRO TO ML

## Supervised learning



# A typical formulation

- Input or features:
- Output or target:
- Training example:
- Training set==
   List of m examples:

$$\{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$$

 $\mathcal{X} \mathcal{Y}$ 

• Space:



Living area (feet <sup>2</sup> )	# bedrooms	Price $(1000$ \$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
	•	÷
#### Linear fit



#### **Objective function:**

How do we find the solution to this problem?

We are given some data. We are given a desired form of the line, but we want to find the *best* line.

We need an objective function = training cost =measure of performance that we can **optimize**.

#### Optimization

Many ML techniques need optimization, e.g.

- Minimizing error in a neural network/adaptive system
- Maximizing probability in Bayesian inference
- From simple "steepest descent" to more advanced techniques (Newton, conjugate gradient,...)



#### Solving the linear regression problem

• The LMS algorithm

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \qquad \qquad h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x,$$

- **Define a cost:**  $J(\theta) = \frac{1}{2} \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) y^{(i)})^2.$
- Optimise for the cost  $\theta_j := \theta_j \alpha \frac{\partial}{\partial \theta_j} J(\theta).$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_\theta(x) - y)^2$$
$$= 2 \cdot \frac{1}{2} (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_\theta(x) - y)$$
$$= (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^n \theta_i x_i - y \right)$$
$$= (h_\theta(x) - y) x_j$$

$$\theta_j := \theta_j + \alpha \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}$$

#### For a dataset of size 1

#### Batch vs stochastic

• Batch

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)} \qquad \text{(for every } j\text{)}$$

• Stochastic

Loop {  
for i=1 to m, {  
$$\theta_j := \theta_j + \alpha \left( y^{(i)} - h_{\theta}(x^{(i)}) \right) x_j^{(i)}$$
 (for every  $j$ )  
}

#### Higher model complexity



#### Model choice:

How do we tell what is the right choice? **Theoretical background:** the bias variance dilemma The practical solution: cross validation

Coming later...

## Logistic Regression

- Suppose now that we don't have a regression problem but classification
  - Decide whether a house will sell (1) or not (0)

Living area (feet <sup>2</sup> )	#bedrooms	Price $(1000$ \$s)
2104	3	400
1600	3	330

- y, our output variable, is now a binary number
- We can encode this as:
  - 0 or 1
  - -1 , 1
  - Or anything else according to convenience (which will be obvious)

#### Visualizing a classification problem



#### Logistic Regression

• Can we solve a binary classification problem 0 or 1 as linear regression?

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}},$$





# Why?

• It has some nice property

$$g'(z) = \frac{d}{dz} \frac{1}{1 + e^{-z}}$$
  
=  $\frac{1}{(1 + e^{-z})^2} (e^{-z})$   
=  $\frac{1}{(1 + e^{-z})} \cdot \left(1 - \frac{1}{(1 + e^{-z})}\right)$   
=  $g(z)(1 - g(z)).$ 

## How to solve the problem (find the theta's)?

- A maximum likelihood view
- $P(y=1 \mid x; \theta) = h_{\theta}(x)$ • Assume  $P(y=0 \mid x; \theta) = 1 - h_{\theta}(x)$

$$p(y \mid x; \theta) = (h_{\theta}(x))^{y} (1 - h_{\theta}(x))^{1-y}$$

• Or...

 Assuming the training examples were generated independently

$$L(\theta) = p(\vec{y} \mid X; \theta) \qquad \begin{array}{l} \text{Objective} \\ \text{function} \\ = \prod_{i=1}^{m} p(y^{(i)} \mid x^{(i)}; \theta) \\ = \prod_{i=1}^{m} \left( h_{\theta}(x^{(i)}) \right)^{y^{(i)}} \left( 1 - h_{\theta}(x^{(i)}) \right)^{1-y^{(i)}} \end{array}$$

Objective

$$\ell(\theta) = \log L(\theta) \qquad {}^{i=1} \\ = \sum_{i=1}^{m} y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log(1 - h(x^{(i)}))$$

#### Contd...

• Lets use gradient ascent  $\theta := \theta + \alpha \nabla_{\theta} \ell(\theta)$ 

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \ell(\theta) &= \left( y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ &= \left( y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x) (1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\ &= \left( y (1 - g(\theta^T x)) - (1 - y) g(\theta^T x) \right) x_j \\ &= \left( y - h_{\theta}(x) \right) x_j \end{aligned}$$

$$\theta_j := \theta_j + \alpha \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}$$

# PART 2: INTRO TO PERCEPTRONS

#### What if we take an extreme sigmoidal



$$g(z) = \begin{cases} 1 & \text{if } z \ge 0\\ 0 & \text{if } z < 0 \end{cases}$$

#### The Perceptron: A Simple Learning Neuron

Rosenblatt (1958)



Inputs may be from  $\{-1, +1\}$  or  $\{0, +1\}$ 

We want

output y

## Perceptron Learning Algorithm

One example of a learning algorithm (presents samples one at a time)

For all input vectors in training set:

- 1) Present input vector x
- 2) Calculate y=1 if  $w^T x \ge 0$ , y=0 if  $w^T x < 0$

3) Compare y with target output t

- a) If *t*=1 but *y*=0, set new w = old w + ηx
- b) If *t*=0 but *y*=1, set new w = old w  $\Box$   $\eta$ x
- c) Otherwise (If y=t), do nothing

[punish] [punish] [reward]

Repeat until correct for all input vectors.

Factor  $\eta$  is called the *learning rate* 

#### **Decision Boundary Example**





## Simple Example

"If summer and not raining, play tennis"



#### Simple Example (cont)



#### **Perceptron Limitations**

- Problem must be *linearly separable*
- Classic non-linearly separable problem: XOR problem



- Minsky & Pappert (1969) conjectured this limitation would *not* be overcome.
- But it was...

# From linear discriminants to more complex decision surfaces

#### **Linear Discriminant Functions**

Suppose we wished to decide whether some data

$$\mathbf{X} = (x_1, x_2, \dots, x_d)$$

belonged to one of two categories

One way to do this is to construct a *Discriminant function*. Let  $g(\mathbf{x})$  define the categories as:

$$\omega(\mathbf{x}) = \begin{cases} \omega_1, g(\mathbf{x}) > 0\\ \omega_2, g(\mathbf{x}) \le 0 \end{cases}$$

If g(x) is linear we can write:

$$g(\mathbf{X}) = \mathbf{W}^{t}\mathbf{X} + w_{0} = \sum_{i=1}^{d} w_{i}x_{i} + w_{0}$$

where  $\mathbf{w} = \{w_1, \dots, w_d\}$  are called the weights and  $w_0$  is called the bias or threshold weight.

See ch 5.2 Linear Discriminant Functions Duda Hart

#### Simple Linear Classifier



Each unit shows its effective input-output function.

#### **Decision surface**

 $g(\mathbf{x}) = 0$  defines a *decision surface* which separates points into  $\omega_1$  and  $\omega_2$ . If  $g(\mathbf{x})$  is linear, this decision surface is a *hyperplane*.

The hyperplane divides the space into two regions:

 $R_1$ : g(**x**) > 0, hence **x** is in  $\omega_1$  and

 $R_2$ : g(**x**) ≤ 0, hence **x** is in  $\omega_2$ 

Suppose  $x_1$  and  $x_2$  are both on the decision surface. Then:

$$\mathbf{w}^{t}\mathbf{X}_{1} + w_{0} = \mathbf{w}^{t}\mathbf{X}_{2} + w_{0}$$
 i.e.  $\mathbf{w}^{t}(\mathbf{X}_{1} - \mathbf{X}_{2}) = 0$ .

Therefore w is normal (orthogonal) to the hyperplane.

# States and a state of the state Hyperplane decision surface Let us write $\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$ $\mathcal{R}_I$ $\mathcal{R}_2$ 5a(x) 11 Ŋ $X_I$

where  $\mathbf{x}_{p}$  is normal projection of  $\mathbf{x} \to H$  and r is distance from H (r > 0 on + ve side, r < 0 on -ve)Since  $g(\mathbf{x}_n) = 0$ ,  $\mathbf{x}_{2} \quad g(\mathbf{x}) = \mathbf{w}^{t}\mathbf{x} + w_{0} = \mathbf{w}^{t}(\mathbf{x}_{p} + r\frac{\mathbf{w}}{\|\mathbf{w}\|}) + w_{0} =$  $= \ldots = r ||\mathbf{w}||$  $\Rightarrow r = g(\mathbf{x}) / \|\mathbf{w}\|$ 

 $g(\mathbf{x})$  measures dist from  $\mathbf{x}$  to H.

## **Generalized Linear Discriminants**

We can generalize  $g(\mathbf{x})$  by adding terms  $x_i x_j$  to give a quadratic (nonlinear) discriminant function:

$$g(\mathbf{X}) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j \qquad \left[ \right]$$

Can generalize to cubic, etc.

We can view this as a linear discriminant function in a new space.

Let  $y_i(\mathbf{x})$  define a new variable as a (nonlinear) function of  $\mathbf{x}$ .

$$g(\mathbf{X}) = \sum_{i=1}^{\hat{d}} a_i y_i(\mathbf{X})$$

Note we have absorbed the bias weight in this formulation - A process called *augmentation*.

See ch 5.3 Generalized Linear Discriminant Functions Duda Hart



# Linearly Separable Case (2 category)

Suppose we have n samples  $y_1, ..., y_n$ , each labelled either  $\omega_1$  or  $\omega_2$  and we wish to *learn* a discriminant function  $g(\mathbf{y}) = \mathbf{a}^T \mathbf{y}$ , that correctly classifies the data.

Sample  $\mathbf{y}_i$  is correctly classified if

 $\mathbf{a}^T \mathbf{y}_i > 0$  when  $\mathbf{y}_i$  is labelled  $\omega_1$ 

or

 $\mathbf{a}^T \mathbf{y}_i < 0$  when  $\mathbf{y}_i$  is labelled  $\omega_2$ 

A data set is called *linearly separable* if there exists a vector a which correctly classifies all samples. This is called a *separating vector* or *solution vector*.

[The case of  $g(\mathbf{y}) = \mathbf{a}^T \mathbf{y} + a_0$  with augmentation:  $\mathbf{y}' \equiv [\mathbf{y}, 1]^T, \mathbf{a}' \equiv [\mathbf{a}, a_0]^T$ ]

See Ch 5.4 The Two-Category Linearly-Separable Case Duda Hart

#### How to find a

- Great now we have understood what the **a** must satisfy (linear inequalities) and some properties.
- However, we still do not know how to find a
- We are given some data and their labels and we need a procedure to find a
- We need to find a criterion that when optimized we have a solution vector **a**.
  - Lets call this criterion J(a). Observe it is a scalar function of a : returns a value pending on a
  - If we make good choices of J() we can use **optimization** theory.
    [We will talk about this a lot later in the class.]
  - Learning a classifier is then <u>reduced</u> to an optimization problem
  - We will consider for now a simple approach: Gradient Descent

#### **Gradient Descent**

Simple concept: Consider I am at some point in my function  $J(a_k)$ . I need to move to a new point  $a_{k+1}$ . What is a good point?

**Gradient:** Gradient points in the direction of the greatest rate of increase of the function. (Generalization of derivative in multivariate functions)



Update :  $\mathbf{a}(k+1) = \mathbf{a}(k) - \eta(k)\nabla J(\mathbf{a}(k))$ Learning rate  $\eta(k)$ Gradient vector  $\nabla J(\cdot)$ 

Basic gradient descent algorithm : 1. Initialize  $k \leftarrow 0$ ,  $\mathbf{a}(1)$ ,  $\theta$ ,  $\eta(.)$ 2.  $k \leftarrow k+1$ 3.  $\mathbf{a} \leftarrow \mathbf{a} - \eta(k) \nabla J(\mathbf{a})$ 4. If  $|\eta(k) \nabla J(\mathbf{a})| > \theta$ , repeat from 2.

# Perceptron Revisited

 $J_p(\mathbf{a}) = \sum (-\mathbf{a}^T \mathbf{y})$ 

 $\mathbf{y} \in Y_m$ 

Could J = number of samples misclassified? No – the function is piecewise constant. Why is this bad? The cost function has 0 gradient at most points and wherever non-zero has discontinuities. We need a better cost.

The Perceptron

warning: Data are assumed normalized

> where  $Y_m$  is the set of samples misclassified by **a**. When is  $J_p(a)=0$ ?

Is J<sub>p</sub>(a) always positive?

In the 2-category problem it is possible to <u>normalize</u> the data: negate <sup>*a*<sub>2</sub></sup> all  $\mathbf{y}_i$  labelled as  $\omega_2$ . Then we require:  $\mathbf{a}^T \mathbf{y}_i > 0$  for all samples.

 $a_1$ 

 $y_2$ 

 $y_2$ 

solution

region

solution

region

 $\overline{y}_1$ 

 $J_p(a)$ 

 $a_2$ 

 $\mathbf{y}_{I}$ 

**y**3

 $y_3$ 

 $a_1$ 

#### **Perceptron Algorithm**

$$[\nabla J_{p}(\mathbf{a})]_{i} = \partial J_{p}(\mathbf{a}) / \partial a_{i} = \sum_{\mathbf{y} \in Y_{m}} (-y_{i})$$
 warning: Data are  
assumed normalized  
i.e.  $\nabla J_{p}(\mathbf{a}) = \sum_{\mathbf{y} \in Y_{m}} (-\mathbf{y})$  so gradient descent rule is :

$$\mathbf{a}(k+1) = \mathbf{a}(k) - \eta(k) \sum_{\mathbf{y} \in Y_m} (-\mathbf{y})$$

where  $Y_m$  is set of samples misclassified by  $\mathbf{a}(k)$ .

#### Batch Perceptron Algorithm

1. Initialize 
$$\mathbf{a}, \eta(\cdot)$$
, stopping criterion  $\theta, k \leftarrow 2. k \leftarrow k+1$   
3.  $\mathbf{a} \leftarrow \mathbf{a} + \eta(k) \sum_{\mathbf{y} \in Y_m} \mathbf{y}$   
4. If  $/\eta(k) \sum_{\mathbf{y} \in Y_m} \mathbf{y} |> \theta$ , repeat from 2.

<mark>η(k)</mark> is the learning rate or step size

0

# PART 3: INTRO TO NEURAL NETS

#### **Multilayer Neural Networks**

Linear discriminants are good for many problems but not **general** enough for demanding applications.

We can get more complex decision surfaces with nonlinear preprocessing,

 $y_i = \varphi_i(\mathbf{X})$ 

Where  $\varphi_i(.)$ , is, for example, a polynomial expansion to some order k.

But **too many free** parameters, so we may not have enough data points to fix them.  $\rightarrow$  *learn* which nonlinearities to use.

The best-known method is based on gradient descent: the so-called *backpropagation* algorithm.

#### **Three-Layer Network**



#### Operation

**Step 1**: each *d*-dimensional input vector  $(x_1, \ldots, x_d)$  is presented to input layer of the network and augmented with a bias term  $x_0 = 1$  to give  $x = (x_0, x_1, \ldots, x_d)$  [old school' term for input

'old school' term for input vector was input pattern.

**Step 2**: at each hidden layer we calculate the weighted sum of inputs to give the net activation:

$$net_j = \sum_{i=1}^d x_i w_{ji} + w_{j0} = \sum_{i=0}^d x_i w_{ji} = \mathbf{w}_j^T \mathbf{x}$$

where  $w_{ji}$  is the weight from the input unit *i* to the hidden unit *j* 

**Step 3**: The hidden unit emits the output:  $y_j = f(net_j)$  where f(.) is some nonlinear activation function

#### **Operation** (cont)

**Step 4**: At each output unit we calculate the weighted sum of the hidden layer units it is connected to giving:

$$n\tilde{e}t_k = \sum_{j=1}^{n_H} y_j \tilde{w}_{kj} + \tilde{w}_{k0} = \sum_{j=0}^{n_H} y_j \tilde{w}_{kj} = \tilde{\mathbf{w}}_k^T \mathbf{y}$$

where  $\tilde{w}_{kj}$  is the weight from the hidden unit j to the output unit k

**Step 5**: each output unit emits  $z_k = f(n\tilde{e}t_k)$ where f(.) is again the nonlinear activation function.

We can therefore think of the network as calculating c discriminant functions:

$$z_k = g_k(\mathbf{x})$$
#### **Example: XOR Problem**



#### **General Feedfoward Operation**

General form of output discriminant functions:

$$g_k(\mathbf{x}) \equiv z_k = f\left(\sum_{j=1}^{n_H} \tilde{w}_{kj} f\left(\sum_{i=1}^d w_{ji} x_i + w_{j0}\right) + \tilde{w}_{k0}\right)$$

[Note that the  $\tilde{w}_{ki}$ s are different from the  $w_{ii}$ s.]

Question: can every decision be implemented with a 3-layer network? <u>Answer</u> [Kolmogorov + others]: (in theory) Yes.

Typically needs very nasty functions at each layer.



See 6.2 and onwards Duda Hart

#### Backpropagation

We have seen that we can *approximate* any function but how do we *learn* functions?

Perceptrons (single layer network): each input affects the output via its weight: so we know which weight to change to reduce errors.

Multilayer networks (a.k.a *multilayer perceptron*, *MLP*): hidden units have no "teacher" – so how reduce the error?

This is known as the *credit assignment problem*.

Backpropagation solves the credit assignment problem using a smooth activation function f(.) and gradient descent.

#### Some known non-linearities (activation functions)

Name	Plot	Equation	Derivative
Identity		f(x) = x	f'(x) = 1
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0\\ 1 & \text{for } x \ge 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	f'(x) = f(x)(1 - f(x))
TanH		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$
Rectified Linear Unit (ReLU)		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \ge 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0\\ 1 & \text{for } x \ge 0 \end{cases}$
Parameteric Rectified Linear Unit (PReLU) <sup>[2]</sup>		$f(x) = \begin{cases} \alpha x & \text{for } x < 0\\ x & \text{for } x \ge 0 \end{cases}$	$f'(x) = \begin{cases} \alpha & \text{for } x < 0\\ 1 & \text{for } x \ge 0 \end{cases}$
Exponential Linear Unit (ELU) <sup>[3]</sup>		$f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0\\ x & \text{for } x \ge 0 \end{cases}$	$f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0\\ 1 & \text{for } x \ge 0 \end{cases}$
SoftPlus		$f(x) = \log_e(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

#### Some known non-linearities (activation functions)

Name	Plot	Equation	Derivative
Identity		f(x) = x	f'(x) = 1
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0\\ 1 & \text{for } x \ge 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	f'(x) = f(x)(1 - f(x))
TanH		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$

These were the "classical ones".

#### Some known non-linearities (activation functions)

Rectified Linear Unit (ReLU)	$f(x) = \begin{cases} 0 & \text{for } x < 0\\ x & \text{for } x \ge 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0\\ 1 & \text{for } x \ge 0 \end{cases}$
Parameteric Rectified Linear Unit (PReLU) <sup>[2]</sup>	$f(x) = \begin{cases} \alpha x & \text{for } x < 0\\ x & \text{for } x \ge 0 \end{cases}$	$f'(x) = \begin{cases} \alpha & \text{for } x < 0\\ 1 & \text{for } x \ge 0 \end{cases}$
Exponential Linear Unit (ELU) <sup>[3]</sup>	$f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0\\ x & \text{for } x \ge 0 \end{cases}$	$f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0\\ 1 & \text{for } x \ge 0 \end{cases}$
SoftPlus	$f(x) = \log_e(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

The newer ones. Some have contributed to the new AI revolution.

#### Backpropagation: why needed?

We have seen that we can *approximate* any function but how do we *learn* functions?

Perceptrons (single layer network): each input affects the output via its weight: so we know *which* weight to change to reduce errors.

Multilayer networks (a.k.a *multilayer perceptron*, *MLP*): hidden units have no "teacher" – so how can we reduce the error?

This is known as the *credit assignment problem*.

Backpropagation solves the credit assignment problem using a smooth activation function f(.), gradient descent and the chain rule.

## **Backpropagation Network**



#### **Backpropagation: Outline**

Step 1: start with an untrained network (random weights)

**Step 2**: present training data to network and calculate outputs:  $\mathbf{z}(n) = g(\mathbf{x}(n))$ This generates a training error:

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n} \sum_{k=1}^{c} \left( t_k(n) - z_k(n) \right)^2 = \frac{1}{2} \sum_{n} \|\mathbf{t}(n) - \mathbf{z}(n)\|^2$$

where **w** represents the set of all weights in the network and  $\mathbf{t}(n)$  are the target outputs

**Step 3**: change the weights in the direction of the negative gradient:

$$\Delta \mathbf{w} = -\eta \frac{\partial J}{\partial \mathbf{w}}$$
 or  $\Delta w_{pq} = -\eta \frac{\partial J}{\partial w_{pq}}$ 

where  $\eta$  is the learning rate (step size).

Step 4: iterate until convergence

MLSP

See 6.3 Duda Hart

#### But I know things...

I just need to be able to calculate

 $\frac{\partial J}{\partial w_{pq}}$ 

Indeed. For example, if we need to find:

 $\frac{\partial J}{\partial \tilde{w}_{kj}}$ 

we can use the *chain rule* to differentiate the training error, and then use general form of the forward pass, and use the chain rule again.

$$g_k(\mathbf{x}) \equiv z_k = f\left(\sum_{j=1}^{n_H} \tilde{w}_{kj} f\left(\sum_{i=1}^d w_{ji} x_i + w_{j0}\right) + \tilde{w}_{k0}\right)$$

True...but there is a more "elegant" way.

MLSP

See 6.2 and onwards Duda Hart

# Backpropagation: hidden-to-output

Let us calculate the gradient for a 3-layer network. We work backwards, starting at the **hidden-to-output** weights. The *chain rule* gives:

$$\frac{\partial J}{\partial \widetilde{w}_{kj}} = \sum_{n} \frac{\partial J}{\partial n \widetilde{e} t_k(n)} \frac{\partial n \widetilde{e} t_k(n)}{\partial \widetilde{w}_{kj}} = -\sum_{n} \widetilde{\delta}_k(n) \frac{\partial n \widetilde{e} t_k(n)}{\partial \widetilde{w}_{kj}}$$

where  $\tilde{\delta}_k(n) = -\frac{\partial J}{\partial n \tilde{e} t_k(n)}$  is termed the *sensitivity* of J to the net activation of k. Applying the chain rule again:

$$\tilde{\delta}_k(n) = -\frac{\partial J}{\partial n \tilde{e} t_k(n)} = -\frac{\partial J}{\partial z_k(n)} \frac{\partial z_k(n)}{\partial n \tilde{e} t_k(n)} = \underbrace{(t_k(n) - z_k(n))}_{\text{error}} f'(n \tilde{e} t_k(n))$$

where f'(.) is the *derivative* of f(.).

# Backpropagation: hidden-to-output

Let us calculate the gradient for a 3-layer network. We work backwards, starting at the **hidden-to-output** weights. The *chain rule* gives:

$$\frac{\partial J}{\partial \widetilde{w}_{kj}} = \sum_{n} \frac{\partial J}{\partial n \widetilde{e} t_k(n)} \frac{\partial n \widetilde{e} t_k(n)}{\partial \widetilde{w}_{kj}} = -\sum_{n} \widetilde{\delta}_k(n) \frac{\partial n \widetilde{e} t_k(n)}{\partial \widetilde{w}_{kj}}$$

where  $\tilde{\delta}_k(n) = -\frac{\partial J}{\partial n \tilde{e} t_k(n)}$  is termed the *sensitivity* of J to the net activation of k. Applying the chain rule again:

$$\tilde{\delta}_k(n) = -\frac{\partial J}{\partial n \tilde{e} t_k(n)} = -\frac{\partial J}{\partial z_k(n)} \frac{\partial z_k(n)}{\partial n \tilde{e} t_k(n)} = \underbrace{(t_k(n) - z_k(n))}_{\text{error}} f'(n \tilde{e} t_k(n))$$

where f'(.) is the *derivative* of f(.).

Since 
$$n \widetilde{e} t_k(n) = \sum_j \widetilde{w}_{kj} y_j(n) \Rightarrow \frac{\partial n \widetilde{e} t_k(n)}{\partial \widetilde{w}_{kj}} = y_j(n)$$

Hence:

$$\Delta \widetilde{w}_{kj} = -\eta \partial J / \partial \widetilde{w}_{kj} = \eta \sum_{n} \widetilde{\delta}_k(n) y_j(n) = \eta \sum_{n} \left( t_k(n) - z_k(n) \right) f'(n \widetilde{e} t_k(n)) y_j(n)$$

# Backpropagation: input-to-hidden

Again using the chain rule for **input-to-hidden** units we have:

$$\frac{\partial J}{\partial w_{ji}} = \sum_{n} \frac{\partial J}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial net_j(n)} \frac{\partial net_j(n)}{\partial w_{ji}}$$

(Note that  $\frac{\partial J}{\partial y_j(n)}$  involves all outputs k = 1, ..., c.)

$$\frac{\partial J}{\partial y_j(n)} = \frac{\partial}{\partial y_j(n)} \left[ \frac{1}{2} \sum_{k=1}^c \left( t_k(n) - z_k(n) \right)^2 \right] = -\sum_{k=1}^c \left( t_k(n) - z_k(n) \right) \frac{\partial z_k(n)}{\partial y_j(n)}$$
$$= -\sum_{k=1}^c \left( t_k(n) - z_k(n) \right) \frac{\partial z_k(n)}{\partial n \tilde{e} t_k(n)} \frac{\partial n \tilde{e} t_k(n)}{\partial y_j(n)} = -\sum_{k=1}^c \left( t_k(n) - z_k(n) \right) f'(n \tilde{e} t_k(n)) \tilde{w}_{kj}$$
$$= -\sum_{k=1}^c \widetilde{\delta}_k(n) \widetilde{w}_{kj}(n)$$

since  $\tilde{\delta}_k(n) = (t_k(n) - z_k(n)) f'(n\tilde{e}t_k(n))$ 

# Hidden Units (cont)

 $\frac{\partial y_j(n)}{\partial net_j(n)} = f'(net_j(n)),$ 

thus we can define

We also have

$$\delta_j(n) \equiv -\frac{\partial J}{\partial net_j(n)} = -\frac{\partial J}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial net_j(n)} = f'(net_j(n)) \sum_{k=1}^c \tilde{w}_{kj} \tilde{\delta}_k(n)$$

as the *hidden unit sensitivities*. Note that the output sensitivities are propagated back to the hidden unit sensitivities; hence the name "**back-propagation of errors**".

Finally we have:

$$\frac{\partial net_j(n)}{\partial w_{ji}} = x_i(n)$$

So the update rule for the input-to-hidden weights is:

$$\Delta w_{ji} = -\eta \frac{\partial J}{\partial w_{ji}} = \eta \sum_{n} x_i(n) \delta_j(n) = \eta \sum_{n} \left[ \sum_{k=1}^c \tilde{w}_{kj} \tilde{\delta}_k(n) \right] f'(net_j(n)) x_i(n)$$

#### **Backwards propagation**



The backpropagation algorithm can be easily generalized to any network with feed-forward connections, e.g. more layers, different nonlinearities in different layers, etc.

#### Criticisms of MLPs

MLPs provide one way to achieve the required expressive power needed to build general classifiers. However they do have their weaknesses:

- Possibility of multiple (local) minima
- $-g(\mathbf{x})$  is nonlinear in terms of the weights: this makes training slow.
- ad hoc solution (how many units, hidden layers, etc?)

Other popular discriminant learning structures:

- Support Vector Machines (SVMs), left as "self-learning"
- Convolutional Networks (CNNs), Sen Wang will cover these

#### **TensorFlow Demo**

- <u>https://playground.tensorflow.org/</u>
- Start with the bottom left (linear activation)
  - Observations: linearly separable
  - Start with a very simple net 1 layer 1 node (linear activations)
- Move to upper right still separable but not linear in input space
  - Options to fix: change layer size, add layers, change activation, change inputs?

# PART 4: SUPPORT VECTOR MACHINES





# **Support Vector Machines**

# Support Vector Machines

Salastin SVMs solve a problem in linear or non-linear space by projecting the input space into a new (possibly infinite) space. **Principle of Support Vector Machines** 



# **Support Vector Machines**

SVMs aim to solve the linearly separable problem. First map feature space into high (possibly  $\infty$ —) dimensional space.

Then find separating hyperplane with **maximal margin**. Recall, intuitively we are more confident in classifying point far away from the decision boundary.



Learning takes the form of a constrained optimization scheme.

#### Ch 5.11 Duda

## Why max margin

- Salar For example these possible hyperplanes. Which one you will choose?  $x_2$ direction 2
  - "Direction 1" has a narrow margin thus on unseen (test) data it has higher likelihood of error.
  - Clearly 'direction 2' is preferred



 $x_1$ 

# SVMs: Maximizing the Margin

Suppose that we have a margin  $\gamma$ , such that  $\omega^{(i)}a^Ty_i \ge \gamma$  for all points, i = 1, 2, ..., n. We therefore want to solve the following:

max 
$$\gamma$$
, such that:  $\omega^{(i)} a^T y_i \ge \gamma$   
and:  $||a|| = 1$ 

This is a messy optimization problem. However, equivalently we can solve:

#### min $||a||^2$ , such that: $\omega^{(i)}a^Ty_i \ge 1$

That is, we search for the minimum size weight vector that is able to separate the data with a margin of  $\gamma = 1$ .

This form of the problem is a constrained quadratic optimization problem. It is convex and (relatively) easy to solve.

## The Support Vectors

Hearming. Interestingly the Max margin solution only depends on a subset of the training data – those that lie exactly on the margin (why?). These are called the **support vectors** (SVs).

#### Also

- the support vectors also define the equation of the optimal hyperplane: a ~ weighted sum of  $\omega^{(i)}y_i$
- The hyperplane is unique the SVs are not unique. Why?



# **SVM Generalization Error**

- strand and a strange of the state of the sta Cross Validation (CV) - Break the data into a training set and a testing set. Use the training data to learn the classifier then evaluate on the test data
  - Leave-one-out CV: choose one data point at random as the testing set.

 $\{x_1, x_2, ..., x_p, x_{p+1}, ..., x_n\}$ 

• Note the LOOCV error will be unaffected unless  $x_p$  is a support vector. Therefore we have:

$$error_{LOOCV} \in \frac{\# \text{ of SVs}}{n}$$

 Therefore the number of SVs tells us how confident we are in the SVM.

Ch 5.11 Duda

# PART 5: FEATURE ENGINEERING

#### Principal Component Analysis

MLP can learn "functions" of data but with high dimensional inputs they need some help!

This topic concerns decomposing signals into useful low dimensional subsets:

- For feature space selection in classification
- For redundancy reduction
- To avoid overfitting
- For signal separation

The key aim is to find a linear transform of the data that better represents the underlying information

e.g. Fourier transform of an image concentrates information into low frequencies

Suppose I want to characterize fish population. Measure length/breadth and plot.



Subtract mean from each axis (**center**) Note relationship between data matters only.



Suppose I want to characterize fish population. Measure length/breadth and plot.

![](_page_100_Figure_2.jpeg)

Subtract mean from each axis (**center**) Note relationship between data matters only.

![](_page_100_Figure_4.jpeg)

I can **move the axis!** But how to chose them? I can rotate them.

![](_page_100_Figure_6.jpeg)

What is my objective: position one axis in such a way that it accounts for the largest proportion of the data's variance.

Inspired by <u>http://www.cmbi.ru.nl/edu/bioinf4/prac-</u> microarray/stats/PCA%20graphical%20explanation.htm

Suppose I want to characterize fish population. Measure length/breadth and plot.

![](_page_101_Figure_2.jpeg)

What is my objective: position one "new" axis in such a way that it accounts for the largest proportion of the data's variance.

Inspired by <u>http://www.cmbi.ru.nl/edu/bioinf4/prac-</u> microarray/stats/PCA%20graphical%20explanation.htm

What might you **call** this new axis? **size = length + breadth** However, we could make one of the **variables** more important.

#### size = 0.75 x length + 0.25 x breadth

[these are weights and are important! They tell us the "significance" of each of the original variables.

#### What about the second axis of the ellipse?

Objective: account for as much of the remaining variation as possible but must also be uncorrelated (**orthogonal**) with the first. [trivial in 2D]

#### Apart from size, how else do the above fish differ?

Not much, minor differences in shape. If we ignore 2<sup>nd</sup> axis (after rotation) we would lose information about the different shapes, but since they are all very similar in shape we <u>wouldn't</u> lose much information at all.

Thus, in the above example we can reduce the data's dimensionality from two (length and breadth) to one (size), with little information loss.

#### **PCA & Subspace Projections**

Introduced by Pearson in 1901 "On lines and planes of closest fit to a system of points in space." (a.k.a. Karhunen-Loéve transform (KLT), or the Hotelling transform,...).

Suppose our data consists of d-dimensional vectors  $x^{(n)} \in R^d$ and we want a low-dimensional approximation for the data.

Let  $u_i$  be an **orthonormal** basis for  $R^d$  (i.e.  $U^T U = I$ ) then we can approximate x by:

let  $z_i = \mathbf{u_i}^T \mathbf{x}$ 

$$\tilde{\mathbf{x}} = \sum_{i=1}^{M} (\mathbf{u}_{i}^{T} \mathbf{x}) \mathbf{u}_{i} + \sum_{j=M+1}^{d} b_{j} \mathbf{u}_{j}$$
Each b<sub>j</sub> is a constant

111

#### **Subspace Projections**

x has d degrees of freedom while  $\tilde{\mathbf{x}}$  has M deg. of freedom.

**Principal Component Analysis** – choose  $u_i$  and  $b_j$  to best approximate x in the LSE sense, ie., minimize  $E_M$ :

$$E_{M} = \frac{1}{2} \sum_{n=1}^{N} \left\| \mathbf{x}^{(n)} - \widetilde{\mathbf{x}}^{(n)} \right\|^{2} = \frac{1}{2} \sum_{n=1}^{N} \left\| \mathbf{U}^{T} \mathbf{x}^{(n)} - \mathbf{U}^{T} \widetilde{\mathbf{x}}^{(n)} \right\|^{2} = \frac{1}{2} \sum_{n=1}^{N} \sum_{j=M+1}^{d} (z_{j}^{(n)} - b_{j})^{2}$$

Taking the derivative with respect to *b<sub>i</sub> gives*:

$$\sum_{n=1}^{N} \left( z_{j}^{(n)} - b_{j} \right) = 0 \implies b_{j} = \frac{1}{N} \sum_{n=1}^{N} z_{j}^{(n)} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{u}_{j}^{T} \mathbf{x}^{(n)} = \mathbf{u}_{j}^{T} \overline{\mathbf{x}}$$

Subspace Projections (cont.)  
So we can write:  

$$E_{M} = \frac{1}{2} \sum_{n=1}^{N} \sum_{j=M+1}^{d} \left( \mathbf{u}_{j}^{T} \left( \mathbf{x}^{(n)} - \overline{\mathbf{x}} \right) \right)^{2}$$

$$= \frac{1}{2} \sum_{j=M+1}^{d} \sum_{n=1}^{N} \mathbf{u}_{j}^{T} \left( \mathbf{x}^{(n)} - \overline{\mathbf{x}} \right) \left( \mathbf{x}^{(n)} - \overline{\mathbf{x}} \right)^{T} \mathbf{u}_{j} = \frac{1}{2} \sum_{j=M+1}^{d} \mathbf{u}_{j}^{T} \mathbf{R}_{\mathbf{x}} \mathbf{u}_{j}$$

where  $\mathbf{R}_{\mathbf{x}}$  is the sample covariance matrix for  $\{\mathbf{x}^{(n)}\}$ 

Minimizing  $E_M$  with respect to  $u_i$  is satisfied by:

 $\wedge$ 

$$\mathbf{R}_{\mathbf{x}} \mathbf{u}_{j} = \lambda \mathbf{u}_{j}, \text{ for } j = 1, ..., \mathbf{M}$$

i.e. the M basis vectors are the principal eigenvectors of the sample covariance matrix.

#### PCA Algorithm and Projection on Principal Directions

The PCA algorithm is summarized as follows:

- 1. subtract mean  $\overline{x}$  from data
- 2. Calculate sample covariance matrix,  $\mathbf{R}_{\mathbf{x}}$  for  $\{\mathbf{X}_{n} \overline{\mathbf{X}}\}$
- 3. Perform eigenvalue decomposition:  $\mathbf{R}_{x} = \mathbf{U} \Lambda \mathbf{U}^{T}$
- 4. Approximate data by the first M components that have the largest eigenvalue:

$$\mathbf{x}_n \approx \overline{\mathbf{x}} + \sum_{i=1}^m (\mathbf{u}_i^T (\mathbf{x}_n - \overline{\mathbf{x}})) \mathbf{u}_i$$

**Recall** by design eigenvectors are **ORTHOGONAL** with each other  $\mathbf{u_i}^T \mathbf{u_j} = \mathbf{0}$  (i # j) and have length 1 ( $\mathbf{u_i}^T \mathbf{u_j} = \mathbf{1}$ ).

#### PCA example: eigenfaces

In face recognition a common practice is to first project data (after alignment) onto a low dimensional PCA space,

e.g. images from images AT&T Laboratories Cambridge.

![](_page_107_Figure_3.jpeg)

 $U_1$   $U_2$   $U_3$   $U_4$ Eigenfaces capture appearance and lighting conditions quite well.
### An interactive demo

<u>http://setosa.io/ev/principal-component-analysis/</u>

### Kernel PCA

(*n*) **PCA** method is built upon the eigenanalysis of (*n*) number of data points)

$$R_X = \frac{1}{n} X^T X = \frac{1}{n} \sum_{i=1}^n \underline{x}_i \underline{x}_i^T$$

There is an equivalent built upon the eigenanalysis of

$$K = XX^{T} = \begin{bmatrix} \underline{x}_{1}^{T} \underline{x}_{1} & \dots & \underline{x}_{1}^{T} \underline{x}_{n} \\ \vdots & \ddots & \vdots \\ \underline{x}_{n}^{T} \underline{x}_{1} & \dots & \underline{x}_{n}^{T} \underline{x}_{n} \end{bmatrix}$$

known as the Gram matrix.

The kernel PCA method is the kernelized version of this, where inner products are replaced by kernel operations.

### Kernels

- Self > Achieve projections in higher dimensions as we saw in GLD without actually explicitly defining the projection
  - Examples of kernels
    - Radial basis Functions:

$$K(\underline{x},\underline{z}) = \exp\left(-\frac{\left\|\underline{x}-\underline{z}\right\|^2}{\sigma^2}\right)$$

• Polynomial:

$$K(\underline{x},\underline{z}) = (\underline{x}^T \underline{z} + 1)^q, \ q > 0$$

• Hyperbolic Tangent:

$$K(\underline{x},\underline{z}) = \tanh(\beta \underline{x}^T \underline{z} + \gamma)$$

for appropriate values of  $\beta$ ,  $\gamma$ .



### Kernel PCA: The algorithm

• Compute the Gram matrix.

$$K(i, j) = K(\underline{x}_i, \underline{x}_j)$$
,  $i, j = 1, 2, ..., n$ 

• Compute the *m* dominant eigenvalues / eigenvectors.

$$\lambda_k, \underline{a}_k, k = 1, 2, ..., m$$

• Perform normalization to unity.

$$1 = n\lambda_k \underline{a}_k^T \underline{a}_k, \ k = 1, 2, \dots, m$$

• Given a vector <u>x</u>, perform the following "nonlinear mapping".

$$y(k) = \sum_{i=1}^{n} a_k(i) K(\underline{x}_i, \underline{x}), \ k = 1, 2, ..., m$$

Remark

elfer anno

• The kernel PCA is equivalent with performing a (linear) PCA in a Reproducing kernel Hilbert space (RKHS) *H*, after a mapping

$$\underline{x} \to \phi(\underline{x}) \in H$$

• It can be shown that the dominant eigenvectors of

$$\frac{1}{n}\sum_{i=1}^{n}\phi(\underline{x}_{i})\phi^{T}(\underline{x}_{i})$$

are given in terms of the dominant eigenvectors of the Gram matrix, i.e.,

$$\underline{v}_k = \sum_{i=1}^n a_k(i)\phi(\underline{x}_i), \ k = 1, 2, ..., m$$

Hence the projection of  $\phi(\underline{x})$  on  $\underline{v}_k$  is given by:

 $\left\langle \underline{v}_{k}, \phi(\underline{x}) \right\rangle = \sum_{i=1}^{n} a_{k}(i) \left\langle \phi(\underline{x}_{i}), \phi(\underline{x}) \right\rangle = \sum_{i=1}^{n} a_{k}(i) K(\underline{x}_{i}, \underline{x}),$ using the properties of the RKHS. 120

## PART 6: BIAS AND VARIANCE DILEMMA

### Why we care about generalisation?

- ✤ We can always train an algorithm to zero "training" error.
  - e.g. neural networks can approximate any function
  - > Does this really matter?
- What matters is how the algorithm will perform on new unseen data that are NOT the same as the training data
  - > How will the algorithm then behave?
  - If the algorithm has understood the underlying structure
     then it should do well
- We say then that the algorithm "generalizes well" (we now formalize this using the bias and variance dilemma)









### Overfitting in ML

Given a limited data set maximising the likelihood  $L(\mathbf{w})$  may lead to *overfitting*. If the model order (dimensions of  $\mathbf{w}$ ) is large enough we can 'exactly' fit model to data.



### Overfitting in ML

States Overfitting relates to Model Order Selection. As with estimating power spectra the problem comes from a *bias - variance* trade-off. Consider the expected error between an optimal estimator F(x) (ie. the true function) and an empirical estimator g(x). [We use a trick.]





### Example in regression

Identify high bias / variance cases:



### Ignore the data **→**

- \* Big approximation errors (high bias)
- No variation between data sets (no variance)

Use all data 🗲

- No approximation errors (zero bias)
  - \* Variation between data sets (high variance)

### Overfitting in ML

These examples tell us:

- Too **complex** a model  $\rightarrow$  **high variance**.
- Too **simple** a model  $\rightarrow$  **high bias** (the simplest model is the fixed value).

### Possible solution on diagnosis for models and data (size)

<u>**Cross-validation**</u>: Break the data into **3 pieces**: *training*, *validation* and *testing* sets. Set the testing set apart, do not touch it, till the end.

Use the <u>training</u> data to learn the model parameters. But identify the best model order (and other parameters such as step size, regularizers, type of activation function etc) using the <u>validation</u> ("out-of-sample") set.

If you are happy with the performance on the validation set, <u>then</u> you can now take the testing set and run the algorithm on the <u>testing</u> set. Then report results on the training and testing sets.

# Diagnosing bias/variance: the practical picture

• Finding a good model size



### Bias vs variance: on data size

• Typical curve with high variance



- Test error decreases as m increases → larger training set may help
- There is a gap between training & test error

### Bias vs. variance: on data size

• Typical curve with high bias



• We are in trouble: training error is high too

• Small gap between the two errors

### Take away messages

- Supervised learning relies on examples to learn
  - Unsupervised does not
    - Semi-supervised combines them
- It is easier to classify when data are **linearly** separable
- We can always **project** to new features to solve linear problems
  - By design (e.g. via polynomial functions)
  - Or we can **summarise** them statistically (e.g. PCA)
  - Or we can **learn** them (via NNs)
- We cannot create **magic** 
  - It is not training error we care for but generalization to unseen data and hence testing error
    - We have **theoretical** tools (bias and variance dilemma) and **practical** tools to assess this

### Thank you

- Sotirios (Sotos) Tsaftaris
- Email: <u>S.Tsaftaris@ed.ac.uk</u>
- Web: <u>https://vios.science</u>

 Additional slides on machine learning and deep learning and particularly in augmentation are here: <u>https://edin.ac/2vSP6T3</u> (Tutorial at ICASSP 2019 on Data augmentation techniques for deep learning )

## PART 7: Convolutional Neural Networks & Deep Learning

# Convolutional Neural Networks & Deep Learning

# From "hand-designed" feature spaces to data-driven ones...



# What if we could "extract" automatically good features....











http://fortune.com/ai-artificial-intelligence-deep-machine-learning/

### **Convolutional Networks (CNNs)**

Consider the problem of building a classifier that is insensitive to translation (or scale, rotation, etc.). A Convolutional Network encodes the invariance within the MLP structure [LeCun 1998].



Figure from: LeCun, Boser, Denker, Henderson, Howard, Hubbard, Jackel. (December 1989). "Backpropagation Applied to Handwritten Zip Code Recognition". Neural Computation. **1** (4): 541-51.

### **Convolutional Networks (CNNs)**

Consider the problem of building a classifier that is insensitive to translation (or scale, rotation, etc.). A Convolutional Network encodes the invariance within the MLP structure [LeCun 1998].









For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a "new image" of size 28x28x6!

ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



### **Convolutional Networks (CNNs)**

- Connections are restricted: hidden units are connected identically to neighbours to encode shift/delays
- Training using back prop. but with ties weights across shifted units (*weight sharing*)
- The resulting MLP has much fewer weights to train than a traditional MLP

Figure from Goodfellow, Bengio and Courville, Deep Learning, MIT Press, 2016



### **Convolutional Networks (CNNs)**

 Growing "influence" with depth →

 (the top layer is influenced by a large span of inputs despite only connected to 3 previous units)



## Pooling

- Max pooling helps with invariance and down sampling reduces number of parameters
  - Reduces dimensionality of representations

Single depth slice



Max pooling with 2×2 filters and stride 2



### Pooling on the layers

Important to remember that pooling operates on each activation map independently



### **CNNs & Deep Learning**

- CNNs are often considered the ancestors of Deep Learning.
- The idea is to use MLP with many (≥ 3) hidden layers. This involves lots of 'tricks' to make training work: convolution, subsampling, pooling of outputs,..., pretraining (helps to start the weights with good initial values)





but exhibit state-of-the-art performance... e.g. see - http://www.cs.nyu.edu/~yann/research Application to face detection and pose estimation

### Deep learning has many hyperparameters

- Network design
- Activation function
- Optimisation choices
- Lots of 'tricks' to make training work: <u>convolution</u>, <u>subsampling</u>, <u>pooling of</u> <u>outputs</u>,..., <u>good initialisation</u>, <u>pretraining</u> (helps to start the weights with good initial values), batch normalisation, learning rate schedules, Adam, RMSprop, SELu, PreLu, noise injection, label smoothing,...,

CS 4620 Intelligent Systems

Changing random stuff until your program works is "hacky" and "bad coding practice."

But if you do it fast enough it is "Machine Learning" and pays 4x your current salary.

..., ...

#### MLSP

Image source: I saw it on https://www.reddit.com/r/MachineLearning/, but supposedly it is taken from a course at CS 4620. No copyright claimed. Original source unknown.

, ..., ..., ..., ..., ..., ..., ...,

### More Demos

<u>http://cs.stanford.edu/people/karpathy/convnetjs/</u>
## Some additional info

 Additional slides on machine learning and deep learning and particularly in augmentation are here: <u>https://edin.ac/2vSP6T3</u> (Tutorial at ICASSP 2019 on Data augmentation techniques for deep learning )