

# ON THE STABILITY OF DEEP NETWORKS

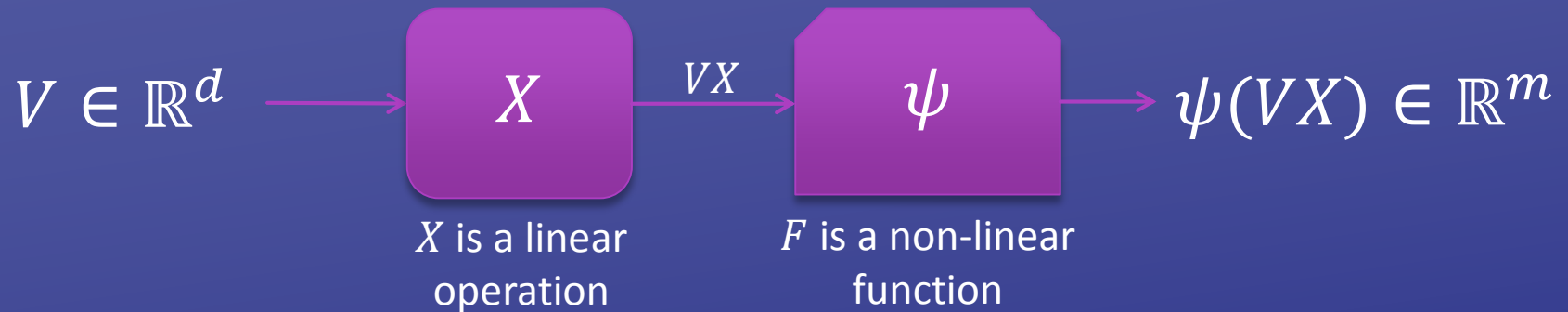
AND THEIR RELATIONSHIP TO COMPRESSED  
SENSING AND METRIC LEARNING

**RAJA GIRYES AND GUILLERMO SAPIRO**  
**DUKE UNIVERSITY**

Mathematics of Deep Learning  
International Conference on Computer Vision (ICCV)  
December 12, 2015

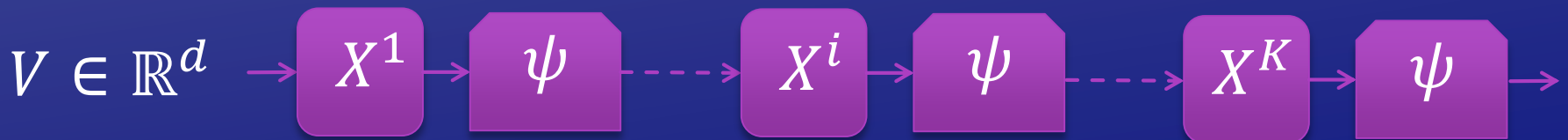
# DEEP NEURAL NETWORKS (DNN)

- One layer of a neural net



- Concatenation of the layers creates the whole net

$$\Phi(X^1, X^2, \dots, X^K) = \psi(\psi(\psi(VX^1)X^2) \dots X^K)$$



# THE NON-LINEAR PART

- Usually  $\psi = g \circ f$ .



- $f$  is the (point-wise) activation function

ReLU  
 $f(x) = \max(x, 0)$

A red line graph of the ReLU function, which is zero for negative x and increases linearly for positive x.

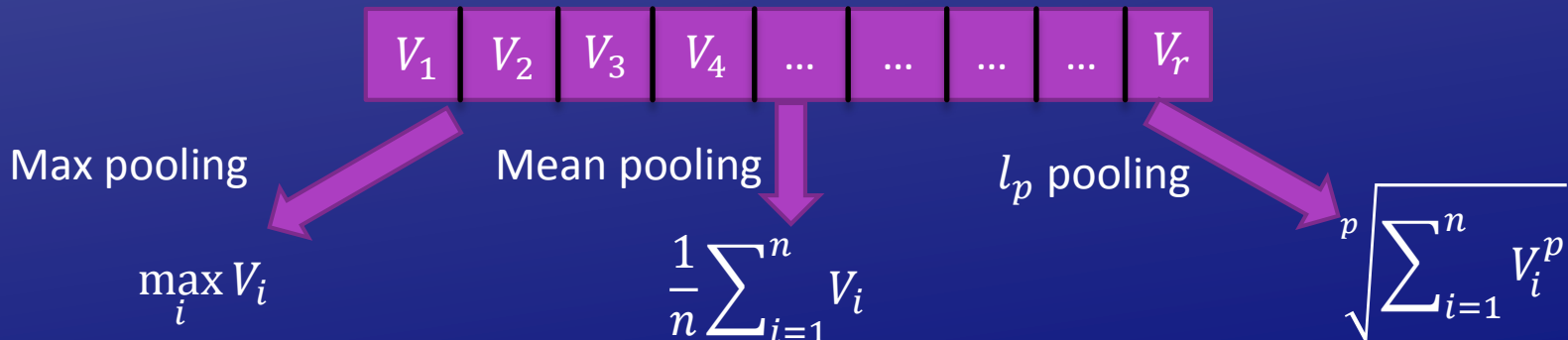
Sigmoid  
 $f(x) = \frac{1}{1 + e^{-x}}$

A red S-shaped curve graph of the Sigmoid function, ranging from 0 to 1.

Hyperbolic tangent  
 $f(x) = \tanh(x)$

A red S-shaped curve graph of the Hyperbolic tangent function, ranging from -1 to 1.

- $g$  is a pooling or an aggregation operator.



# WHY DNN WORK?

What is so special with the DNN structure?

What is the role of the depth of DNN?

Why a local training of each layer is a good choice?

How many training samples do we need?

What is the role of pooling?

What is the role of the activation function?

What happens to the data throughout the layers?

# SAMPLE OF RELATED EXISTING THEORY

- Universal approximation for any measurable Borel functions [Hornik et. al., 1989, Cybenko 1989].
- Depth of a network provides an exponential complexity compared to the number parameters [Montúfar. al. 2014]
- Pooling stage provides shift invariance [Bruna et. al. 2013]
- Relation of pooling and phase retrieval [Bruna et. al. 2014].
- Deeper networks have more local minima that are close to the global one and less saddle points [Saxe et. al. 2014], [Dauphin et. al. 2014] [Choromanska et. al. 2015] [Haeffele & Vidal, 2015]

DNN keep  
the  
important  
information  
of the data.

Gaussian mean  
width is a good  
measure for the  
complexity of  
the data.

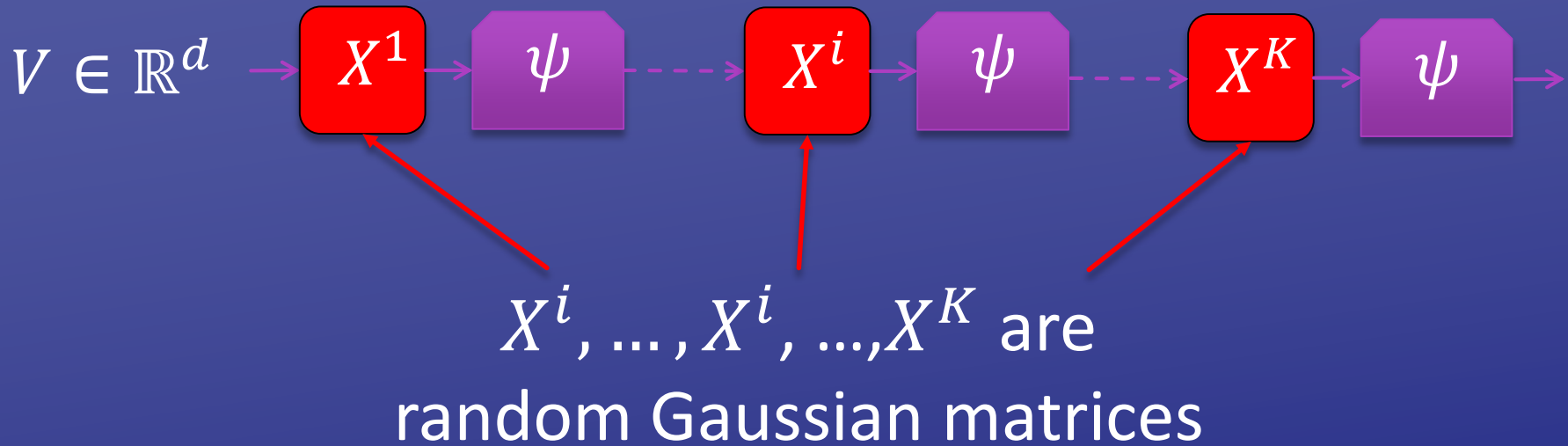
## Take Home Message

Important goal  
of training:  
Classify the  
boundary points  
between the  
different classes  
in the data.

Random  
Gaussian  
weights are  
good for  
classifying the  
average points  
in the data.

Deep learning  
can be viewed  
as a metric  
learning.

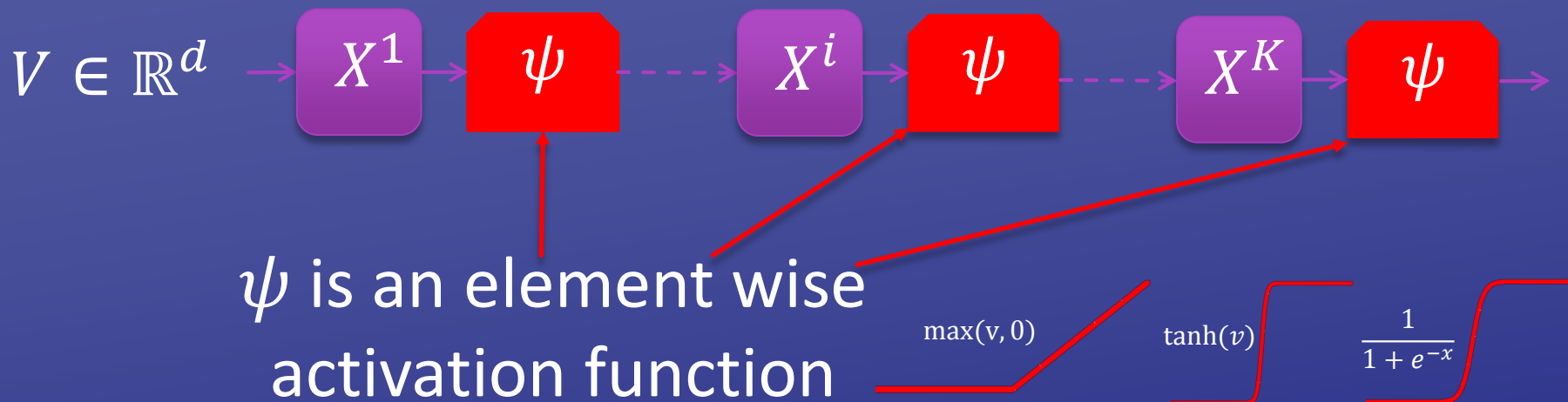
# ASSUMPTIONS – GAUSSIAN WEIGHTS



- Infusion of random weights reveals internal properties of a system



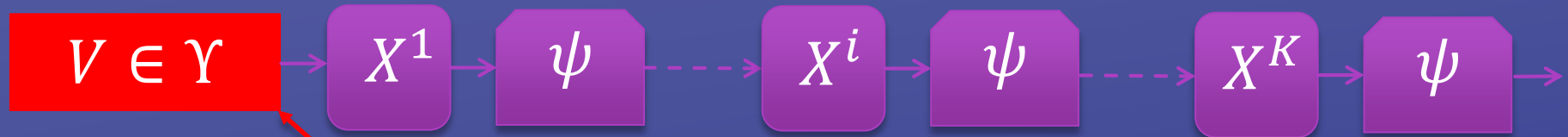
# ASSUMPTIONS – NO POOLING



- Pooling provides invariance [Boureau et. al. 2010, Bruna et. al. 2013].
- We assume that all equivalent points in the data were merged together and omit this stage.
- Reveals the role of the other components in the DNN.



# ASSUMPTIONS – LOW DIMENSIONAL DATA



$\Upsilon$  is a low dimensional set

Gaussian  
Mixture  
Models  
(GMM)

Low Rank  
Matrices

Signals with  
Sparse  
Representations

Low  
Dimensional  
Manifolds

DNN keep  
the  
important  
information  
of the data.

Gaussian mean  
width is a good  
measure for the  
complexity of  
the data.

# Gaussian Mean Width

Important goal  
of training:  
Classify the  
boundary points  
between the  
different classes  
in the data.

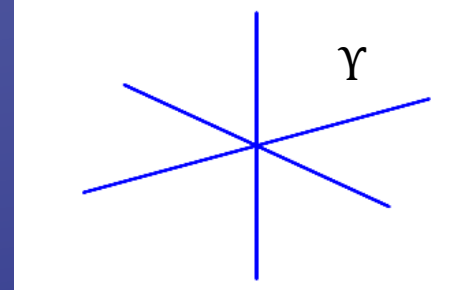
Random  
Gaussian  
weights are  
good for  
classifying the  
average points  
in the data.

Deep learning  
can be viewed  
as a metric  
learning.

# WHAT HAPPENS TO SPARSE DATA IN DNN?

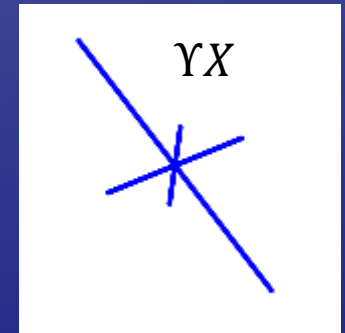
- Let  $\Upsilon$  be sparsely represented data

- Example:  $\Upsilon = \{V \in \mathbb{R}^3: \|V\|_0 \leq 1\}$

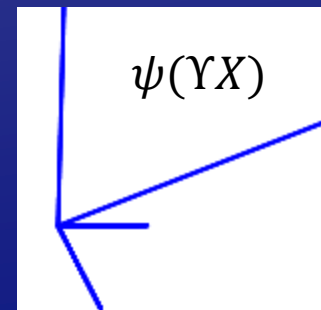


- $\Upsilon X$  is still sparsely represented data

- Example:  $\Upsilon X = \{V \in \mathbb{R}^3: \exists W \in \mathbb{R}^3, V = XW, \|W\|_0 \leq 1\}$



- $\psi(\Upsilon X)$  not sparsely represented
- But is still low dimensional

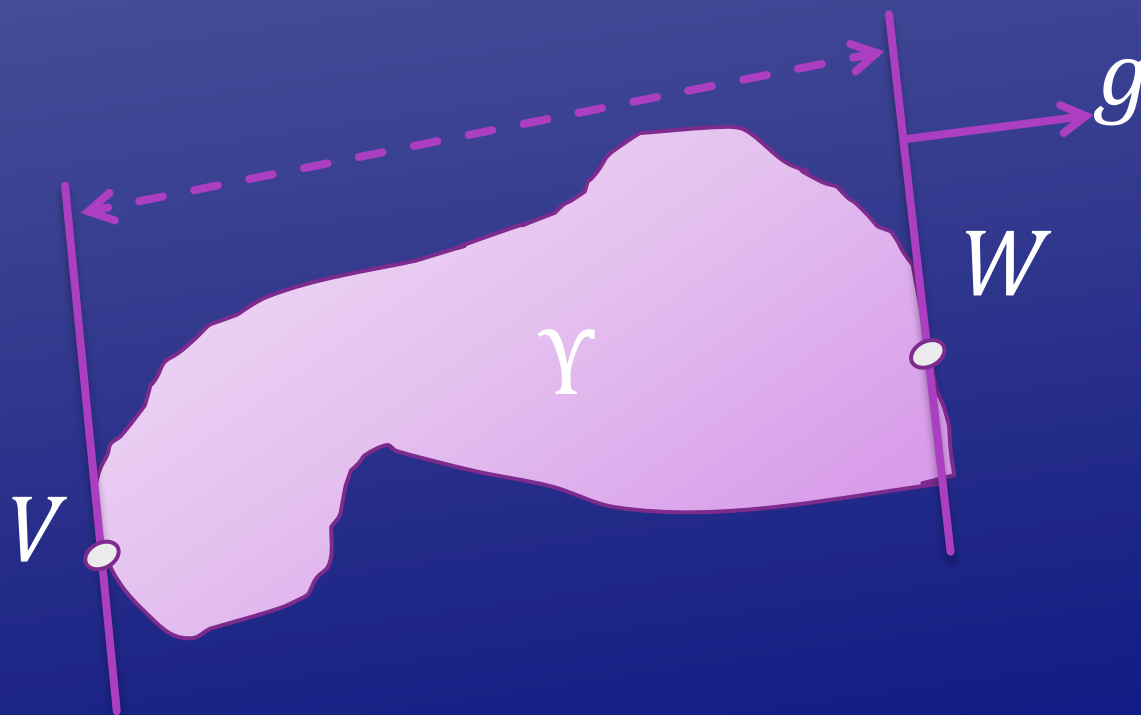


# GAUSSIAN MEAN WIDTH

- Gaussian mean width:

$$\omega(Y) = E \sup_{V, W \in Y} \langle V - W, g \rangle, \quad g \sim N(0, I).$$

The width of  
the set  $Y$  in  
the direction  
of  $g$ :



# MEASURE FOR LOW DIMENSIONALITY

- Gaussian mean width:

$$\omega(Y) = E \sup_{V, W \in Y} \langle V - W, g \rangle, \quad g \sim N(0, I).$$

- $\omega^2(Y)$  is a measure for the dimensionality of the data.
- Examples:

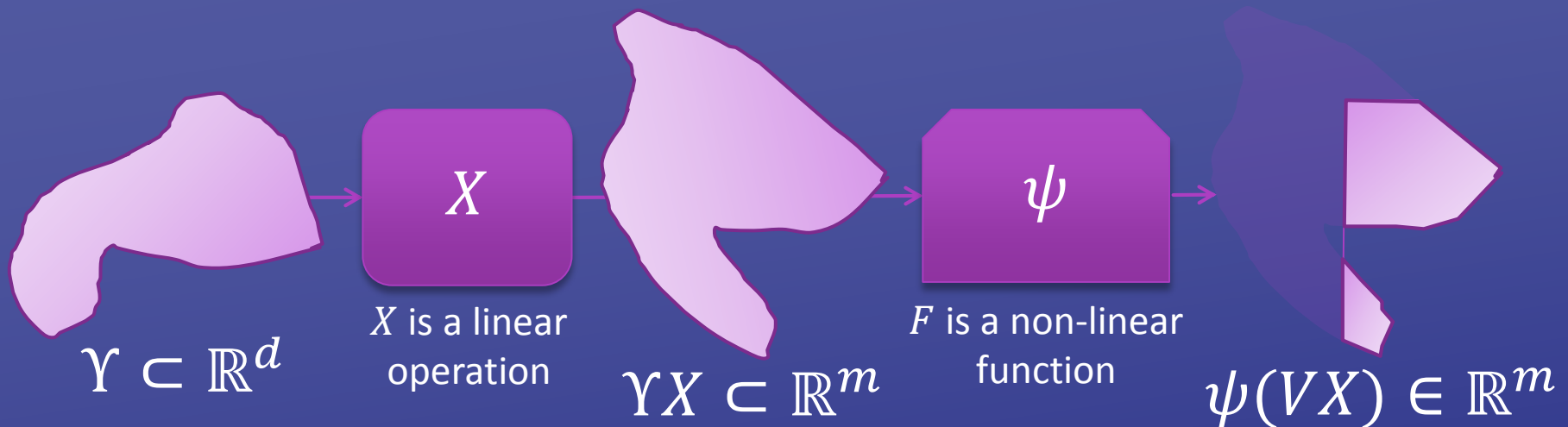
If  $Y \subset \mathbb{B}^d$  is a Gaussian Mixture Model with  $k$  Gaussians then

$$\omega^2(Y) = O(k)$$

If  $Y \subset \mathbb{B}^d$  is a data with  $k$ -sparse representations then

$$\omega^2(Y) = O(k \log d)$$

# GAUSSIAN MEAN WIDTH IN DNN



Theorem 1: small  $\frac{\omega^2(\Upsilon)}{m}$  imply  $\omega^2(\Upsilon) \approx \omega^2(\psi(VX))$

Small  $\omega^2(\Upsilon)$



Small  $\omega^2(\psi(VX))$



It is sufficient to provide proofs only for a single layer

DNN keep  
the  
important  
information  
of the data.

Gaussian mean  
width is a good  
measure for the  
complexity of  
the data.

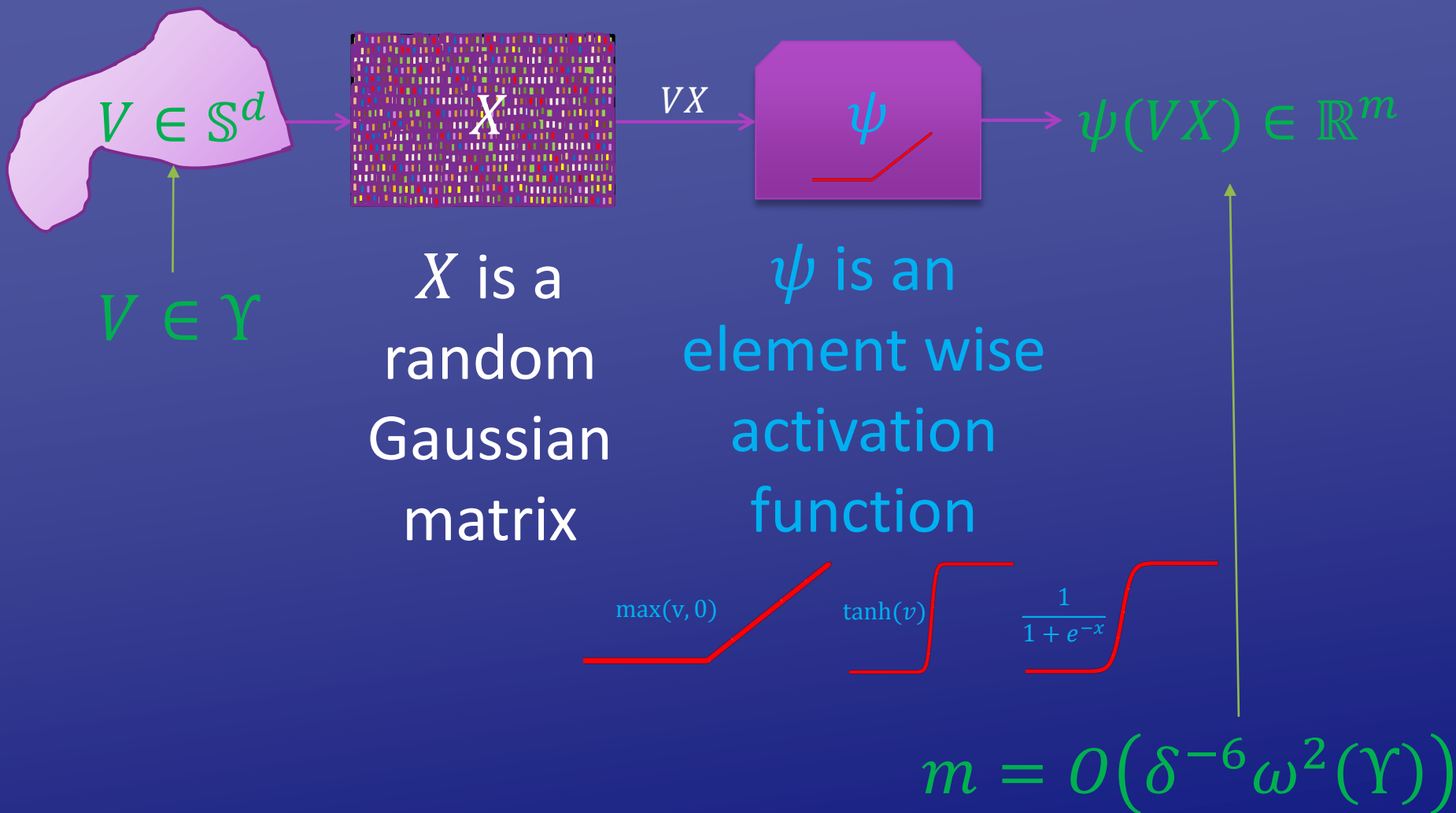
Important goal  
of training:  
Classify the  
boundary points  
between the  
different classes  
in the data.

# Stability

Random  
Gaussian  
weights are  
good for  
classifying the  
average points  
in the data.

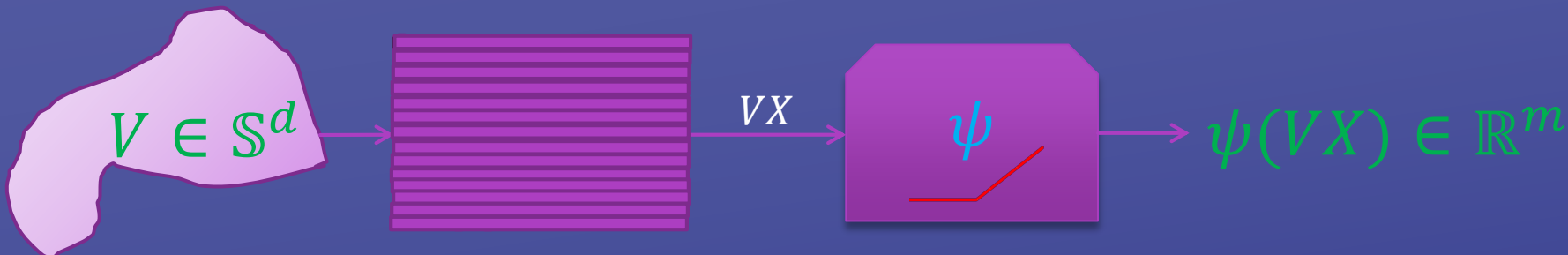
Deep learning  
can be viewed  
as a metric  
learning.

# ASSUMPTIONS

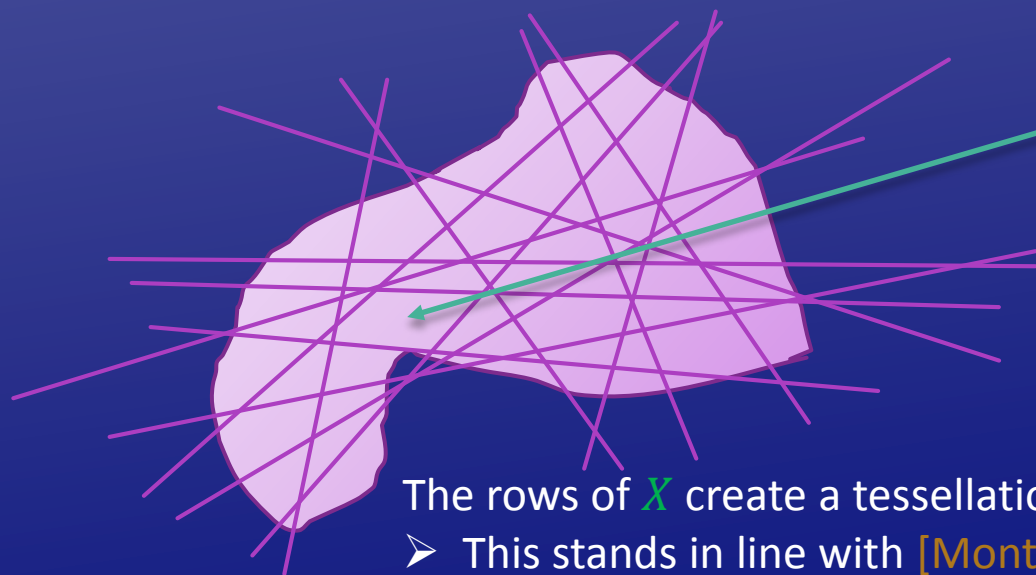




# ISOMETRY IN A SINGLE LAYER



Theorem 2:  $\psi(\cdot X)$  is a  $\delta$ -isometry in the Gromov-Hausdorff sense between the sphere  $\mathbb{S}^{d-1}$  and the Hamming cube [Plan & Vershynin, 2014, Giryes, Sapiro & Bronstein 2015].

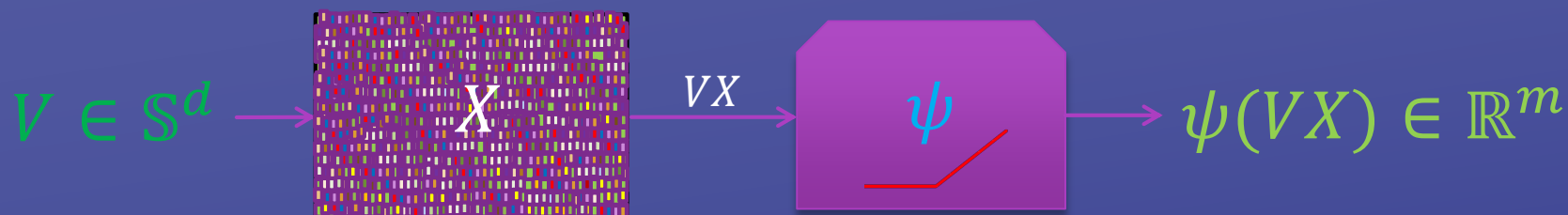


- If two points belong to the same tile then their distance  $< \delta$
- ➔ Each layer of the network keeps the main information of the data

The rows of  $X$  create a tessellation of the space.

➤ This stands in line with [Montúfar et. al. 2014]

# ONE LAYER STABLE EMBEDDING



Theorem 3: There exists an algorithm  $\mathcal{A}$  such that

$$\|V - \mathcal{A}(\psi(VX))\| < O\left(\frac{\omega(Y)}{\sqrt{m}}\right) = O(\delta^3)$$

[Plan & Vershynin, 2013, Giryes, Sapiro & Bronstein 2015].

- After  $K$  layers we have an error  $O(K\delta^3)$
- Stands in line with [Mahendran and Vedaldi, 2015].
- DNN keep the important information of the data

DNN keep  
the  
important  
information  
of the data.

Gaussian mean  
width is a good  
measure for the  
complexity of  
the data.

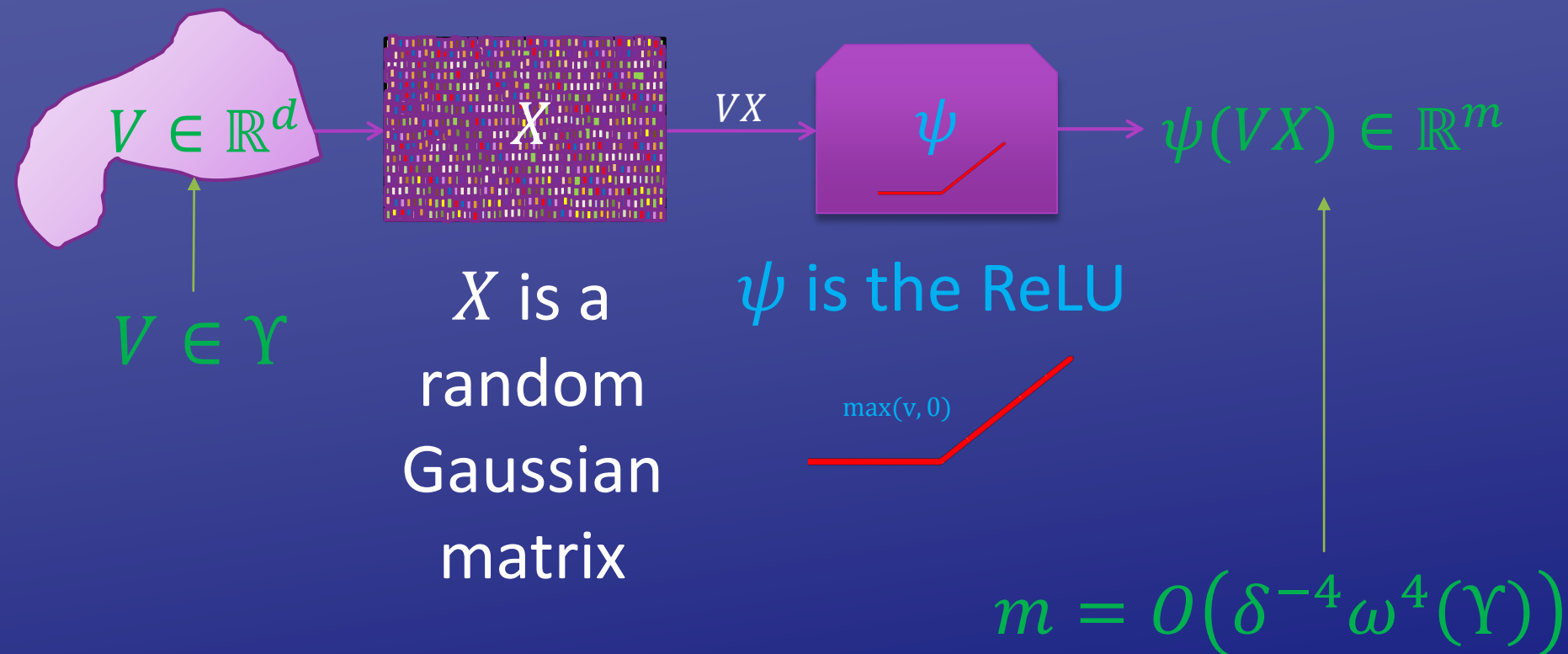
# DNN with Gaussian Weights

Important goal  
of training:  
Classify the  
boundary points  
between the  
different classes  
in the data.

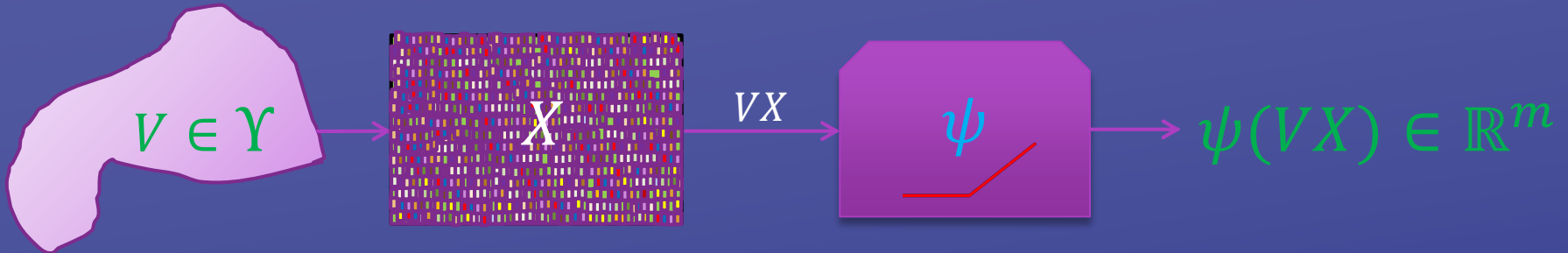
Random  
Gaussian  
weights are  
good for  
classifying the  
average points  
in the data.

Deep learning  
can be viewed  
as a metric  
learning.

# ASSUMPTIONS



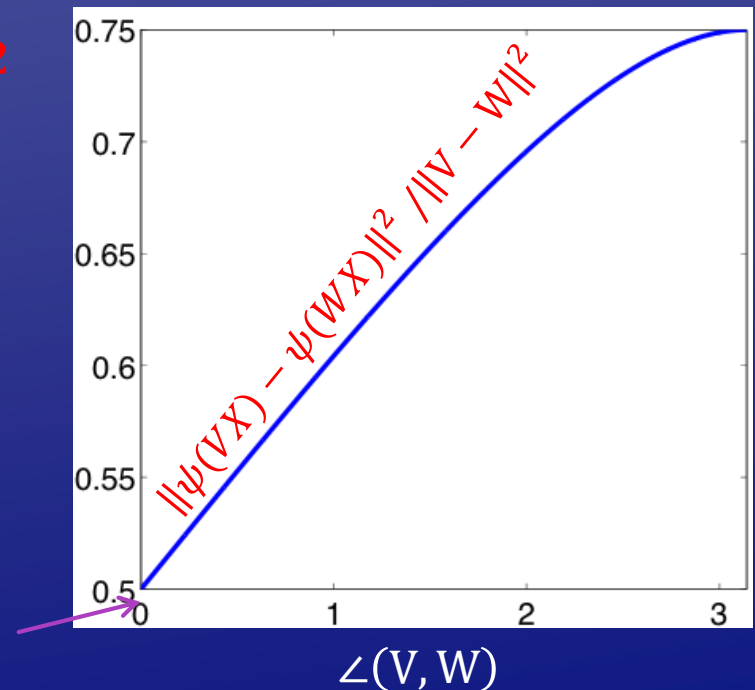
# DISTANCE DISTORTION



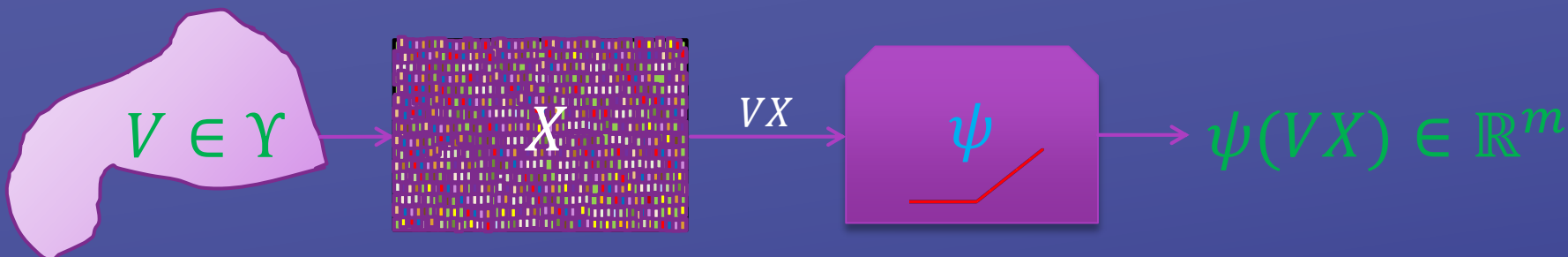
Theorem 4: for  $V, W \in \mathcal{Y}$

$$\left| \|\psi(VX) - \psi(WX)\|^2 - \frac{1}{2}\|V - W\|^2 - \frac{\|V\|\|W\|}{\pi} (\sin \angle(V, W)) \right|$$

The smaller  $\angle(V, W)$  the smaller the distance we get between the points



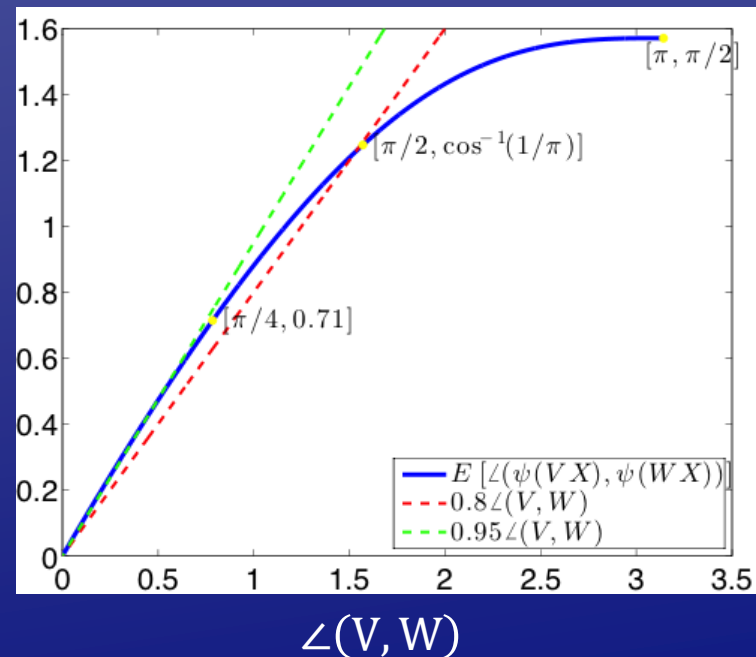
# ANGLE DISTORTION



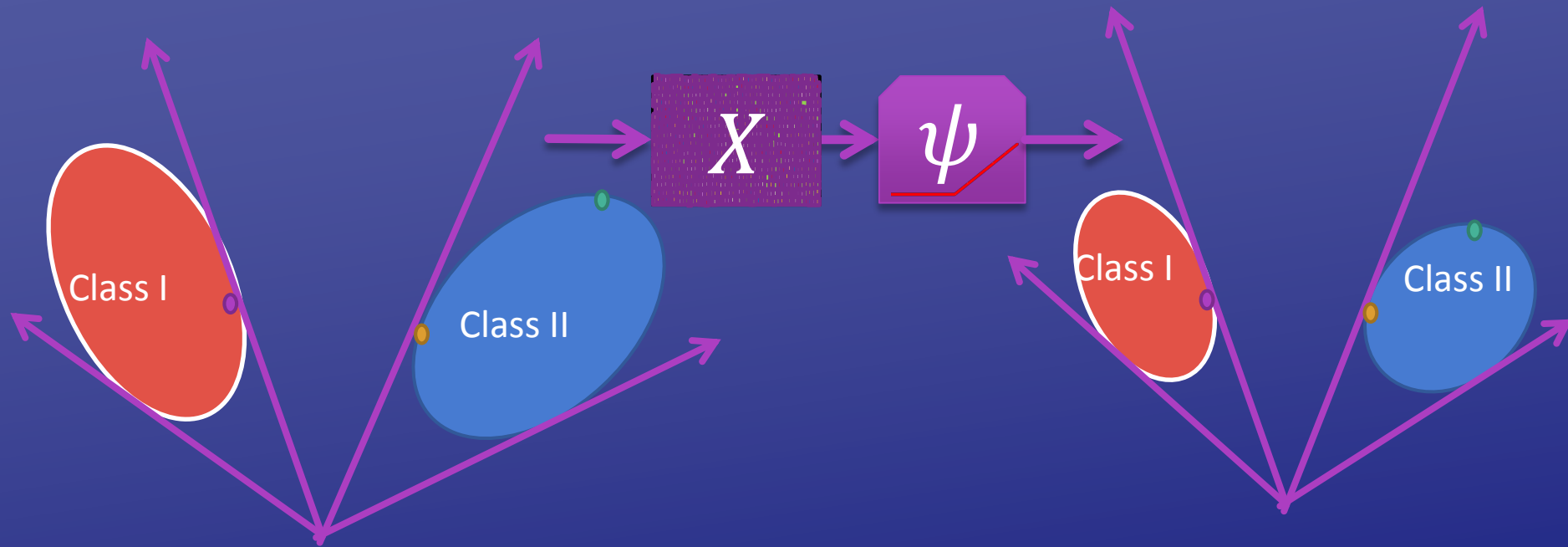
Theorem 5: for  $V, W \in \Upsilon$

$$\left| \cos \angle(\psi(VX), \psi(WX)) - \cos \angle(V, W) \right| - \frac{1}{\pi} (\sin \angle(V, W))$$

Behavior of  $\angle(\psi(VX), \psi(WX))$



# DISTANCE AND ANGLES DISTORTION



Points with small angles between them become closer than points with larger angles between them

# TRAINING DATA SIZE

- Stability in the network implies that close points in the input are close also at the output
- ➡ Having a good network for an  $\varepsilon$ -net of the input set  $Y$  guarantees a good network for all the points in  $Y$ .
- ➡ Using Sudakov minoration the number of data points is

$$\exp(\omega^2(Y)/\varepsilon^2).$$

- ➡ Though this is not a tight bound, it introduces the Gaussian mean width  $\omega(Y)$  as a measure for the complexity of the input data and the required number of training samples.



# POOLING AND CONVOLUTIONS

- We test empirically this behavior on convolutional neural networks (CNN) with random weights and the MNIST, CIFAR-10 and ImageNet datasets.
- The behavior predicted in the theorem remains also in the presence of pooling and convolutions.

# Role of Training

DNN keep the important information of the data.

Gaussian mean width is a good measure for the complexity of the data.

Important goal of training: Classify the boundary points between the different classes in the data.

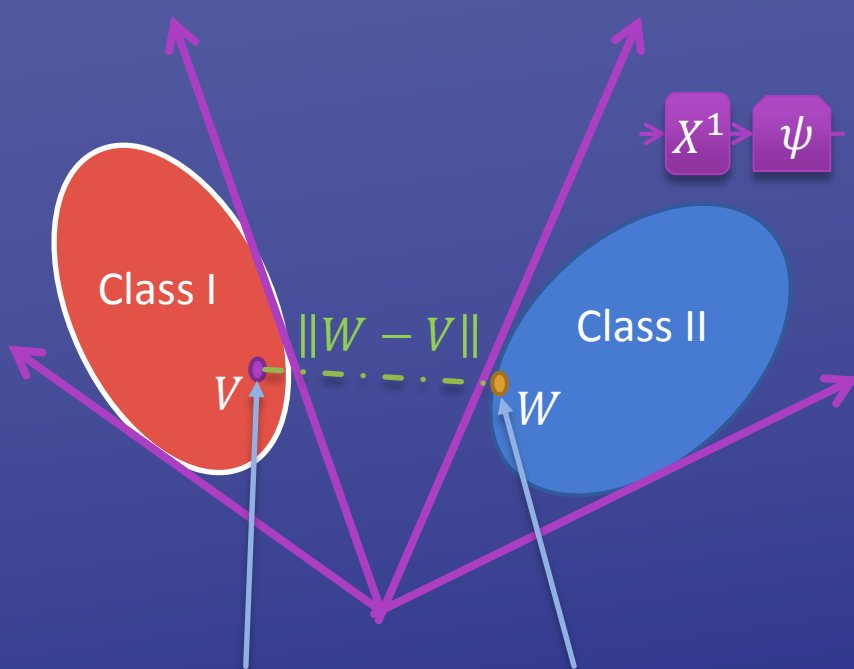
Random Gaussian weights are good for classifying the average points in the data.

Deep learning can be viewed as a metric learning.

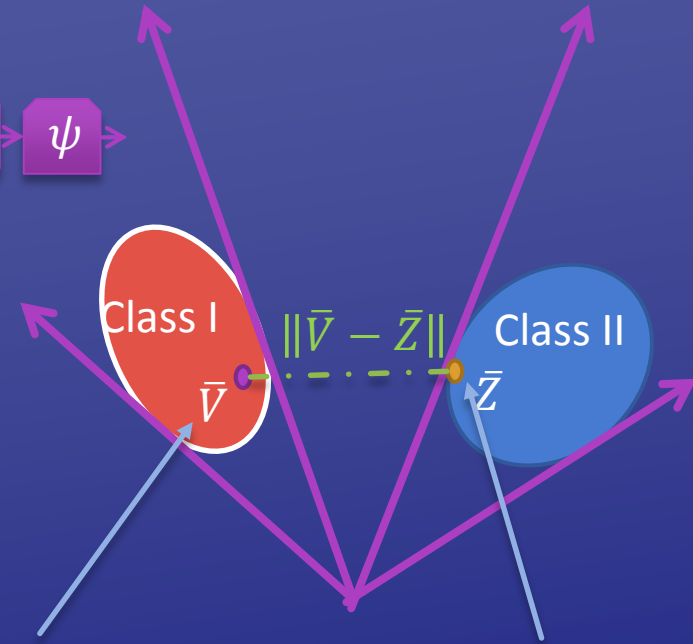
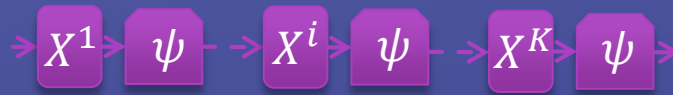
# ROLE OF TRAINING

- Having a theory for Gaussian weights we test the behavior of DNN after training.
- We looked at the MNIST, CIFAR-10 and ImageNet datasets.
- We will present here only the ImageNet results.
- We use a state-of-the-art pre-trained network for ImageNet [Simonyan & Zisserman, 2014].
- We compute inter and intra class distances.

# INTER BOUNDARY POINTS DISTANCE RATIO



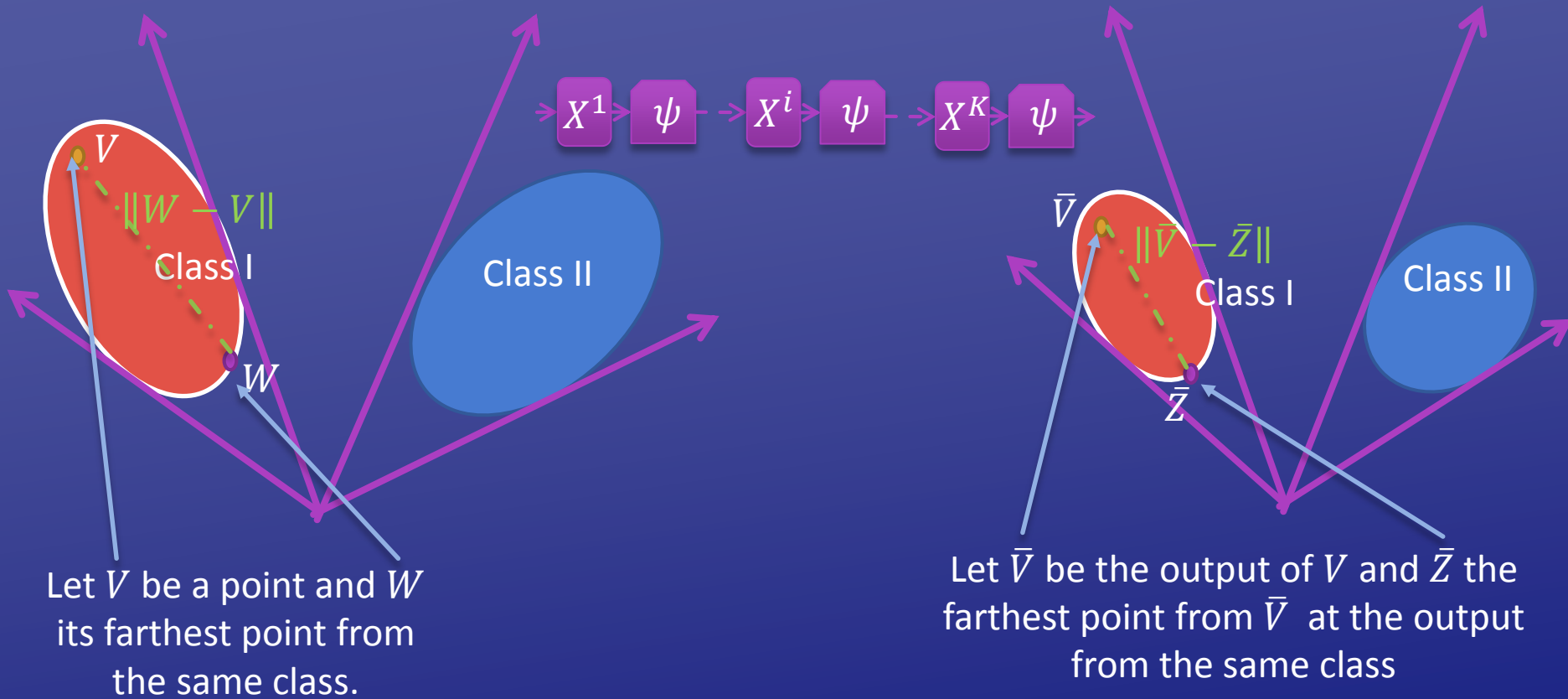
$V$  is a random point and  $W$  its closest point from a different class.



$\bar{V}$  is the output of  $V$  and  $\bar{Z}$  the closest point to  $\bar{V}$  at the output from a different class.

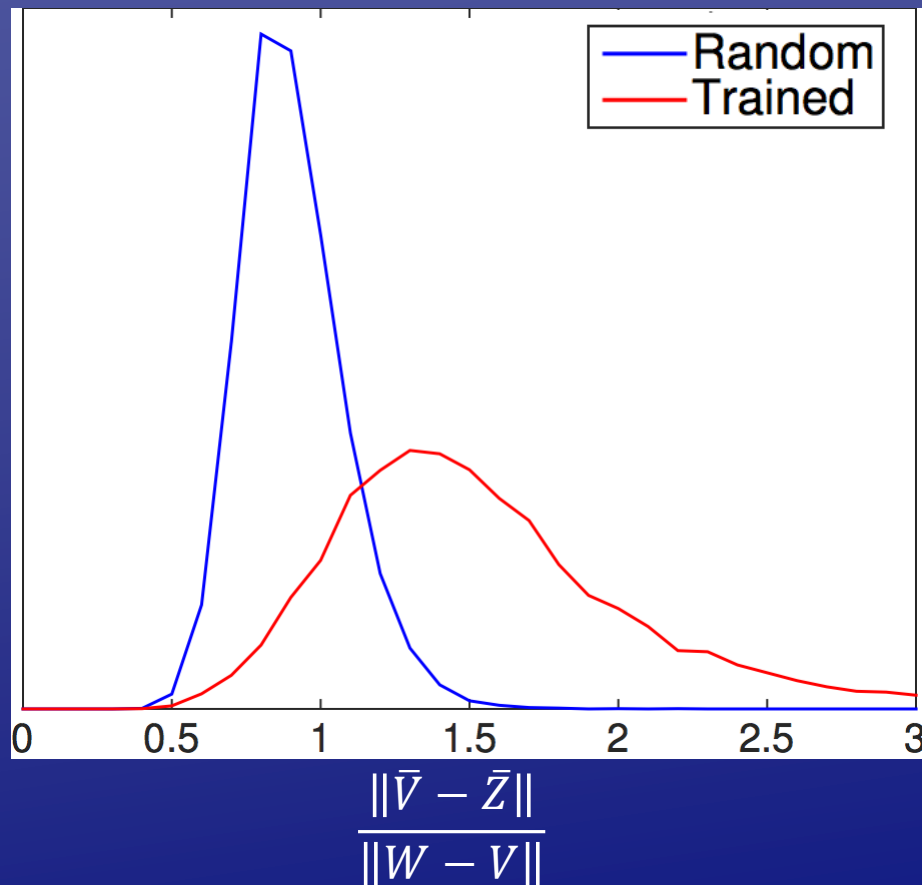
Compute the distance ratio:  $\frac{\|\bar{V} - \bar{Z}\|}{\|W - V\|}$

# INTRA BOUNDARY POINTS DISTANCE RATIO

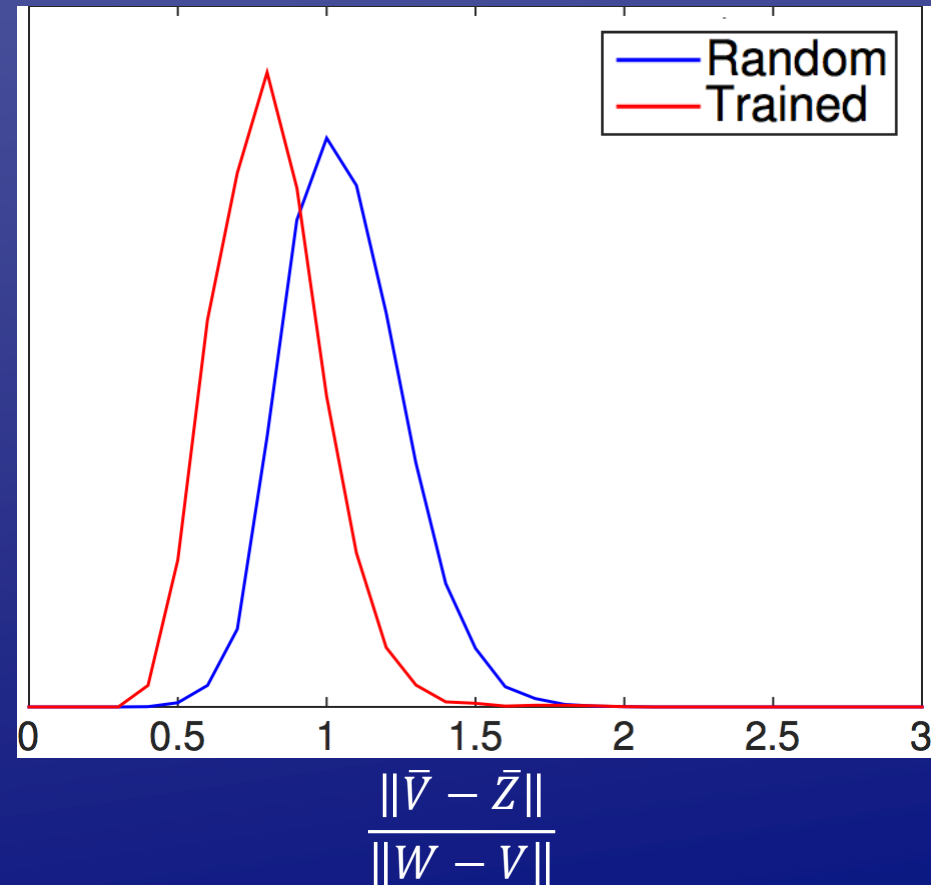


# BOUNDARY DISTANCE RATIO

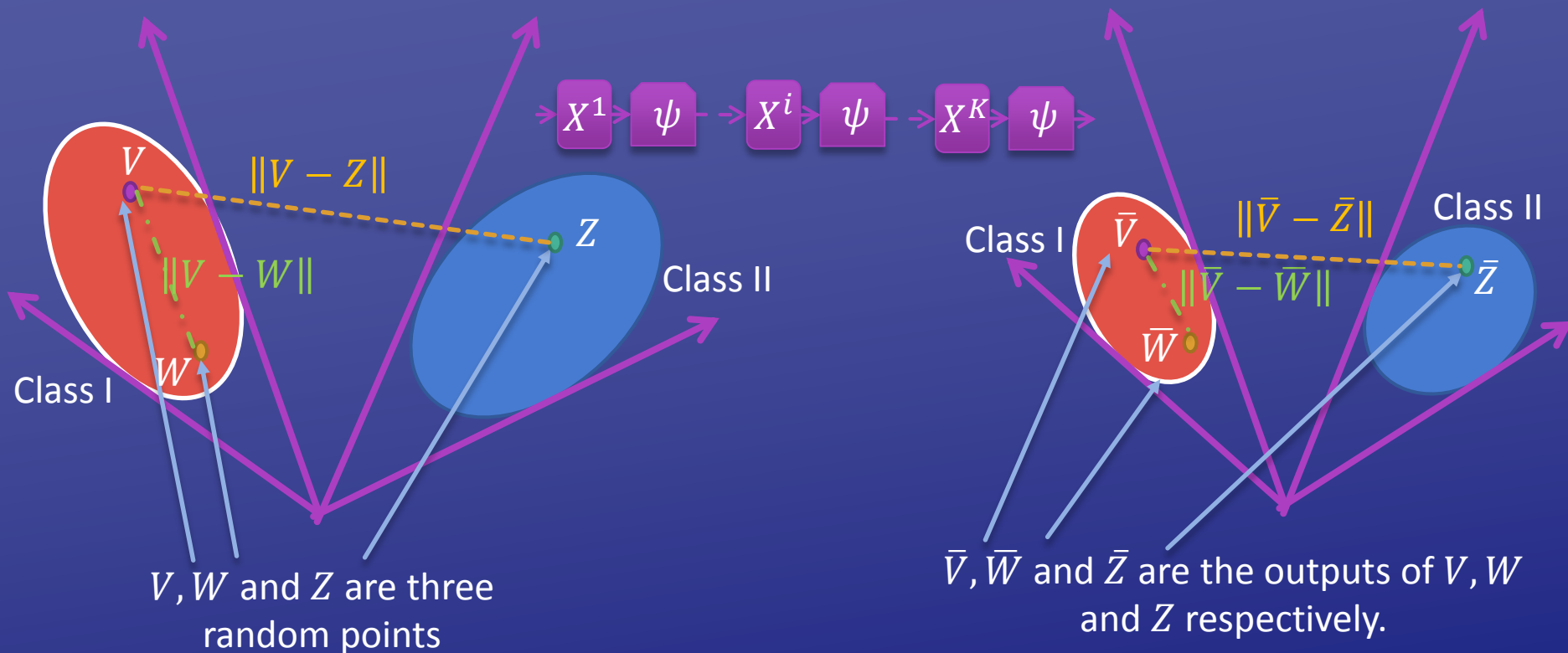
Inter-class



Intra-class



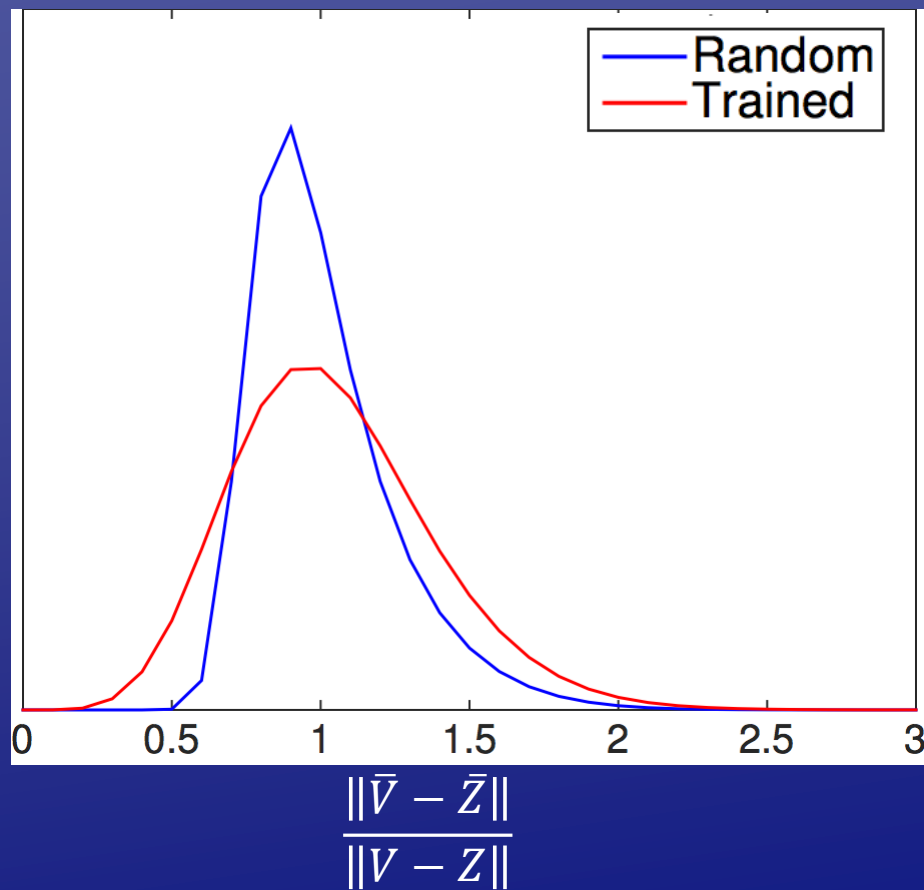
# AVERAGE POINTS DISTANCE RATIO



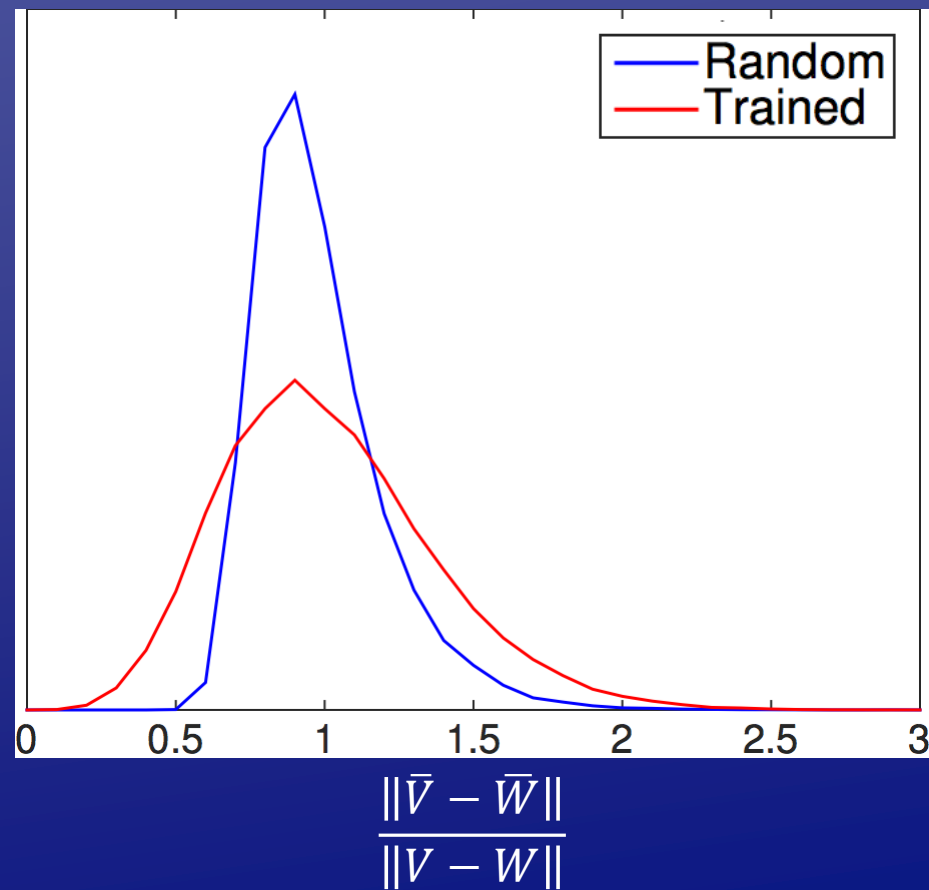
Compute the distance ratios:  $\frac{\|\bar{V} - \bar{W}\|}{\|V - W\|}, \frac{\|\bar{V} - \bar{Z}\|}{\|V - Z\|}$

# AVERAGE DISTANCE RATIO

Inter-class



Intra-class





# ROLE OF TRAINING

- On average distances are preserved in the trained and random networks.
- The difference is with respect to the boundary points.
- The inter distances become larger.
- The intra distances shrink.

DNN keep  
the  
important  
information  
of the data.

Gaussian mean  
width is a good  
measure for the  
complexity of  
the data.

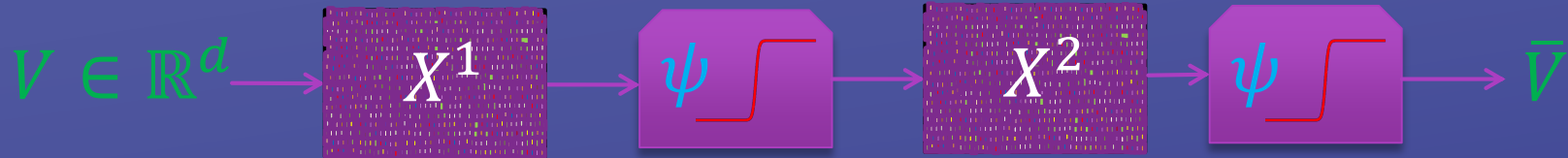
Important goal  
of training:  
Classify the  
boundary points  
between the  
different classes  
in the data.

## DNN as Metric Learning

Random  
Gaussian  
weights are  
good for  
classifying the  
average points  
in the data.

Deep learning  
can be viewed  
as a metric  
learning.

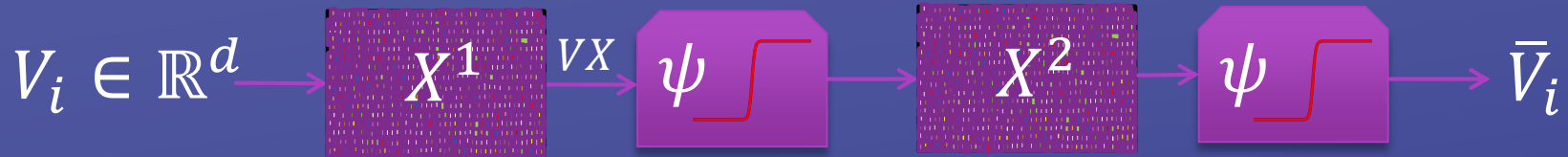
# ASSUMPTIONS



$X$  is fully  
connected  
and trained

$\psi$  is the  
hyperbolic tan

# METRIC LEARNING BASED TRAINING



- Cosine Objective:

$$\min_{X^1, X^2} \sum_{i,j \in \text{Training Set}} \left( \frac{\bar{V}_i^T \bar{V}_j}{\|\bar{V}_i\| \|\bar{V}_j\|} - \vartheta_{i,j} \right)^2$$

$$\vartheta_{i,j} = \begin{cases} \lambda + (1 - \lambda) \frac{V_i^T V_j}{\|V_i\| \|V_j\|} & i, j \in \text{same class} \\ -1 & i, j \in \text{different class} \end{cases}$$

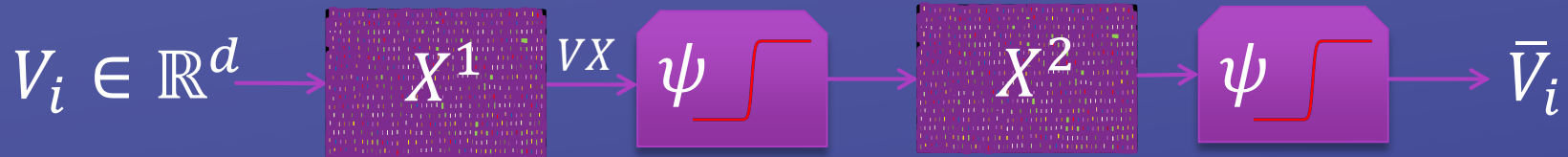
Classification term

Metric preservation term

$i, j \in \text{same class}$

$i, j \in \text{different class}$

# METRIC LEARNING BASED TRAINING



$$l_{ij} = \begin{cases} 1 & i, j \in \text{same class} \\ -1 & i, j \in \text{different class} \end{cases} \quad l_{ij} = \begin{cases} \text{average intra class distance} & i, j \in \text{same class} \\ \text{average inter class distance} & i, j \in \text{different class} \end{cases}$$

- Euclidean Objective:

$$\min_{X^1, X^2} \frac{\lambda}{|Training Set|} \sum_{i, j \in Training Set} [l_{ij} \|\bar{V}_i - \bar{V}_j\| - t_{ij}]_+ + \frac{1-\lambda}{|Neighbours|} \sum_{V_i, V_j \text{ are neighbours}} \left| \|\bar{V}_i - \bar{V}_j\| - \|V_i - V_j\| \right|$$

Classification term

Metric learning term

# ROBUSTNESS OF THIS NETWORK

- Metric learning objectives impose stability
- Similar to what we have in the random case
- Close points at the input are close at the output
- Using the theory of  $(T, \epsilon)$ -robustness, the generalization error scales as  $\sqrt{\frac{T}{|Training\ set|}}$
- $T$  is the covering number.
- Also here, the number of training samples scales as  $\exp(\omega^2(\Upsilon)/\epsilon^2)$ .

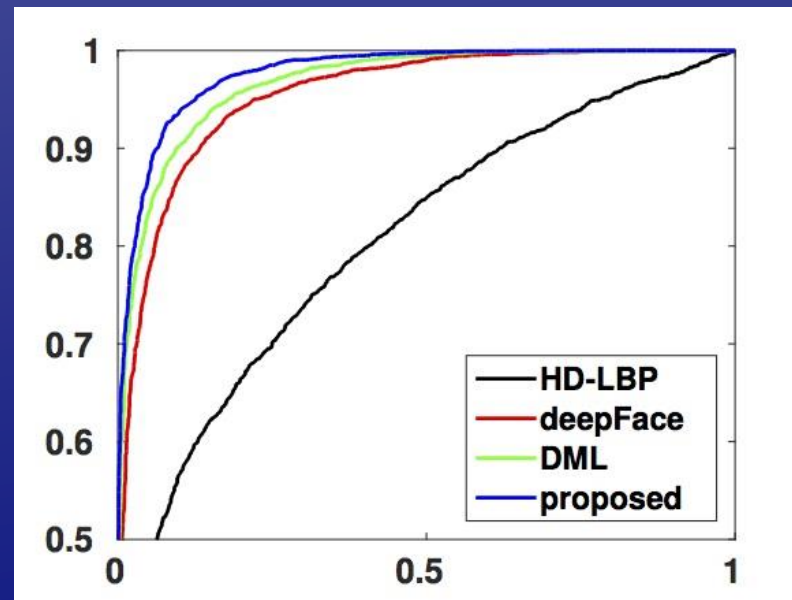
# RESULTS

- Better performance with less training samples

MNIST  
Dataset

#Training/class	30	50	70	100
original pixels	81.91%	86.18%	86.86%	88.49%
LeNet	87.51%	89.89%	91.24%	92.75%
Proposed 1	92.32%	94.45%	95.67%	96.19%
Proposed 2	<b>94.14%</b>	<b>95.20%</b>	<b>96.05%</b>	<b>96.21%</b>

Faces in  
the wild  
ROC curve



[Huang, Qiu, Sapiro,  
Calderbank, 2015]

DNN keep  
the  
important  
information  
of the data.

Gaussian mean  
width is a good  
measure for the  
complexity of  
the data.

## Take Home Message

Important goal  
of training:  
Classify the  
boundary points  
between the  
different classes  
in the data.

Random  
Gaussian  
weights are  
good for  
classifying the  
average points  
in the data.

Deep learning  
can be viewed  
as a metric  
learning.



# ACKNOWLEDGEMENTS



Alex Bronstein  
Tel Aviv University



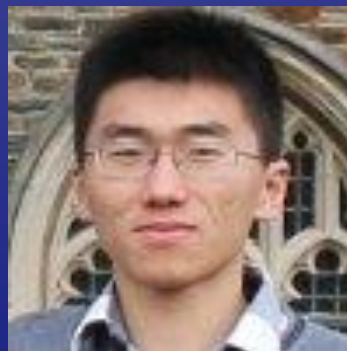
Guillermo Sapiro  
Duke University



Robert Calderbank  
Duke University



Qiang Qiu  
Duke University



Jiaji Huang  
Duke University

## Grants

**NSF**

**ONR**

**NGA**

**NSSEFF**

**ARO**

**ERC**

# QUESTIONS?

[SITES.DUKE.EDU/RAJA](http://SITES.DUKE.EDU/RAJA)

## REFERENCES

R. GIRYES, G. SAPIRO, A. M. BRONSTEIN, *DEEP NEURAL NETWORKS WITH RANDOM GAUSSIAN WEIGHTS: A UNIVERSAL CLASSIFICATION STRATEGY?*

J. HUANG, Q. QIU, G. SAPIRO, R. CALDERBANK, *DISCRIMINATIVE GEOMETRY-AWARE DEEP TRANSFORM*

J. HUANG, Q. QIU, G. SAPIRO, R. CALDERBANK, *DISCRIMINATIVE ROBUST TRANSFORMATION LEARNING*