# Approximation of functions

A multilayer perceptron trained through BPA can be an engine for learning the non-linear I/O mappings

$$f : \mathbb{R}^{m_0} \longrightarrow \mathbb{R}^{m_L}$$

Space of I/P neurons

Space of o/p neurons

Qn: What is the min # of hidden layers needed for a MLP to learn the I/O mapping?

# Universal Approximation Theorem

(Key: 1 single layer of hidden neurons will suffice)

Let $\varphi(\cdot)$ be a nonconstant, bounded and monotone non-decreasing, continuous function. Let $I_{m_0}$ denote a $m_0$-dim. unit hyper cube $[0,1]^{m_0}$. The space of cont. functions on $I_{m_0}$ is denoted by $C(I_{m_0})$. Then, given any function $f \in C(I_{m_0})$ and $\varepsilon > 0$ $\exists$ an integer $m_1$ & sets of constants $\alpha_i, b_i, w_{ij}$ $i = 1, \ldots, m_1$ & $j = 1, \ldots, m_0$

$$F(x, \ldots x_{m_0}) = \sum_{i=1}^{m_1} \alpha_i \varphi\left( \sum_{j=1}^{m_0} w_{ij} x_j + b_i \right) \text{ is an}$$

$\underbrace{\text{approximation}}$ of the true $f(\cdot)$ / $\| F(\cdot) - f(\cdot) \|_P < \varepsilon$

← typically $p = 2$

# Bounds on Approximation errors

Barron (1993) established the approx. properties of a MLP
Suppose $f$ is the fn. representing the data points
The n/w gives us the approx. function $\hat{F}$.

$$F \sim \hat{f}$$

Let $\tilde{f}(\underline{w})$ denote the multi-dimensional Fourier transform
of $f(\underline{x})$; $\underline{x} \in \mathbb{R}^{m_o}$; $\underline{w}$ is the frequency
vector

$$f(\underline{x}) = \int_{\mathbb{R}^{m_0}} \tilde{f}(\underline{\omega}) \, e^{j\underline{\omega}^T \underline{x}} \, d\underline{\omega}$$

For a complex valued fn $\tilde{f}(\underline{\omega})$ for which $\underline{\omega}\,\tilde{f}(\underline{\omega})$ is integrable, let

$$C_f = \int_{\mathbb{R}^{m_0}} |\tilde{f}(\underline{\omega})| \, \|\underline{\omega}\|^{\frac{1}{2}} \, d\underline{\omega}$$

first absolute moment of the Fourier mag. distribution of $f$

$C_f$ measures smoothness of $f$

Consider a ball $B_r = \{ \underline{x} : \| \underline{x} \| \le r \}$ $r > 0$

Theorem: For every cont. function $f(\underline{x})$ with finite first moment $C_f$ and every $m_1 \ge 1$ $\exists$ a linear combination of sigmoid based functions $F(\underline{x})$ / when $f(\underline{x})$ is observed over ips $\underline{x}$ $\{ \underline{x}_i \}_{i=1}^{N}$ restricted to $B_r$, the empirical risk

$$\mathcal{E}_{av} = \frac{1}{N} \sum_{i=1}^{N} \left( f(\underline{x}_i) - F(\underline{x}_i) \right)^2 \le \frac{\left( 2 r C_f \right)^2}{m_1}$$

$$\mathcal{E}_{av}(N) \leq O\left(\frac{C_f^2}{m_1}\right) + O\left(\frac{m_0 \, m_1}{N} \log N\right)$$

$$m_1 \leq C_f \left(\frac{N}{m_0 \log N}\right)^{\frac{1}{2}}$$

$$\mathcal{E}_{av}(N) = O\left(\left(\frac{1}{N}\right)^{\frac{2s}{2s + m_0}}\right)$$

$s$ is a measure of smoothness

$$O\left(\left(\frac{1}{N}\right)^{\frac{1}{m_0}}\right)$$

$$O\left(\frac{1}{N^{1/m_0}}\right) \leftarrow \text{density}$$

Sampling density $\propto$ $N^{\frac{1}{m_0}}$

$\implies$ We need to densely sample the data points to learn the function well.

$\implies$ Higher dim $\implies$ Exponential growth in the complexity of the n/w algo.

# A few other technicalities

1) Jacobians
2) Hessians

## Jacobian :

Let $W$ denote the total # of free parameters i.e., the synaptic weights and biases of a MLP to form a vector $\underline{W}$.

Let $N$ denote the examples presented to the n/w.

Using the BPA, we can get $F(\underline{W}, \{\underline{x}\})$

Suppose $\{\underline{x}_n\}_{n=1}^{N}$

$$J := \left[ \frac{\partial F(\underline{w}, \underline{x}_n)}{\partial w_{ij}} \right] \qquad \begin{array}{l} i = 1, \ldots, W \\ j = 1, \ldots, N \end{array}$$

Empirically   rank $(J)$ can decide the efficiency
of the BPA

If J is not full rank (rank deficiency)
$\implies$ longer training times

## Hessians

$$H := \left[ \frac{\partial^2 \mathcal{E}_{av}(\underline{w})}{\partial \underline{w}^2} \right]$$

Why do we need to study Hessians?

1) Eigen values of Hessian have a role in the BPA dynamics.
   - Inverse of $H$ can provide a basis for pruning/deleting insignificant wts.

2) Lead to 2nd order opt. methods.

The Hessian of an error surface has:

(Empirical)

(a) Small # of small & large sized eigen values

(b) large # of med. sized eigen values

whose composition depend

(1) non zero mean of i/p signals & induced neural o/p signals
(non zero)

(2) Correlations between various attributes of a data vector

(3) Wide variations in the 2nd order derivatives of the J
w.r.t $\underline{W}$ from one layer to the other.

# Complexity regularization

In the context of the back prop. algo., we may want to minimize the foll. risk

$$R(\underline{w}) = \mathcal{E}_{av}(\underline{w}) + \lambda \, \mathcal{E}_c(\underline{w})$$

↗ standard performance metric

↖ Scalar

↶ Complexity part

$\mathcal{E}_c(\underline{w})$ is the complexity penalty measured in terms of $\underline{w}$

$$\mathcal{E}_c(\underline{w}) = \|\underline{w}\|^2$$

$\Big($ Choice of min $\mathcal{E}_c(\underline{w})$, force some wts. to zero & permit larger weights $\Big)$

# Alternative Strategies to learning

$$\mathcal{E}_{av}(\underline{w} + \Delta \underline{w}) = \mathcal{E}_{av}(\underline{w}) + \underline{g}^T(\underline{w}) \Delta \underline{w}$$
$$+ \frac{1}{2!} \Delta \underline{w}^T H \Delta \underline{w}$$
$$\swarrow 2^{nd} \text{ order}$$
$$+ h.o.t.$$

a small perturbation $\Delta \underline{w}$ around $\underline{w}$.

We can set
$$\frac{\partial \mathcal{E}_{av}(\underline{w} + \Delta \underline{w})}{\partial \Delta w} = 0 \implies \boxed{\Delta \underline{w} = - H^{-1} \underline{g}}$$

# Optimum brain Surgeon Algo

$$\min_{\Delta \underline{w}} \quad \frac{1}{2} \Delta \underline{w}^T \mathbb{H} \Delta \underline{w}$$

1) We set the $i$th comp. of $\Delta \underline{w}$ to zero & min. over all Synaptic wt. vectors that remain

2) Do this min. again over all indices '$i$'

# Cross validation Strategies

**Need:** We may get a low training error over a data set, but the error may be pronounced over a different data set. To get a reasonable performance, we divide the training sets into groups, train & test over the groups. "Empirical soln" in the absence of a "theoretical" set up

Eg: 1) <u>Leave 1 out strategy</u>: Train on $(N-1)$ samples & test on the other. Do over all $\binom{N}{1}$ choices

2) <u>K- fold cross validation</u>: Divide $N$ samples into $K$ equal groups, Train on $(K-1)$ groups & test on the other. Repeat over all the choices

# Convolution networks

**Motivation :** Structural layout of the MLP algo. with preprocessing steps that are neurobiologically inspired.

**Ex:** Work from Hubel & Wiesel related to specialized neurons in the visual cortex of a cat.

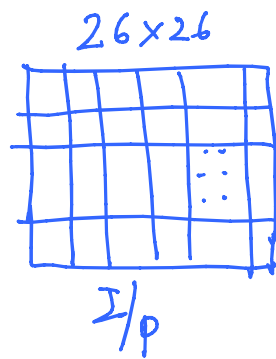Simple cells — locally sensitive

Complex cells — Orientation etc.

GOAL : To design a MLP to recognize 2D/3D shapes/objects with high invariance to translation/rotation/scaling.

1) Feature extraction : Each neuron takes its inputs from local receptive fields in previous layers, forcing to extract local features.
Once the features are extracted, its exact location is less important as long the relative position w.r.t. other features is preserved.

2) <u>Feature mapping:</u> Each computational layer is composed of multiple feature maps, with each feature map being in the form of an appropriate geometry to the signal (e.g., plane for images etc.) & constrained to share the same <u>synaptic weights</u>

a) Shift Invariance : Forced into the feature map through convolution with a kernel of a small size

b) Reduction in the # of free parameters. This is ensured via <u>wt. sharing.</u>

3) Sub Sampling: Each conv. layer performs some local operations followed by a non-linear activation $\psi(\cdot)$ & then sub sampling

$\Longrightarrow$ that reduces the __resolution__ Tolerance to sensitivity of the feature maps to shifts & distortion.

26×26

I/p

4 masks of size 3×3

4 @ 24×24
feature maps

($\psi(\cdot)$ built in)

**Pyramidal Effect !**
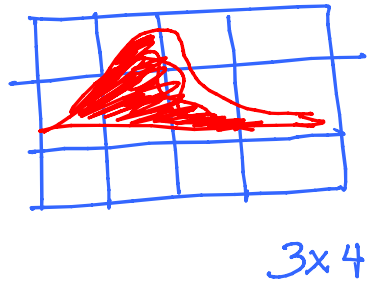
Subsampling over a 2×2 field

4@ 12×12

Conv + Subsam

Conv. followed by subsampling is inspired by _simple cells_ followed by Complex cells as described by Hubel & Wiesel.

As # of Sub-sampling ↑ Spatial resolution ↓
layers increase compared to what you
had in the prev. layer

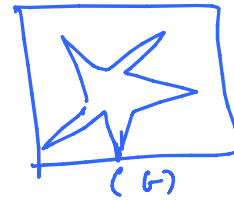Let us understand the computations through an example

1) Image:



3x4

I/p to an array of integers

| 0 | 1 | 1 | 0 |
| 1 | 1 | 1 | 1 |
| 0 | 0 | 0 | 0 |

2) $\underline{Filter}$ the image by masks / feature masks to get the features ( low pass to high pass )

More Low pass



(a)



(b)

High Pass

Types of masks :

$$\begin{array}{|c|c|} \hline 1 & 1 \\ \hline 1 & 1 \\ \hline \end{array} \cdots \quad \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 1 & 1 \\ \hline \end{array} \cdots \quad \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array}$$

2×2

Eg :

3×4 $\bigotimes$ 2×2 (Mask) filter $\longrightarrow$ 2×3 array @ o/p

$$\begin{array}{|c|c|c|c|} \hline 1 & 2 & 3 & 4 \\ \hline 5 & 6 & 7 & 8 \\ \hline 9 & 10 & 11 & 12 \\ \hline \end{array} \quad \bigotimes \quad \begin{array}{|c|c|} \hline 1 & 1 \\ \hline 1 & 1 \\ \hline \end{array} \quad \longrightarrow \quad \begin{array}{|c|c|c|} \hline 14 & 18 & \cdot \\ \hline \cdot & \cdot & \cdot \\ \hline \end{array}$$

Low Pass

3) Once the features are extracted, we feed this to a non-linear act. fn $ReLU(x) = \max(0, x)$

$ReLU(x)$ graph

4) Pooling (Subsampling)

2nd
Sub block
of size 2×2

1st
2×2 subblock

| 1 | 3 | 10 | 1 |
| 4 | 9 | 6 | 8 |
| 1 | 0 | 1 | 7 |
| 8 | 9 | 3 | 4 |

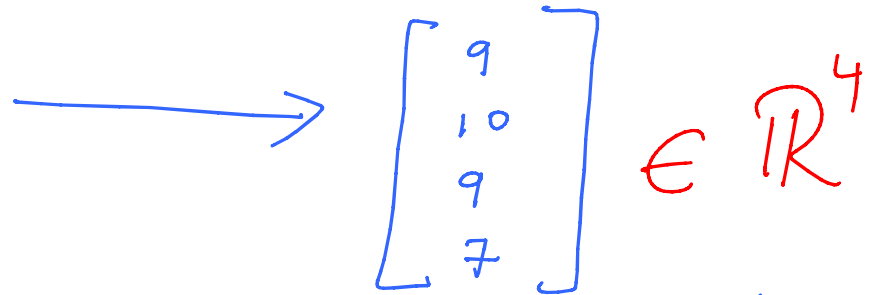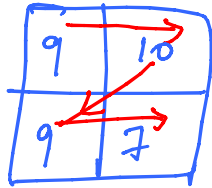Rect. Feat. Map

Max Pooling over a 2×2 sub block (Non-overlapping)

| 9 | 10 |
| 9 | 7 |

Subsampled image / map.

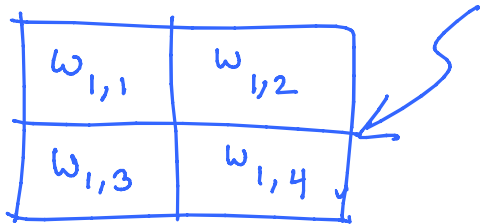5) Iterate steps (2) - (4) in a pyramidal way to get the final size as desired.

Rastering

$$\begin{bmatrix} 9 \\ 10 \\ 9 \\ 7 \end{bmatrix} \in \mathbb{R}^4$$

a feature vector

desired response → ↓ fed into the MLP

Masks can be adaptive

| $w_{1,1}$ | $w_{1,2}$ |
|---|---|
| $w_{1,3}$ | $w_{1,4}$ |

Conditions on the elements within the mask

$\Longrightarrow$ Regularized n/w for learning I/o mappings
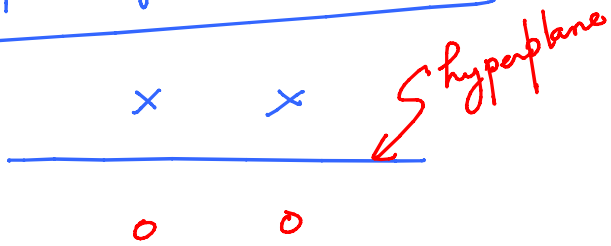
# Cover's Theorem

Ref 1: Thomas Cover, Feb 1964
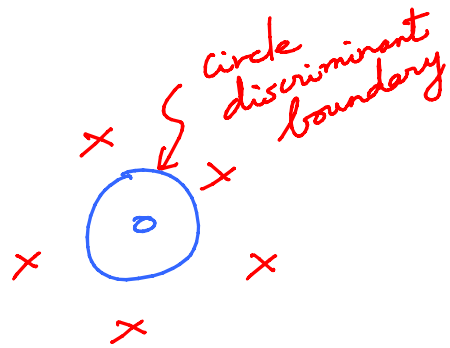
Ref 2: Nils Nilsson, Learning, Machines, 1965

## Motivation:

a) How do we quantify the complexity of a neural network architecture?

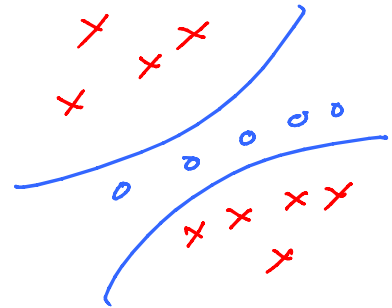b) Need a counting measure for a discrete set of mappings.

## Examples of dichotomies



(a) Linear dichotomy

(b) Spherical dichotomy

(c) Quadratic dichotomy

Consider a fixed finite set of input vectors
$\{ \underline{x}_1, \dots, \underline{x}_p \}$; $\underline{x}_i \in \mathbb{R}^N$

Can we attempt to compute the dichotomies for a perceptron?

The # of linearly realizable dichotomies on the set of points depends on a mild condition called 'general position'

General position demands that no subset of size $< N$ on $\{ \underline{x}_1, \underline{x}_2, \dots, \underline{x}_p \}$ is linearly dependent

**Theorem:** Let $\{\underline{x}_1, \dots \underline{x}_p\}$ be vectors in $\mathbb{R}^N$ that are in a general position. The # of distinct dichotomies applied to these points can be realized by a hyperplane

is
$$C(P, N) = 2 \sum_{k=0}^{N} \binom{P-1}{k}$$

**Proof:** We start with $P$ points in general position. Let us assume that there are $C(P, N)$ dichotomies on them.

↑ counting fn.

Suppose we add an extra point to this set.

We need $C(P+1, N)$ dichotomies.

**Idea:** Set up a recursion to link $C(P+1, N)$ with $C(P, N)$.

Let $(G_1, G_2, \ldots, G_P)$ be a dichotomy realizable by a hyperplane over the set of $P$ inputs.

$G_i \in \{-1, 1\}$ for every $i = 1, \ldots, P$. There is a set of weights $\underline{w}$ so that

$$\left( \text{sign} \left( \underline{w}^T \underline{x}_1 \right), \text{sign} \left( \underline{w}^T \underline{x}_2 \right), \cdots, \text{sign} \left( \underline{w}^T \underline{x}_p \right) = (b_1, b_2, \cdots, b_p)$$

Using one such $\underline{w}$, we get a dichotomy over $P+1$ points

i.e., $\left( \text{sign} \left( \underline{w}^T \underline{x}_1 \right), \cdots, \text{sign} \left( \underline{w}^T \underline{x}_p \right), \text{sign} \left( \underline{w}^T \underline{x}_{p+1} \right) \right)$

$\nearrow b_1$ $\qquad\qquad\qquad\qquad \uparrow b_p$

For every linearly realized dichotomy over $P$ points, there is at least one linearly realized dichotomy over $P+1$ points.
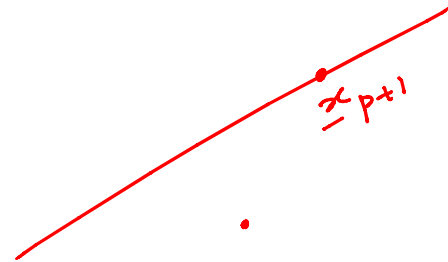
Different dichotomies over $P$ points define different dichotomies over $P+1$ points since they differ some where over first $P$ coordinates

Note that the additional dichotomy $(b_1, \ldots, b_p, -\text{sign}(\underline{w}^T \underline{x}_{p+1}))$ is also possible by reversing the sign of the last coordinate.
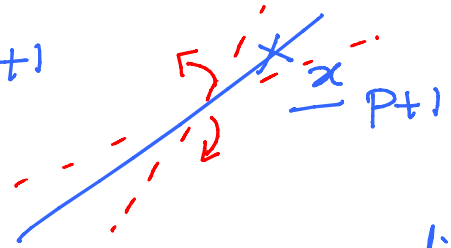
$$\implies \quad C(P+1, N) > C(P, N)$$

$$\text{Let} \quad C(P+1, N) = C(P, N) + E \quad \text{extra dichotomies possible.}$$

There are $\underline{2 \text{ cases}}$ to consider

$\underline{x}_{p+1}$

## Case A:

One of the weight vectors $\underline{w}$ that generates $(b_1, b_2, \ldots, b_p)$ passes through $\underline{x}_{p+1}$



By adjusting the angle of the hyperplane, we can adjust $\text{sign}\left(\underline{w}^T \underline{x}_{p+1}\right)$ to $+1$ or $-1$ i.e., both $(b_1 \ldots b_p +1)$ and $(b_1 \ldots b_p -1)$ are also possible!

**Case B:** No hyperplane passes through $\underline{x}_{p+1}$ and generates $b_1, \dots, b_p$ on first $p$ vectors.

$\Rightarrow$ Point lies on one side of the (old) dichotomy

$E$ is the # of dichotomies over $p$ points that are realized by a hyperplane passing through a fixed point $\underline{x}_{p+1}$. By forcing the hyperplane to pass through $\underline{x}_{p+1}$, we are going to $N-1$ dimensions instead of $N$

$\underline{x}_{x_{p+1}}$

Geometrically, if a point is on the x-axis, the hyperplane has $N-1$ axes left to work on the problems. If it is not on the x-axis, then rotate the axes of the space to get the point on the x-axis and there is no effect on the geometry of the problems

$$\therefore \quad E = C(P, N-1)$$

$$\therefore \quad C(P+1, N) = C(P, N) + C(P, N-1) \qquad \textcircled{1}$$

Let us consider the **boundary conditions**

$$C(1, N) = 2 \checkmark$$

(There is 1 point in $\mathbb{R}^N$ & can be realized by 2 labels)

$$C(P, 1) = 2P \checkmark$$

(There are P points in $\mathbb{R}^1$)

Consider $P = 3$, we have

$O - o$
$x - i$

```
o  o  o /
o  o / x
o / x / o
o / x  x
```

```
x  o  o
x  o  x
x  x  o
x  x  x
```

Observe that except
$\underline{o \, x \, o}$ and $\underline{x \, o \, x}$, we
can have a hyperplane
that can shatter the points!

Let us prove the result through induction.

(We are doing induction over 'P')

Base case :
Case 1 : $C(1, N) = 2$ as expected

There is $1$ point in $N$ dim. Follows from one of the boundary conditions

Induction : $C(P+1, N) = 2 \sum_{k=0}^{N} \binom{P-1}{k} + 2 \sum_{k=0}^{N-1} \binom{P-1}{k}$

$\underbrace{\qquad}_{C(P, N)}$ $\underbrace{\qquad}_{C(P, N-1)}$

Ponder on $\binom{P-1}{0}$ Case

Meaning of '0' dim.

Given by statement of Cover's theorem

Now, simplifying,

$$= 2 \sum_{k=0}^{N} \binom{P-1}{k} + 2 \sum_{k=0}^{N} \binom{P-1}{k-1} \quad \left( \because \binom{P-1}{-1} = 0 \right)$$

$$= 2 \left[ \sum_{k=0}^{N} \left[ \binom{P-1}{k} + \binom{P-1}{k-1} \right] \right]$$

$$= 2 \sum_{k=0}^{N} \binom{P}{k} \qquad \left( \because \binom{n}{k} = \binom{n-1}{k} + \binom{n-1}{k-1} \right.$$
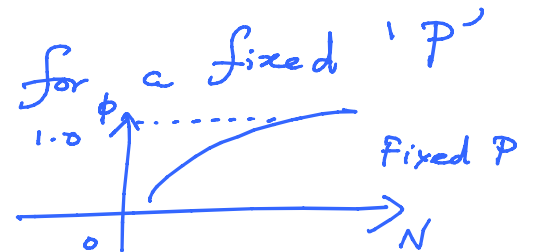$$\left. \text{Basic identity from elementary Combinatorics} \right)$$

# Implications:

(linear decision boundary)

Let us consider the prob. of having a perceptron that can provide a linear dichotomy over $P$ points

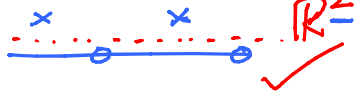$$p = \frac{C(P, N)}{2^P} \leftarrow \text{total \# of dichotomies } \left(\begin{array}{l}\text{not necessarily}\\ \text{linear}\end{array}\right)$$

$$= 2^{1-P} \sum_{k=0}^{N} \binom{P-1}{k}$$

**Home Work:** Plot $p$ vs. dim $N$ for a fixed '$P$'

Observe the __Concavity__ of $p$

$\mathbb{R}^1$

(Linear decision boundary is not possible)

$\mathbb{R}^2$ ✓ Yes, a linear decision boundary is possible

Let us revisit the XOR problem

---

Recall:     Cover's theorem has 2 important Consequences

1) The use of a non-linear function i.e., a hidden function defined by $\varphi$ $\underline{(\underline{x})}$ $\longleftarrow$ acts on the input vector

2) High dimensionality of the hidden / feature space compared to the input space
$N :=$ dim. feature space
$M := $ dim. input space $\Bigg\}$     $N > M$

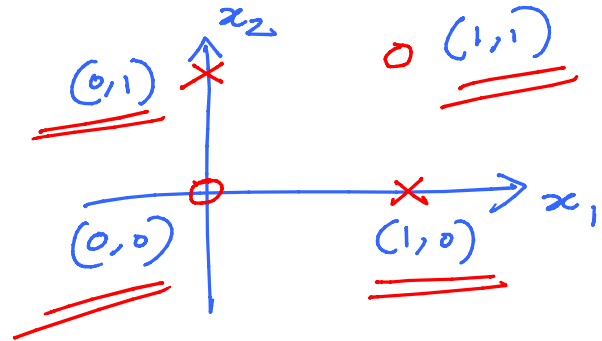Recall: A dichotomy is $\phi-$ separable if $\exists$ a
N-dim. vector $\underline{w}$ / non-linear fn.

$$\underline{w}^T \phi(\underline{x}) > 0 \qquad \underline{x} \in \text{Class 1}$$

$$\underline{w}^T \phi(\underline{x}) \leq 0 \qquad \underline{x} \in \text{Class 2}$$
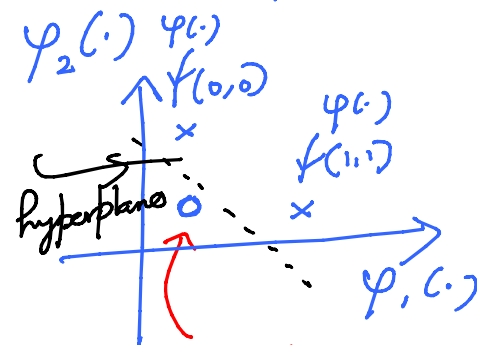
We have the XOR problem

Let us consider the pair of Gaussian hidden functions.

$$\varphi_1(\underline{x}) = \exp\left(-\|\underline{x} - \underline{t}_1\|^2\right); \quad \underline{t}_1 = [1, 1]^T$$

$$\varphi_2(\underline{x}) = \exp\left(-\|\underline{x} - \underline{t}_2\|^2\right); \quad \underline{t}_2 = [0, 0]^T$$

Let us tabulate the comp. evaluations of $\underline{x}$ over $\varphi_1$ & $\varphi_2$

| $x$ | $\varphi_1(\underline{x})$ | $\varphi_2(\underline{x})$ |
|---|---|---|
| (1, 1) | 1 | 0.1353 |
| (0, 1) | 0.3678 | 0.3678 |
| (0, 0) | 0.1353 | 1 |
| (1, 0) | 0.3678 | 0.3678 |



Both points (0,1) & (1,0) map to the same point.