



#### The Alan Turing Institute

### Introduction to Machine Learning

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

Turing Fellow http://tsaftaris.com S.Tsaftaris@ed.ac.uk

DRAFT, released for enhanced learning; final version to be made available electronically on the day.



1

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.
- Material subject to change.



# 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

# What is machine learning?

• What do you think?

#### 5

# Some 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
- <u>http://www.youtube.com/watch?v=OIOtOmyySQo</u>
- Google car
- <u>http://www.youtube.com/watch?v=cdgQpa1pUUE</u>

### 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.

7

# **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>



 $\mathbb{A}\mathbb{P}$ 

9

# 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...



# Counting objects (new school)



Giuffrida et al. The Plant Journal Pheno-Deep Counter: a unified and versatile deep learning architecture for leaf counting

12

# 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

• How would I look in 30 years?







• A much harder task is how would my brain look?







# Doing more with less

· Build systems that learn from few data, few annotations





- 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.























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
÷		•





### Solving the linear regression problem

• The LMS algorithm

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

• Define a cost:  $J(\theta) = \frac{1}{2} \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^2.$ 

• Optimise for the cost 
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta).$$

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

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

For a dataset of size 1

33

#### Batch vs stochastic







• Can we solve the classification problem as linear regression?

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

• The logistic function (sigmoid function)  $\Rightarrow$   $g(z) = \frac{1}{1 + e^{-z}}$ 

36

Why?

It has some nice property

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

37

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

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

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

$$L(\theta) = p(\vec{y} \mid X; \theta)$$

Objective function

• Assuming the training examples were generated independently  $\ell(\theta) = \log L(\theta)$ m $= \prod_{i=1}^{m} p(y^{(i)} | x^{(i)}; \theta)$  $= \prod_{i=1}^{m} (h_{\theta}(x^{(i)}))^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$ 

• Or...

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

Contd...  
• Lets use gradient ascent  

$$\begin{aligned}
& \hat{\theta} := \theta + \alpha \nabla_{\theta} \ell(\theta) \\
& = \left(y \frac{1}{g(\theta^{T}x)} - (1 - y) \frac{1}{1 - g(\theta^{T}x)}\right) \frac{\partial}{\partial \theta_{j}} g(\theta^{T}x) \\
& = \left(y \frac{1}{g(\theta^{T}x)} - (1 - y) \frac{1}{1 - g(\theta^{T}x)}\right) g(\theta^{T}x)(1 - g(\theta^{T}x) \frac{\partial}{\partial \theta_{j}} \theta^{T}x) \\
& = (y(1 - g(\theta^{T}x)) - (1 - y)g(\theta^{T}x))x_{j} \\
& = (y - h_{\theta}(x))x_{j}
\end{aligned}$$

# What if we take an extreme sigmoidal



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







### 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 <i>t</i> =1 but <i>y</i> =0, set new w = old w + ηx	[punish]
b) If <i>t</i> =0 but <i>y</i> =1, set new w = old w – $\eta$ x	[punish]

c) Otherwise (If y=t), do nothing

[reward]

Repeat until correct for all input vectors. Factor  $\eta$  is called the *learning rate* 







### **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 GLD

### **Linear Discriminant Functions**

Suppose we wished to decide whether some data

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

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



### **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.



# Generalized Linear Discriminants

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

$$g(\mathbf{X}) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j$$
 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  $\mathbf{x}$ .

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

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  $g(\mathbf{y}) = \mathbf{a}^{\mathsf{T}}\mathbf{y} + a_0$  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 reduced 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.





# 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.



#### Operation

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

**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* 

66

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

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

$$n\widetilde{e}t_{k} = \sum_{j=1}^{n_{H}} y_{j}\widetilde{w}_{kj} + \widetilde{w}_{k0} = \sum_{j=0}^{n_{H}} y_{j}\widetilde{w}_{kj} = \widetilde{\mathbf{w}}_{k}^{T}\mathbf{y}$$

where  $\widetilde{w}_{ki}$  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:

$$z_k = g_k(\mathbf{X})$$



### **General Feedfoward Operation**

General form of output discriminant functions:

$$g_{k}(\mathbf{X}) \equiv z_{k} = f\left(\sum_{j=1}^{n_{H}} \widetilde{w}_{kj} f\left(\sum_{i=1}^{d} w_{ji} x_{i} + w_{j0}\right) + \widetilde{w}_{k0}\right)$$

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

Question: can every decision be implemented with a 3-layer network?

Answer [Kolmogorov + others]: (in theory) Yes.

Typically needs very nasty functions at each layer.



#### 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

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	$\square$	$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
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		$f(x) = \log_e(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$



### **Backpropagation Outline**

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

**Step 2**: present training data to network and calculate outputs: z(n) = g(x(n))

This generates a training error:

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

Where  $\mathbf{w}$  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:

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

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

Step 4: iterate until convergence

**Backpropagation of Errors** 

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

$$\widetilde{\delta}_{k}(n) = -\frac{\partial J}{\partial n \widetilde{e} t_{k}(n)} = -\frac{\partial J}{\partial z_{k}(n)} \frac{\partial z_{k}(n)}{\partial n \widetilde{e} t_{k}(n)} = \underbrace{(t_{k}(n) - z_{k}(n))f'(n \widetilde{e} t_{k}(n))}_{\text{error}}$$

where f'(.) is the derivative of f(.). Noting that:

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

Hence:

$$\Delta \widetilde{w}_{kj} = -\eta \partial J / \partial \widetilde{w}_{kj} = \eta \sum_{n} \widetilde{\delta}_{k}(n) y_{j}(n) = \eta \sum_{n} (t_{k}(n) - z_{k}(n)) f'(n \widetilde{e} t_{k}(n)) y_{j}(n)$$

#### 74

# Backpropagation to Hidden Units

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

$$\frac{\partial J}{\partial w_{ji}} = \sum_{n} \frac{\partial J}{\partial y_{j}(n)} \frac{\partial y_{j}(n)}{\partial net_{j}(n)} \frac{\partial net_{j}(n)}{\partial w_{ji}(n)}$$

Note that  $\partial J / \partial y_i(\mathbf{n})$  involves all outputs k

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

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

Hidden Units (cont)

We also have  $\partial y_i(n) / \partial net_i(n) = f'(net_i(n))$  so we can define

$$\delta_{j}(n) \equiv -\frac{\partial J}{\partial net_{j}(n)} = -\frac{\partial J}{\partial y_{j}(n)} \frac{\partial y_{j}(n)}{\partial net_{j}(n)} = f'(net_{j}(n)) \sum_{k=1}^{c} \widetilde{w}_{kj} \widetilde{\delta}_{k}(n)$$

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

Finally we have:

 $\partial net_i(n) / \partial w_{ii} = x_i(n)$ 

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

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

-	$\sim$
1	n
	U.



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)
- Convolutional Networks (CNNs)

### **TensorFlow Demo**

- 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?



### **Support Vector Machines**

 SVMs solve a problem in linear or non-linear space by projecting the input space into a new (possibly infinite) space.



# **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.

y

Learning takes the form of a constrained optimization scheme.

Ch 5.11 Duda



82

# SVMs: Maximizing the Margin

Suppose that we have a margin  $\gamma$ , such that  $\omega^{(i)}a^T y_i \ge \gamma$  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.

Ch 5.11 Duda

84

# The Support Vectors

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 ω<sup>(i)</sup>y<sub>i</sub>
- The hyperplane is unique the SVs are not unique. Why?



SVMs can be shown to generalize well in terms of cross validation error

Ch 5.11 Duda

### **SVM Generalization Error**

- 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<sub>p</sub> is a support vector. Therefore we have:

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

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

Ch 5.11 Duda

# 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

### PCA – Graphical Intuition

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



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







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





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

### PCA – Graphical Intuition

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  $\mathbf{x}^{(n)} \in \mathbb{R}^d$ and we want a low-dimensional approximation for the data.

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

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

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

#### **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  $b_i$  gives:

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

$\sim$	
y,	4

#### 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:

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

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

# 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.

#### 96

# PCA Algorithm Summary 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  $\mathbf{u_i}^T \mathbf{u_i} = \mathbf{0}$  (i # j) and have length 1 ( $\mathbf{u_i}^T \mathbf{u_i} = \mathbf{1}$ ).

http://setosa.io/ev/principal-component-analysis/

An interactive demo

#### Kernel PCA

(linear) 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.

99

> The kernel PCA algorithm

• Compute the Gram matrix.

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

• Compute the *m* dominant eigenvalues / eigenvectors.

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

• 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".

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

100

#### Remark

• 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:

$$\langle \underline{v}_k, \phi(\underline{x}) \rangle = \sum_{i=1}^n a_k(i) \langle \phi(\underline{x}_i), \phi(\underline{x}) \rangle = \sum_{i=1}^n a_k(i) K(\underline{x}_i, \underline{x}),$$
  
using the properties of the RKHS.

# WHY CARE ABOUT FEATURES AND HOW MANY



5-103

102









# 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







# **DEEP LEARNING**

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

2

92

6

We see





112



# **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].



LeNet: a layered model composed of convolution and subsampling operations followed by a holistic representation and ultimately a classifier for handwritten digits. [ LeNet ]

# **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].

Let us first review 2D convolution  $\rightarrow$ 

Can I write *w conv I* as a matrix multiplication?

Next layer O = W\*I W must have a correct shape and I: a vector (input layer)!



Inspired by Goodfellow et al Deep learning book

116

# **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





# **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

