

Deep Neural Networks: A Universal Classification Strategy?

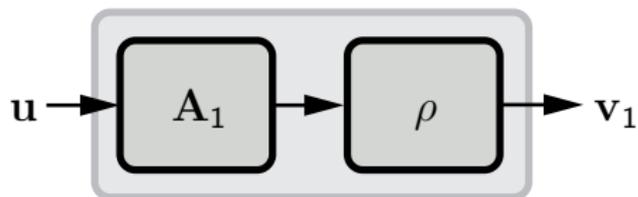
Alex Bronstein

School of Electrical Engineering
Tel Aviv University

February 3, 2016

Based on joint work with Raja Giryes and Guillermo Sapiro

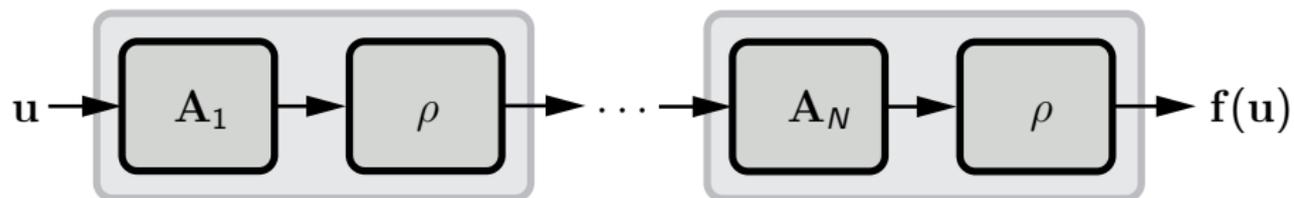
(Deep) neural network



Single layer:

$$v_1 = \rho(A_1 u)$$

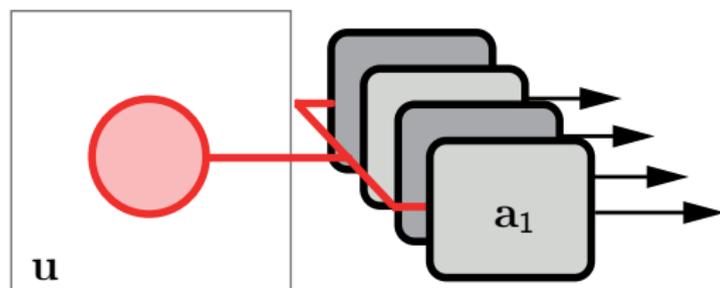
(Deep) neural network



Whole net response:

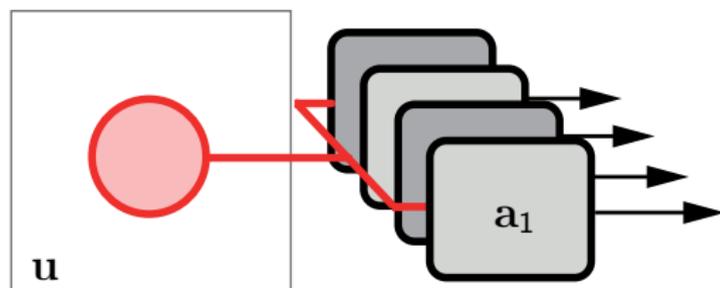
$$\mathbf{f}(\mathbf{u}) = \rho(\mathbf{A}_N \rho(\mathbf{A}_{N-1} \rho(\dots \rho(\mathbf{A}_1 \mathbf{u}) \dots)))$$

Linear part



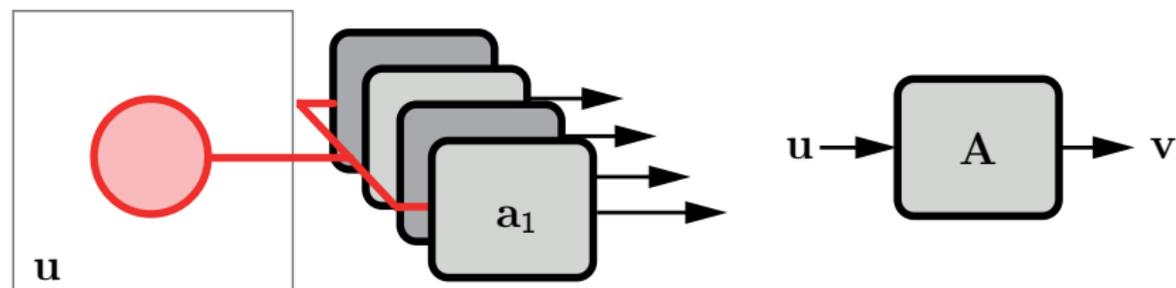
- Convolutional layer: shift-invariant filter bank $\mathbf{v}_i = \mathbf{a}_i * \mathbf{u}_i$

Linear part



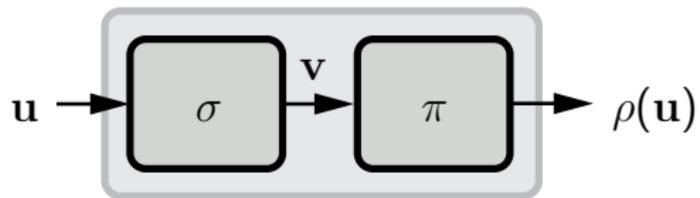
- Convolutional layer: shift-invariant filter bank $\mathbf{v}_i = \mathbf{a}_i * \mathbf{u}$;
 \mathbf{A} is block-Töplitz

Linear part

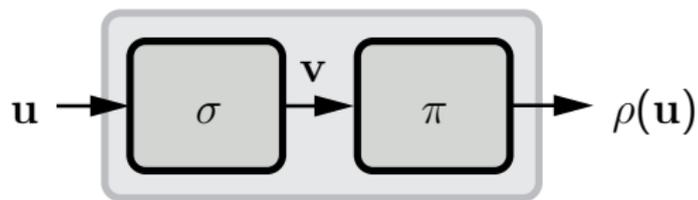


- Convolutional layer: shift-invariant filter bank $\mathbf{v}_i = \mathbf{a}_i * \mathbf{u}$
 \mathbf{A} is block-Töplitz
- Fully-connected layer: $\mathbf{v} = \mathbf{A}\mathbf{u}$

Non-linear part



Non-linear part



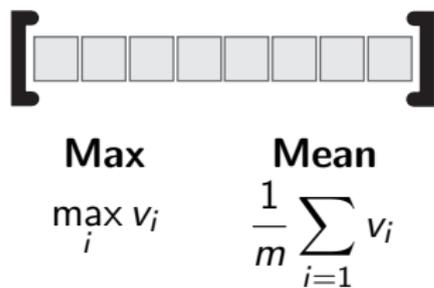
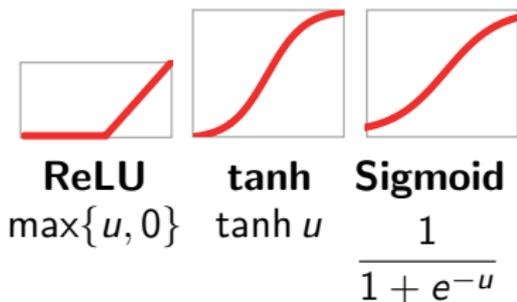
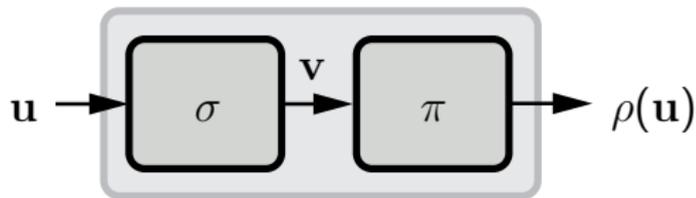
ReLU
 $\max\{u, 0\}$

tanh
 $\tanh u$

Sigmoid
 $\frac{1}{1 + e^{-u}}$

- Element-wise activation function $\sigma(u)$

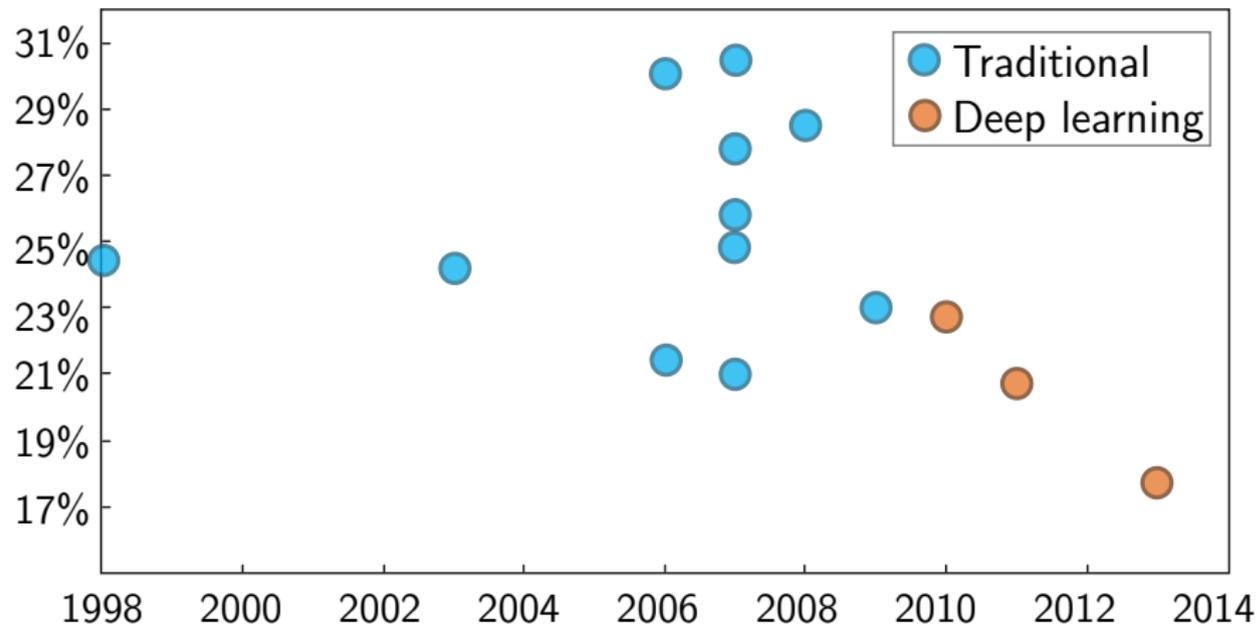
Non-linear part



- Element-wise activation function $\sigma(u)$
- Pooling or aggregation operator $\pi(v)$

Impact of deep learning

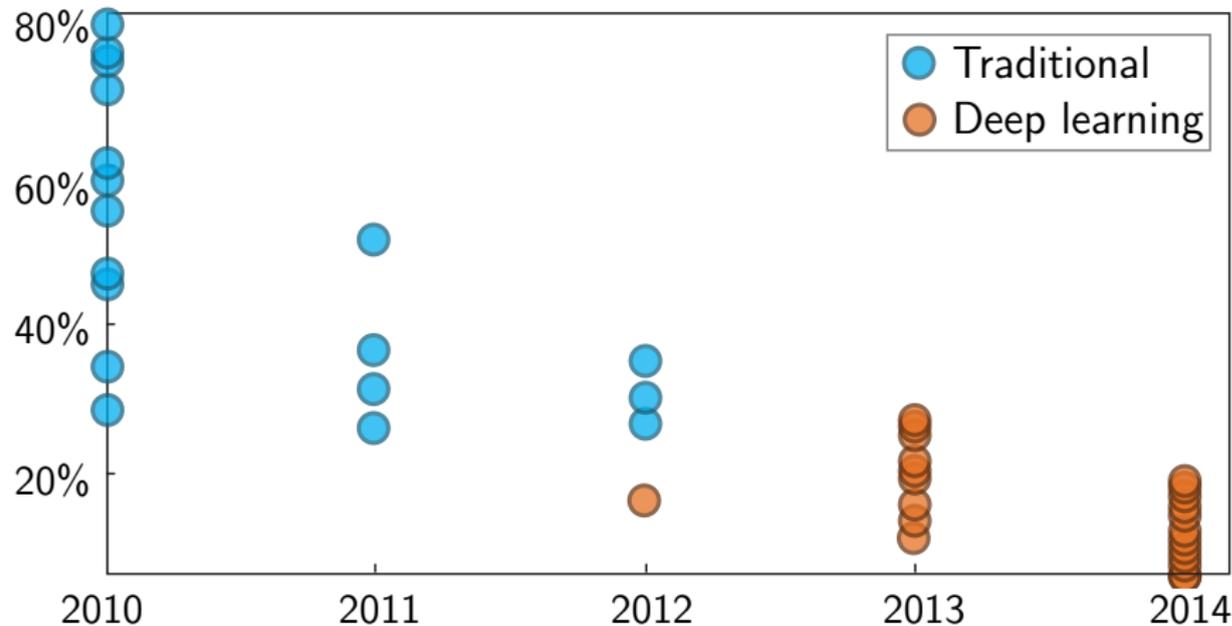
Audio recognition error rates



Source: Clarifi

Impact of deep learning

Visual recognition error rates



Source: Clarifi

Why DNNs work?

Why DNNs work?

- Representation power?

Why DNNs work?

- Representation power?
- Role of depth?

Why DNNs work?

- Representation power?
- Role of **depth**?
- Role of **pooling**?

Why DNNs work?

- Representation power?
- Role of **depth**?
- Role of **pooling**?
- Role of **nonlinearity**?

Why DNNs work?

- Representation power?
- Role of depth?
- Role of pooling?
- Role of nonlinearity?
- How to train?

Why DNNs work?

- Representation power?
- Role of depth?
- Role of pooling?
- Role of nonlinearity?
- How to train?
- How much training data are needed?

Representation power

- DNNs are **universal approximators** of any Borel function¹

¹Cybenko 1989; Hornik 1991; ²Barron 1992

Representation power

- DNNs are **universal approximators** of any Borel function¹
- **Estimation error** of a function f by DNN is²

$$\mathcal{O}\left(\frac{C_f}{K}\right) + \mathcal{O}\left(\frac{nK}{T} \log T\right)$$

C_f = smoothness of f

K = # of degrees of freedom

n = input dimension

T = # of training samples

¹Cybenko 1989; Hornik 1991; ²Barron 1992

Representation power

- DNNs represent **restricted Boltzmann machines** with number of parameters **exponentially greater** than the number of degrees of freedom of the network¹

¹Montúfar & Morton, 2014; ²Montúfar *et al.*, 2014

Representation power

- DNNs represent **restricted Boltzmann machines** with number of parameters **exponentially greater** than the number of degrees of freedom of the network¹
- Deep network with the same number of degrees of freedom divides the input space into **exponentially greater** number of sets²

¹Montúfar & Morton, 2014; ²Montúfar *et al.*, 2014

Representation power

- DNNs represent **restricted Boltzmann machines** with number of parameters **exponentially greater** than the number of degrees of freedom of the network¹
- Deep network with the same number of degrees of freedom divides the input space into **exponentially greater** number of sets²
- **Depth is important!**

¹Montúfar & Morton, 2014; ²Montúfar *et al.*, 2014

Role of pooling

- Pooling provides **shift invariance**¹

¹Bruna, LeCun & Szeliski, 2013,2014; ²Bruna & Mallat, 2013

Role of pooling

- Pooling provides **shift invariance**¹
- **Scattering networks**² : a cascade of wavelet transform, modulus and averaging

¹Bruna, LeCun & Szlam, 2013,2014; ²Bruna & Mallat, 2013

Role of pooling

- Pooling provides **shift invariance**¹
- **Scattering networks**² : a cascade of wavelet transform, modulus and averaging
- Deeper network provides invariance to **more complex transformations**

¹Bruna, LeCun & Szlam, 2013,2014; ²Bruna & Mallat, 2013

Role of pooling

- Pooling provides **shift invariance**¹
- **Scattering networks**² : a cascade of wavelet transform, modulus and averaging
- Deeper network provides invariance to **more complex transformations**
- **Depth is important!**

¹Bruna, LeCun & Szlam, 2013,2014; ²Bruna & Mallat, 2013

Training

- Supervised training by **back propagation** (chain rule) results in **non-convex** optimization problem

Saxe, McClelland & Ganguli, 2014; Dauphin *et al.*, 2014; Choromanska *et al.*, 2015; Haeffele & Vidal, 2015

Training

- Supervised training by **back propagation** (chain rule) results in **non-convex** optimization problem
- In deep networks, **local minima** are nearly as good as global ones

Saxe, McClelland & Ganguli, 2014; Dauphin *et al.*, 2014; Choromanska *et al.*, 2015; Haeffele & Vidal, 2015

Training

- Supervised training by **back propagation** (chain rule) results in **non-convex** optimization problem
- In deep networks, **local minima** are nearly as good as global ones
- Deep networks have less **saddle points**

Saxe, McClelland & Ganguli, 2014; Dauphin *et al.*, 2014; Choromanska *et al.*, 2015; Haefele & Vidal, 2015

Training

- Supervised training by **back propagation** (chain rule) results in **non-convex** optimization problem
- In deep networks, **local minima** are nearly as good as global ones
- Deep networks have less **saddle points**
- **Random initialization** works well

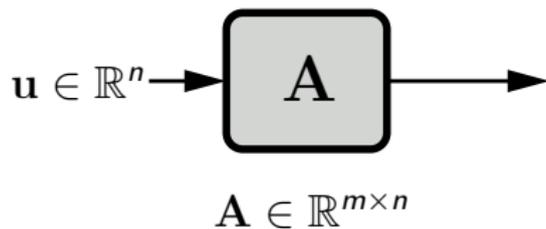
Saxe, McClelland & Ganguli, 2014; Dauphin *et al.*, 2014; Choromanska *et al.*, 2015; Haeffele & Vidal, 2015

Training

- Supervised training by **back propagation** (chain rule) results in **non-convex** optimization problem
- In deep networks, **local minima** are nearly as good as global ones
- Deep networks have less **saddle points**
- **Random initialization** works well
- **Extreme learning** strategies rely only on randomization

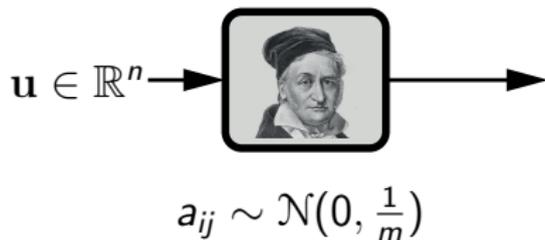
Saxe, McClelland & Ganguli, 2014; Dauphin *et al.*, 2014; Choromanska *et al.*, 2015; Haefele & Vidal, 2015

Assumptions



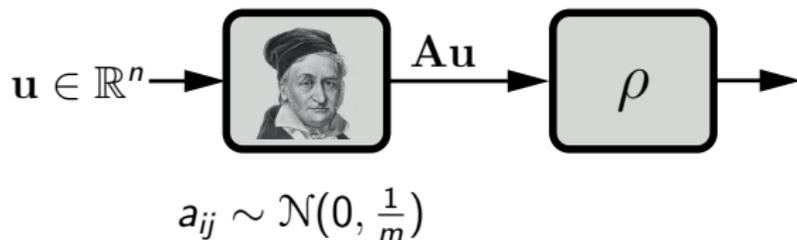
- Fully-connected linear layers

Assumptions



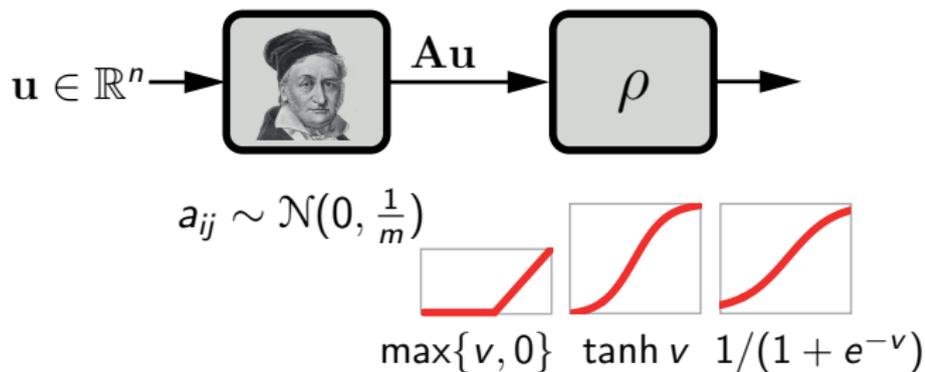
- Fully-connected linear layers with random Gaussian weights

Assumptions



- Fully-connected linear layers with random Gaussian weights
- Element-wise (approximately) truncated linear activation

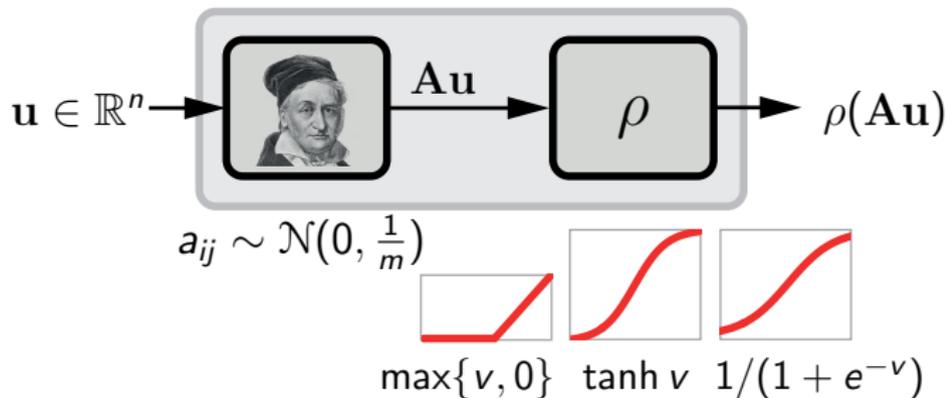
Assumptions



- Fully-connected linear layers with random Gaussian weights
- Element-wise (approximately) truncated linear activation

$$\rho|_{[a,b]} \text{ linear} \quad \rho|_{\mathbb{R} \setminus [a,b]} = \text{const}$$

Assumptions

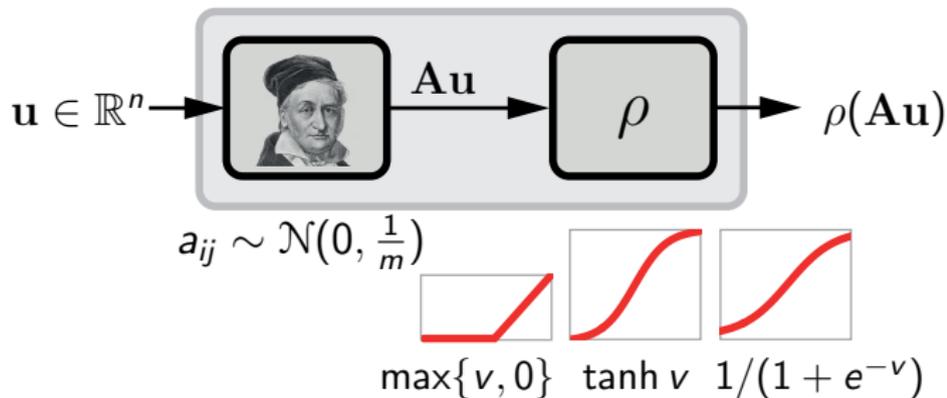


- Fully-connected linear layers with random Gaussian weights
- Element-wise (approximately) truncated linear activation

$$\rho|_{[a,b]} \text{ linear} \quad \rho|_{\mathbb{R} \setminus [a,b]} = \text{const}$$

- No pooling

Assumptions

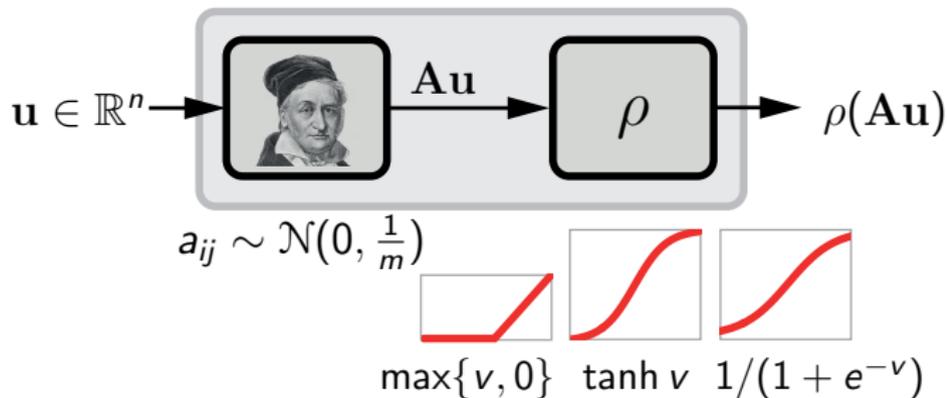


- Fully-connected linear layers with random Gaussian weights
- Element-wise (approximately) truncated linear activation

$$\rho|_{[a,b]} \text{ linear} \quad \rho|_{\mathbb{R} \setminus [a,b]} = \text{const}$$

- No pooling (pooling = invariance)

Assumptions

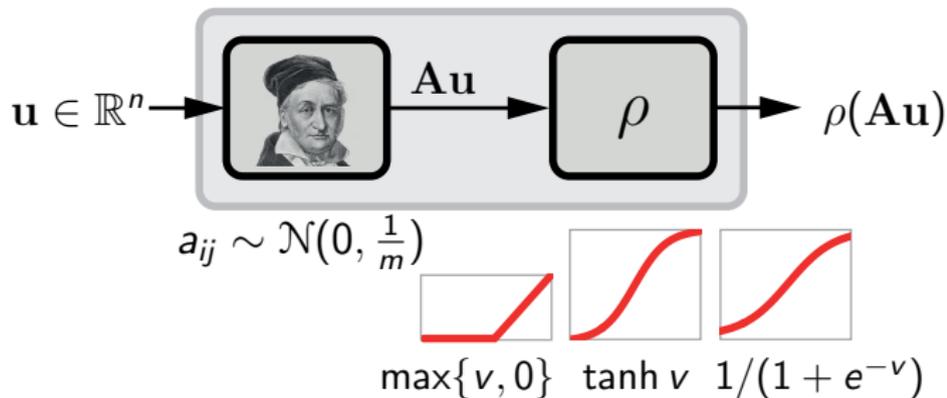


- Fully-connected linear layers with random Gaussian weights
- Element-wise (approximately) truncated linear activation

$$\rho|_{[a,b]} \text{ linear} \quad \rho|_{\mathbb{R} \setminus [a,b]} = \text{const}$$

- No pooling: input already invariant

Assumptions



- Fully-connected linear layers with random Gaussian weights
- Element-wise (approximately) truncated linear activation
 $\rho|_{[a,b]}$ linear $\rho|_{\mathbb{R} \setminus [a,b]} = \text{const}$
- No pooling: input already invariant
- Low dimensional input data

Near-Optimal Signal Recovery From Random Projections: Universal Encoding Strategies?

Emmanuel J. Candes and Terence Tao



Near-Optimal Signal Recovery From Random Projections: Universal Encoding Strategies?

Emmanuel J. Candes and Terence Tao

$$\mathbf{v} = \begin{bmatrix} \text{Portrait of a woman} \end{bmatrix} \mathbf{u}$$
$$a_{ij} \sim \mathcal{N}(0, \frac{1}{m})$$

Random Gaussian weights

$$\mathbf{v} = \left[\begin{array}{c} \text{Portrait of a woman} \end{array} \right] \mathbf{u}$$
$$a_{ij} \sim \mathcal{N}\left(0, \frac{1}{m}\right)$$

A k -sparse vector $\mathbf{u} \in \mathbb{R}^n$ can be reconstructed from $m = \mathcal{O}(k \log(n/k))$ random projections

Random Gaussian weights

$$\mathbf{v} = \left[\begin{array}{c} \text{Portrait of a woman} \end{array} \right] \mathbf{u}$$

$a_{ij} \sim \mathcal{N}(0, \frac{1}{m})$

A k -sparse vector $\mathbf{u} \in \mathbb{R}^n$ can be reconstructed from $m = \mathcal{O}(k \log(n/k))$ random projections

Restricted isometry property (RIP)

Random Gaussian weights

$$\mathbf{v} = \left[\begin{array}{c} \text{Portrait of a woman} \end{array} \right] \mathbf{u}$$

$a_{ij} \sim \mathcal{N}(0, \frac{1}{m})$

A k -sparse vector $\mathbf{u} \in \mathbb{R}^n$ can be reconstructed from $m = \mathcal{O}(k \log(n/k))$ random projections

Restricted isometry property (RIP)

Random projection is **universally good**

Low-dimensional input

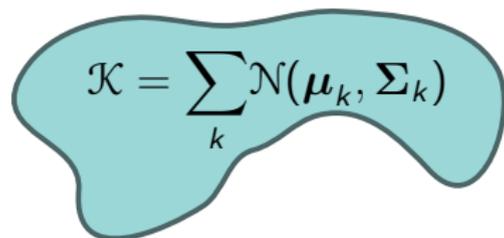
Input data have a small number of degrees of freedom

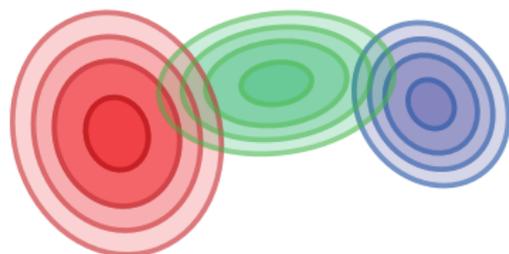
Low-dimensional input

**Input data have a small number of degrees of freedom
but may be embedded in a high-dimensional space**

Low-dimensional input

Input data have a **small number of degrees of freedom** but may be **embedded in a high-dimensional space**

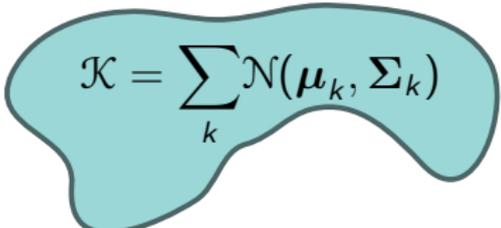

$$\mathcal{K} = \sum_k \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

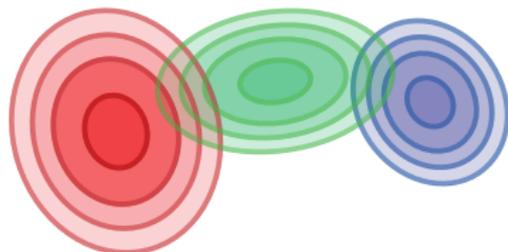


Gaussian mixture model

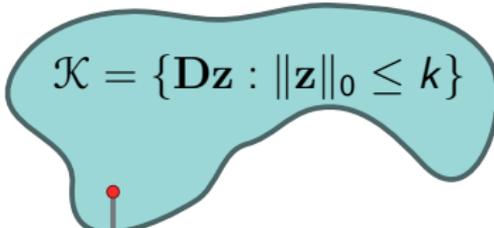
Low-dimensional input

Input data have a **small number of degrees of freedom** but may be **embedded in a high-dimensional space**

$$\mathcal{K} = \sum_k \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$




Gaussian mixture model

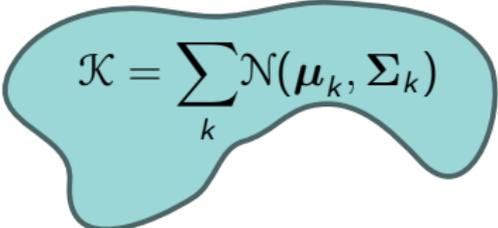
$$\mathcal{K} = \{\mathbf{D}\mathbf{z} : \|\mathbf{z}\|_0 \leq k\}$$


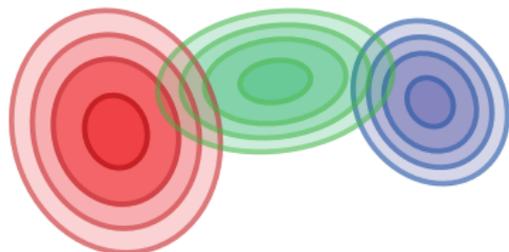


Sparse representation

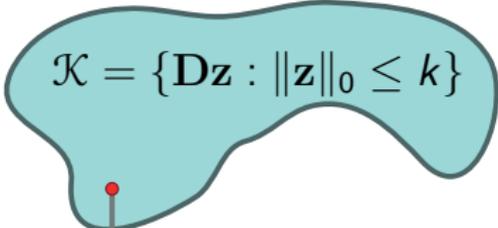
Low-dimensional input

Violated at the output due to DNN nonlinearity!

$$\mathcal{K} = \sum_k \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$




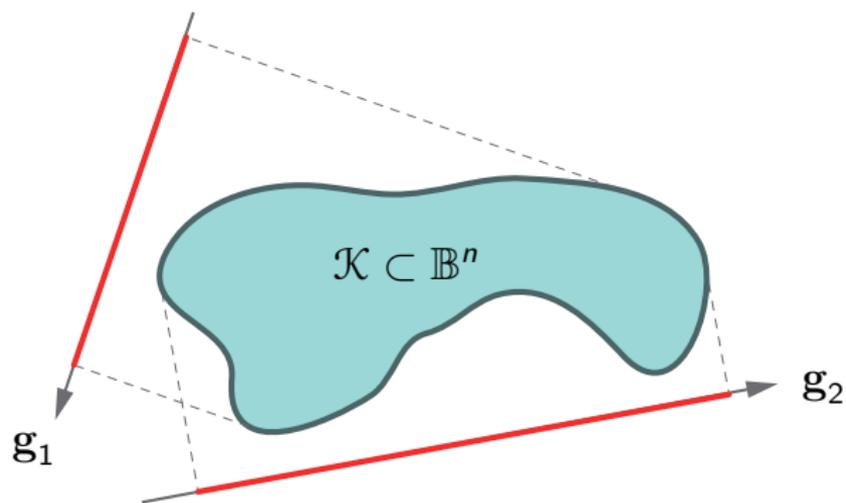
Gaussian mixture model

$$\mathcal{K} = \{\mathbf{Dz} : \|\mathbf{z}\|_0 \leq k\}$$




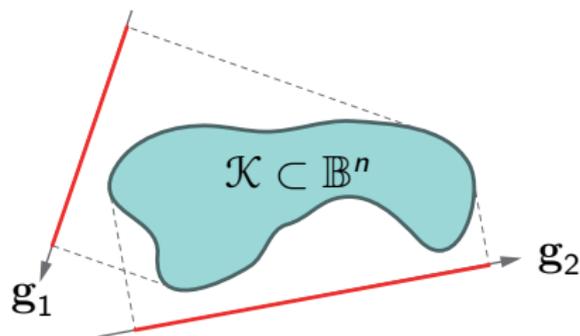
Sparse representation

Gaussian mean width



$$\omega(\mathcal{K}) = \mathbb{E} \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{K}} \langle \mathbf{u} - \mathbf{v}, \mathbf{g} \rangle \quad \mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

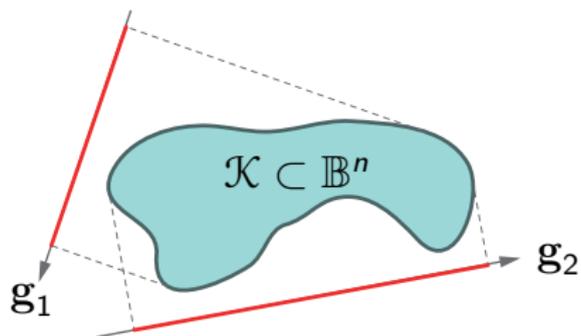
Gaussian mean width



$$\omega(\mathcal{K}) = \mathbb{E} \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{K}} \langle \mathbf{u} - \mathbf{v}, \mathbf{g} \rangle \quad \mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- $\omega^2(\mathcal{K})$ measures **intrinsic data dimension**

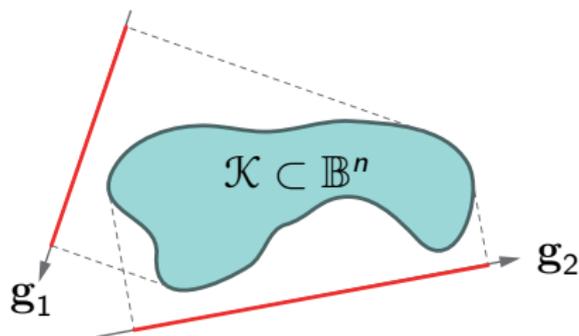
Gaussian mean width



$$\omega(\mathcal{K}) = \mathbb{E} \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{K}} \langle \mathbf{u} - \mathbf{v}, \mathbf{g} \rangle \quad \mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- $\omega^2(\mathcal{K})$ measures **intrinsic data dimension**
- \mathcal{K} is GMM with k Gaussians: $\omega^2(\mathcal{K}) = \mathcal{O}(k)$

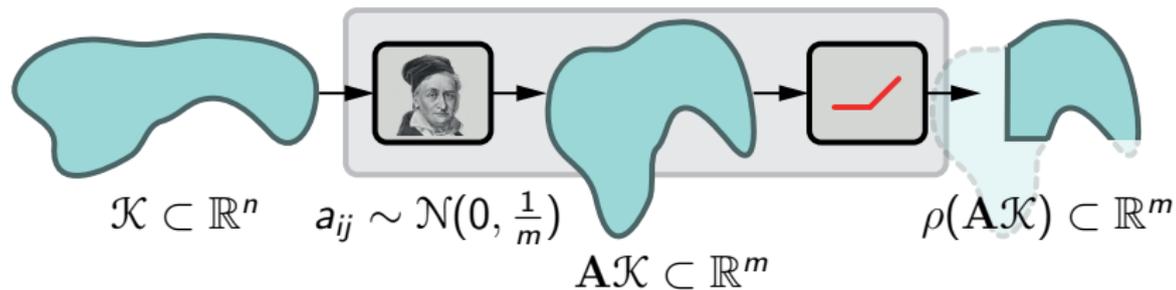
Gaussian mean width



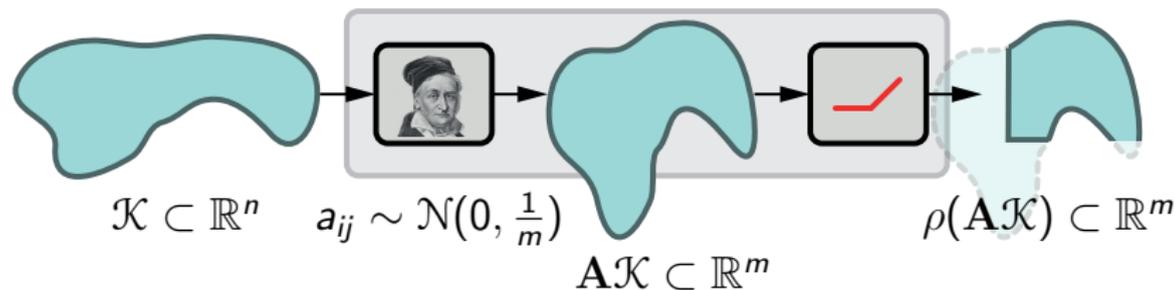
$$\omega^2(\mathcal{K}) = \mathbb{E} \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{K}} \langle \mathbf{u} - \mathbf{v}, \mathbf{g} \rangle \quad \mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- $\omega^2(\mathcal{K})$ measures **intrinsic data dimension**
- \mathcal{K} is GMM with k Gaussians: $\omega^2(\mathcal{K}) = \mathcal{O}(k)$
- \mathcal{K} is k -sparsely representable: $\omega^2(\mathcal{K}) = \mathcal{O}(k \log(n/k))$

Low dimension in – low dimension out

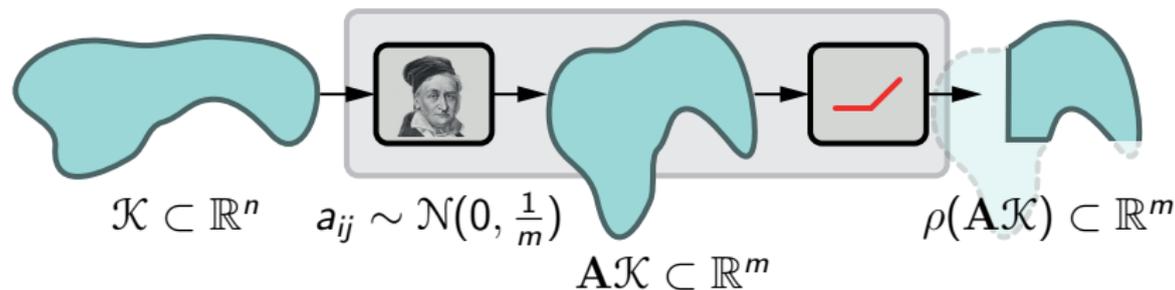


Low dimension in – low dimension out



Theorem: if $\omega^2(\mathcal{K}) \ll m$ then $\omega^2(\rho(\mathbf{A}\mathcal{K})) \approx \omega^2(\mathcal{K})$

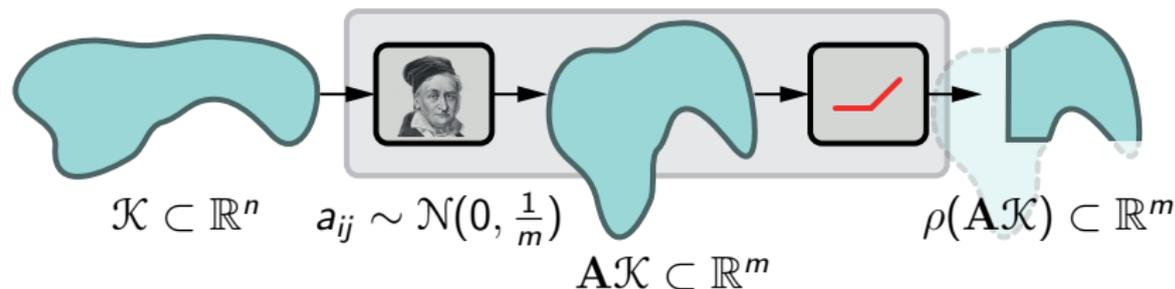
Low dimension in – low dimension out



Theorem: if $\omega^2(\mathcal{K}) \ll m$ then $\omega^2(\rho(\mathbf{AK})) \approx \omega^2(\mathcal{K})$

Proof: covering argument

Low dimension in – low dimension out

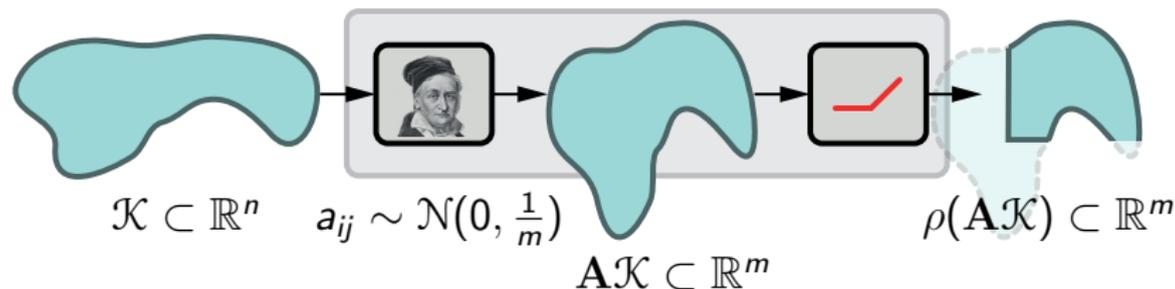


Theorem: if $\omega^2(\mathcal{K}) \ll m$ then $\omega^2(\rho(\mathbf{A}\mathcal{K})) \approx \omega^2(\mathcal{K})$

Proof: covering argument

- Intrinsic data dimension **does not grow** significantly through the network

Low dimension in – low dimension out



Theorem: if $\omega^2(\mathcal{K}) \ll m$ then $\omega^2(\rho(\mathbf{A}\mathcal{K})) \approx \omega^2(\mathcal{K})$

Proof: covering argument

- Intrinsic data dimension **does not grow** significantly through the network
- It is sufficient to analyze a **single layer** in DNN

Stable embedding

Theorem: the map $\mathbf{h} : (\mathcal{K} \subset \mathbb{S}^{n-1}, d_{\mathbb{S}^{n-1}}) \mapsto (\mathbf{h}(\mathcal{K}), d_{\mathbb{H}^m})$ defined by $\mathbf{h}(\mathbf{u}) = \text{sign}(\rho(\mathbf{A}\mathbf{u}))$ is a δ -isometry with $\delta = c m^{-1/6} \omega^{1/3}(\mathcal{K})$

Stable embedding

Theorem: the map $\mathbf{h} : (\mathcal{K} \subset \mathbb{S}^{n-1}, d_{\mathbb{S}^{n-1}}) \mapsto (\mathbf{h}(\mathcal{K}), d_{\mathbb{H}^m})$ defined by $\mathbf{h}(\mathbf{u}) = \text{sign}(\rho(\mathbf{A}\mathbf{u}))$ is a δ -isometry, i.e.,

$$|d_{\mathbb{S}^{n-1}}(\mathbf{u}, \mathbf{v}) - d_{\mathbb{H}^m}(\mathbf{h}(\mathbf{u}), \mathbf{h}(\mathbf{v}))| \leq \delta \quad \forall \mathbf{u}, \mathbf{v} \in \mathcal{K}$$

and every $\mathbf{w} \in \mathbf{h}(\mathcal{K})$ has some $\mathbf{u} \in \mathcal{K}$ such that $d_{\mathbb{H}^m}(\mathbf{h}(\mathbf{u}), \mathbf{w}) \leq \delta$ with $\delta = c m^{-1/6} \omega^{1/3}(\mathcal{K})$

Stable embedding

Theorem: the map $\mathbf{h} : (\mathcal{K} \subset \mathbb{S}^{n-1}, d_{\mathbb{S}^{n-1}}) \mapsto (\mathbf{h}(\mathcal{K}), d_{\mathbb{H}^m})$ defined by $\mathbf{h}(\mathbf{u}) = \text{sign}(\rho(\mathbf{A}\mathbf{u}))$ is a δ -isometry, i.e.,

$$|d_{\mathbb{S}^{n-1}}(\mathbf{u}, \mathbf{v}) - d_{\mathbb{H}^m}(\mathbf{h}(\mathbf{u}), \mathbf{h}(\mathbf{v}))| \leq \delta \quad \forall \mathbf{u}, \mathbf{v} \in \mathcal{K}$$

and every $\mathbf{w} \in \mathbf{h}(\mathcal{K})$ has some $\mathbf{u} \in \mathcal{K}$ such that $d_{\mathbb{H}^m}(\mathbf{h}(\mathbf{u}), \mathbf{w}) \leq \delta$ with $\delta = c m^{-1/6} \omega^{1/3}(\mathcal{K})$

Proof: follows Plan & Vershynin, 2013

Stable embedding

Theorem: the map $\mathbf{h} : (\mathcal{K} \subset \mathbb{S}^{n-1}, d_{\mathbb{S}^{n-1}}) \mapsto (\mathbf{h}(\mathcal{K}), d_{\mathbb{H}^m})$ defined by $\mathbf{h}(\mathbf{u}) = \text{sign}(\rho(\mathbf{A}\mathbf{u}))$ is a δ -isometry, i.e.,

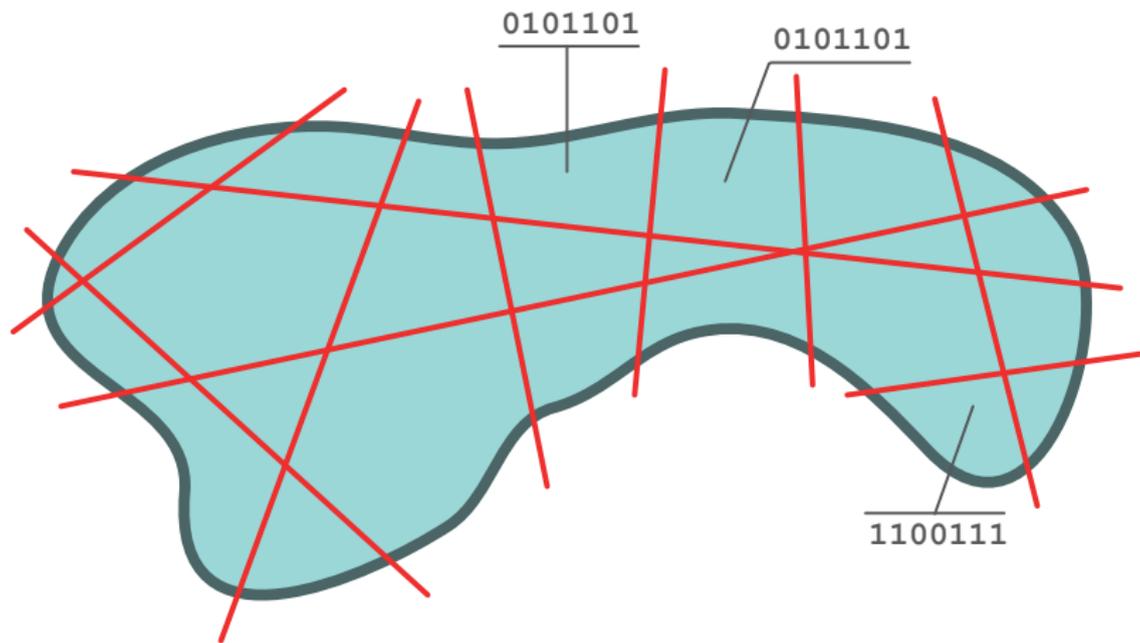
$$|d_{\mathbb{S}^{n-1}}(\mathbf{u}, \mathbf{v}) - d_{\mathbb{H}^m}(\mathbf{h}(\mathbf{u}), \mathbf{h}(\mathbf{v}))| \leq \delta \quad \forall \mathbf{u}, \mathbf{v} \in \mathcal{K}$$

and every $\mathbf{w} \in \mathbf{h}(\mathcal{K})$ has some $\mathbf{u} \in \mathcal{K}$ such that $d_{\mathbb{H}^m}(\mathbf{h}(\mathbf{u}), \mathbf{w}) \leq \delta$ with $\delta = c m^{-1/6} \omega^{1/3}(\mathcal{K})$

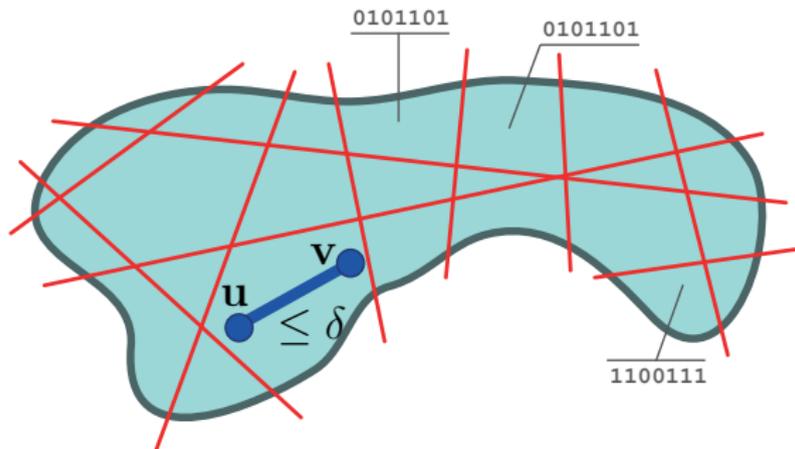
Proof: follows Plan & Vershynin, 2013

DNN layer performs stable embedding in the Gromov-Hausdorff sense

Tessellation of the input space

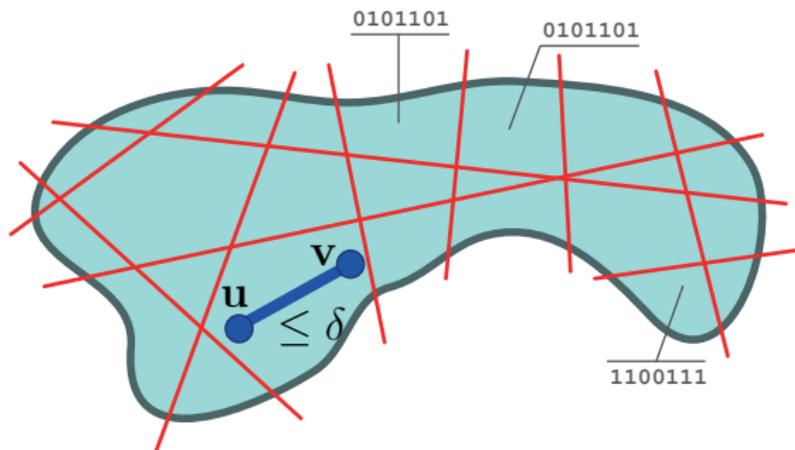


Tessellation of the input space



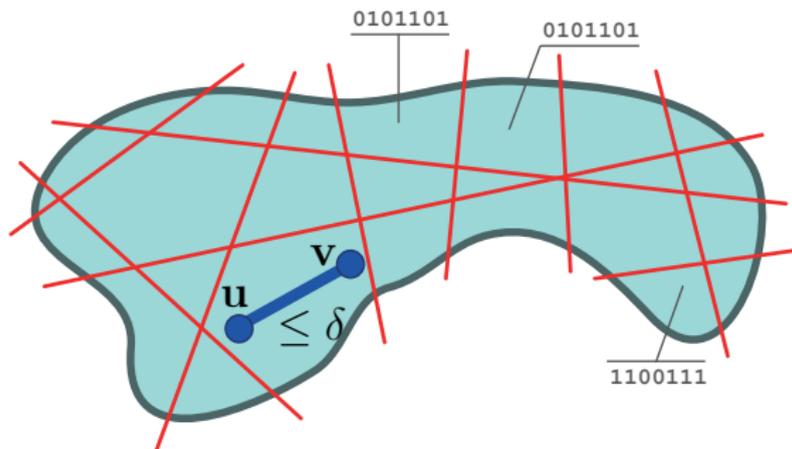
- Cell diameter $\leq \delta$

Tessellation of the input space



- Cell diameter $\leq \delta$
- If $\mathbf{h}(\mathbf{u}) = \mathbf{h}(\mathbf{v})$ then $d_{\mathbb{S}^{n-1}}(\mathbf{u}, \mathbf{v}) \leq \delta$

Tessellation of the input space



- Cell diameter $\leq \delta$
- If $\mathbf{h}(\mathbf{u}) = \mathbf{h}(\mathbf{v})$ then $d_{\mathbb{S}^{n-1}}(\mathbf{u}, \mathbf{v}) \leq \delta$
- **Input metric** can be recovered up to a small distortion

Stable embedding

Theorem: there exists a procedure \mathcal{P} such that

$$\|\mathcal{K} - \mathcal{P}(\rho(\mathbf{A}\mathcal{K}))\| < \mathcal{O}\left(\frac{\omega(\mathcal{K})}{\sqrt{m}}\right) = \mathcal{O}(\delta^3)$$

Stable embedding

Theorem: there exists a procedure \mathcal{P} such that

$$\|\mathcal{K} - \mathcal{P}(\rho(\mathbf{A}\mathcal{K}))\| < \mathcal{O}\left(\frac{\omega(\mathcal{K})}{\sqrt{m}}\right) = \mathcal{O}(\delta^3)$$

Proof: follows Plan & Vershynin, 2013

Theorem: there exists a procedure \mathcal{P} such that

$$\|\mathcal{K} - \mathcal{P}(\rho(\mathbf{A}\mathcal{K}))\| < \mathcal{O}\left(\frac{\omega(\mathcal{K})}{\sqrt{m}}\right) = \mathcal{O}(\delta^3)$$

Proof: follows Plan & Vershynin, 2013

- After N layers the error grows as $\mathcal{O}(N\delta^3)$

Theorem: there exists a procedure \mathcal{P} such that

$$\|\mathcal{K} - \mathcal{P}(\rho(\mathbf{A}\mathcal{K}))\| < \mathcal{O}\left(\frac{\omega(\mathcal{K})}{\sqrt{m}}\right) = \mathcal{O}(\delta^3)$$

Proof: follows Plan & Vershynin, 2013

- After N layers the error grows as $\mathcal{O}(N\delta^3)$
- DNNs **keep important information** of the data

Stable embedding

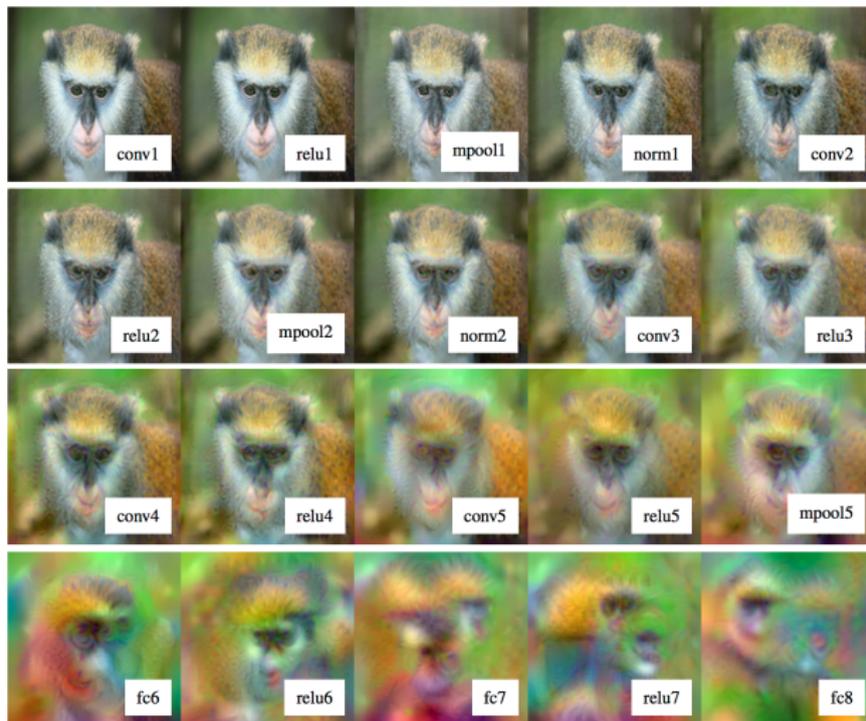
Theorem: there exists a procedure \mathcal{P} such that

$$\|\mathcal{K} - \mathcal{P}(\rho(\mathbf{A}\mathcal{K}))\| < \mathcal{O}\left(\frac{\omega(\mathcal{K})}{\sqrt{m}}\right) = \mathcal{O}(\delta^3)$$

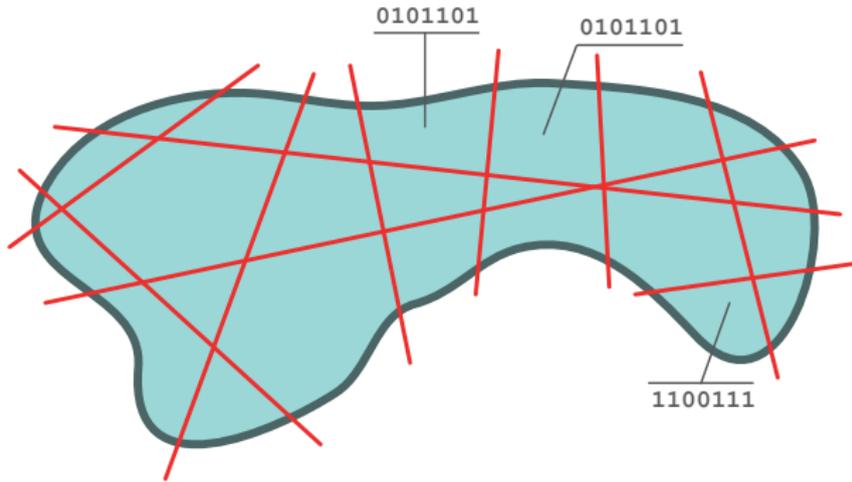
Proof: follows Plan & Vershynin, 2013

- After N layers the error grows as $\mathcal{O}(N\delta^3)$
- DNNs **keep important information** of the data
- **Input can be recovered from output** if output dimension m is big enough

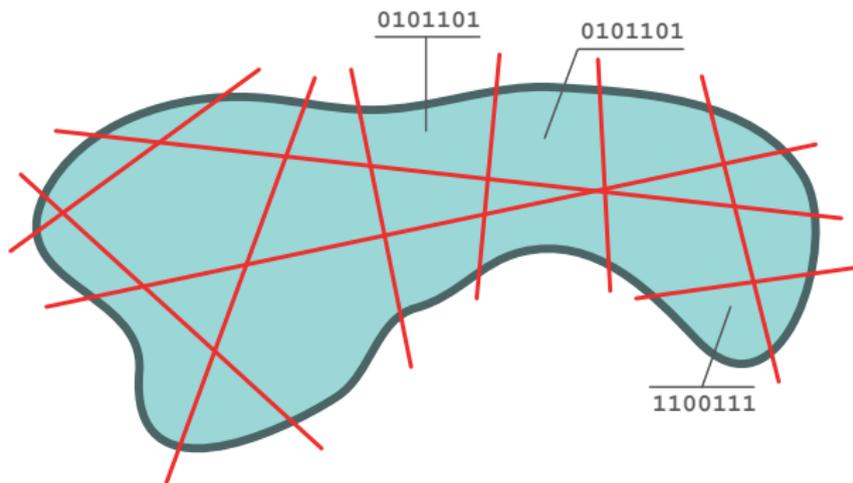
Inverting a CNN



Similarity-preserving hashing

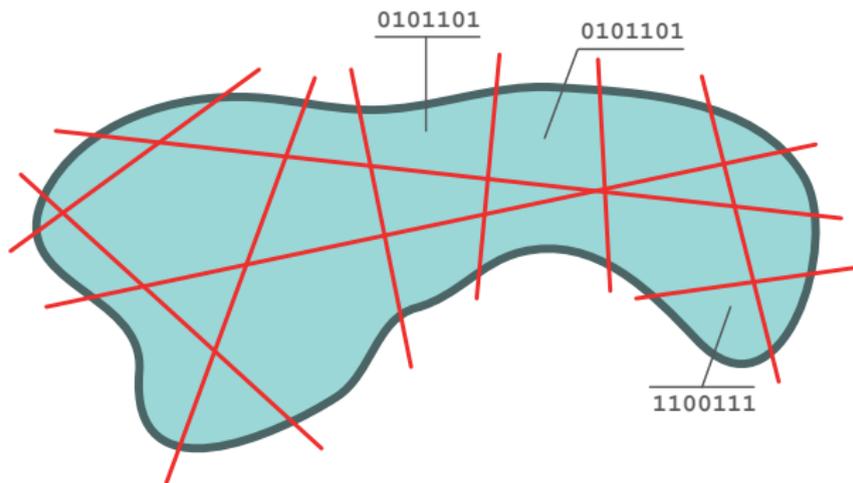


Similarity-preserving hashing



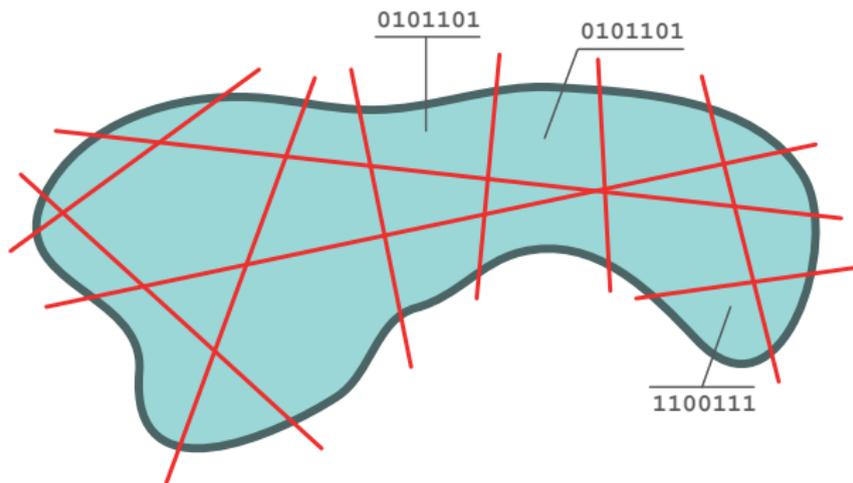
- Single layer = **locality sensitive hashing (LSH)**

Similarity-preserving hashing



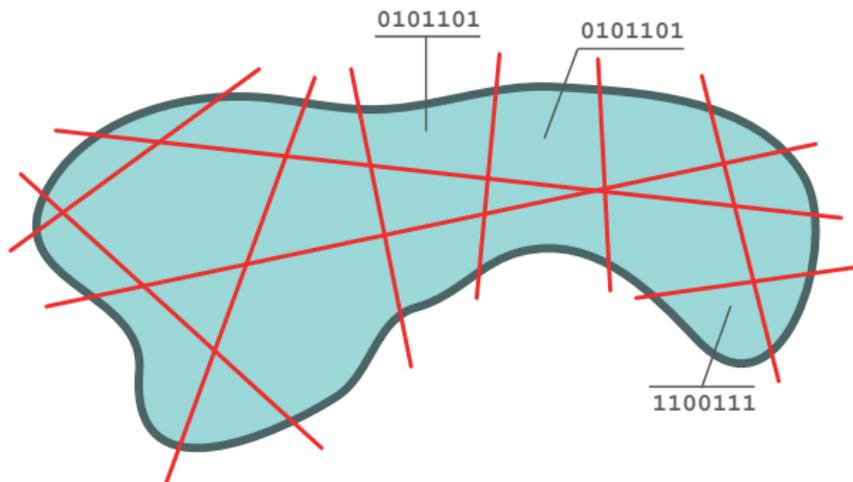
- Single layer = **locality sensitive hashing (LSH)**
- Random weights perform well **universally**

Similarity-preserving hashing



- Single layer = **locality sensitive hashing (LSH)**
- Random weights perform well **universally**
- Can be tuned to specific data by **training**

Similarity-preserving hashing



For **deep** networks, number of cells in the tessellation is **exponentially greater** than the number of degrees of freedom

Image retrieval with similarity-preserving hashing

60K images from 10 different classes taken from [Tiny images](#)

Represented using 384-dimensional [GIST descriptor](#)

[Training](#): 200 images per class; [Testing](#): 59K images

Method / m	12	24	48	
Raw		19.16		
DiffHash	14.72	13.35	12.85	
SSH	15.42	16.75	17.06	
AGH	15.46	15.29	15.15	
KSH	25.79	29.01	30.84	
NN	1 layer	31.48	35.41	36.79
	2 layer	45.42	49.88	50.46

Performance (mAP in %)

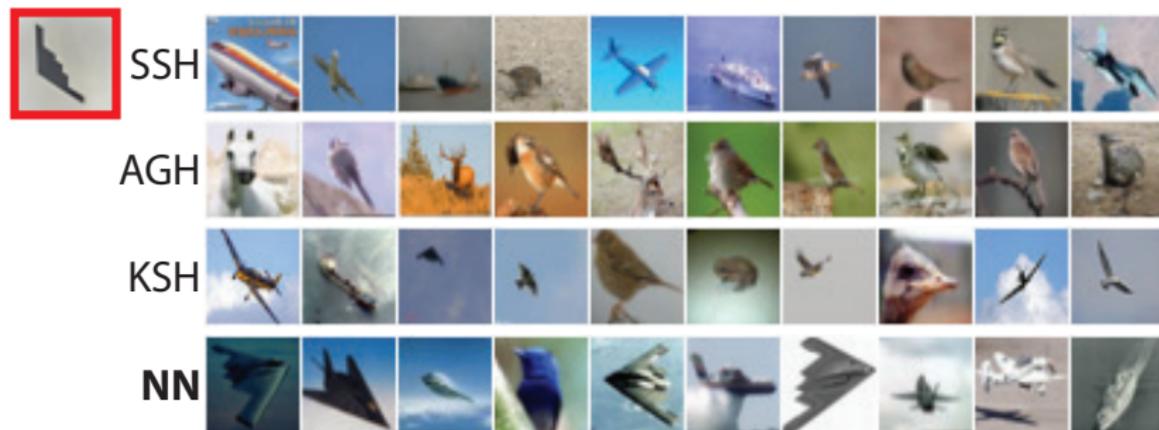
Data: Torralba et al. 2008, Krizhevsky 2009; Methods: Strecha et al. 2011 (diff-hash); Shakhnarovich 2005 (SSH); Liu et al. 2011 (AGH); Liu et al. 2012 (KSH); Masci, B², Schmidhuber 2012 (NN)

Image retrieval with similarity-preserving hashing

60K images from 10 different classes taken from [Tiny images](#)

Represented using 384-dimensional [GIST descriptor](#)

Training: 200 images per class; **Testing:** 59K images



Ranking using 48-bit hashes

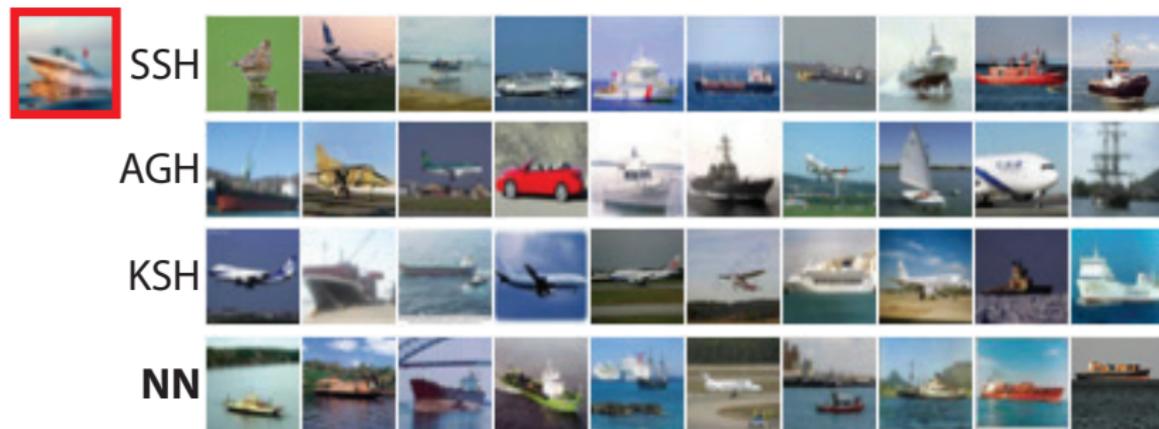
Data: Torralba et al. 2008, Krizhevsky 2009; Methods: Strecha et al. 2011 (diff-hash); Shakhnarovich 2005 (SSH); Liu et al. 2011 (AGH); Liu et al. 2012 (KSH); Masci, B², Schmidhuber 2012 (NN)

Image retrieval with similarity-preserving hashing

60K images from 10 different classes taken from [Tiny images](#)

Represented using 384-dimensional [GIST descriptor](#)

Training: 200 images per class; **Testing:** 59K images



Ranking using 48-bit hashes

Data: Torralba et al. 2008, Krizhevsky 2009; Methods: Strecha et al. 2011 (diff-hash); Shakhnarovich 2005 (SSH); Liu et al. 2011 (AGH); Liu et al. 2012 (KSH); Masci, B², Schmidhuber 2012 (NN)

Image retrieval with similarity-preserving hashing

60K images from 10 different classes taken from [Tiny images](#)

Represented using 384-dimensional [GIST descriptor](#)

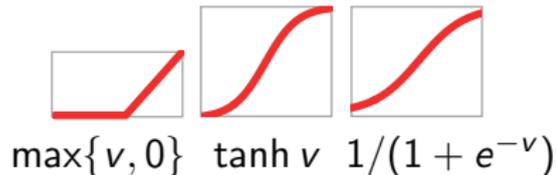
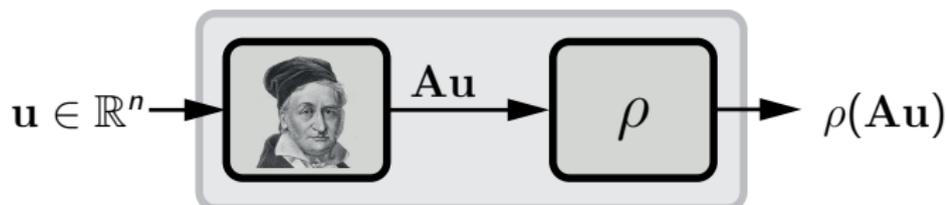
Training: 200 images per class; **Testing:** 59K images



Ranking using 48-bit hashes

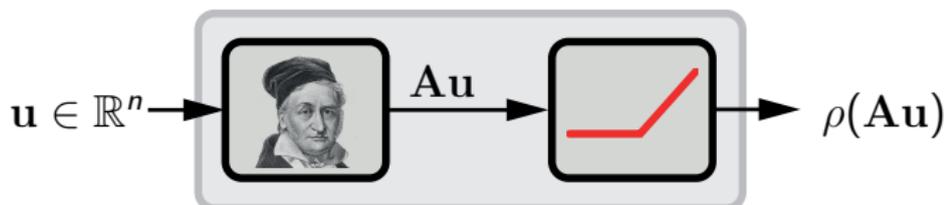
Data: Torralba et al. 2008, Krizhevsky 2009; Methods: Strecha et al. 2011 (diff-hash); Shakhnarovich 2005 (SSH); Liu et al. 2011 (AGH); Liu et al. 2012 (KSH); Masci, B², Schmidhuber 2012 (NN)

Assumptions



- Fully-connected linear layers with random Gaussian weights
- Element-wise approximately truncated linear activation
- No pooling
- Low dimensional input data

Assumptions



- Fully-connected linear layers with random Gaussian weights
- ReLU activation
- No pooling
- Low dimensional input data

Angle distortion

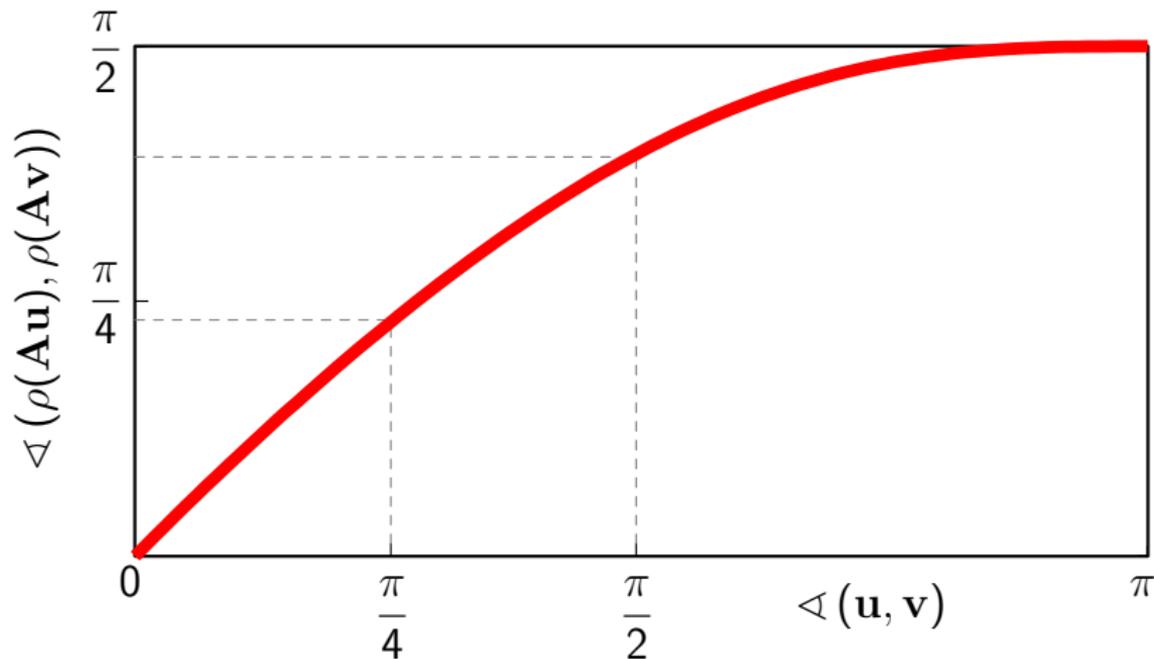
Theorem (concentration of output angle): for $\mathbf{u}, \mathbf{v} \in \mathcal{K}$

$$\cos \angle(\rho(\mathbf{A}\mathbf{u}), \rho(\mathbf{A}\mathbf{v})) \approx \cos \angle(\mathbf{u}, \mathbf{v}) + \psi(\angle(\mathbf{u}, \mathbf{v}))$$

where $\psi(\alpha) = \frac{1}{\pi}(\sin \alpha - \alpha \cos \alpha)$

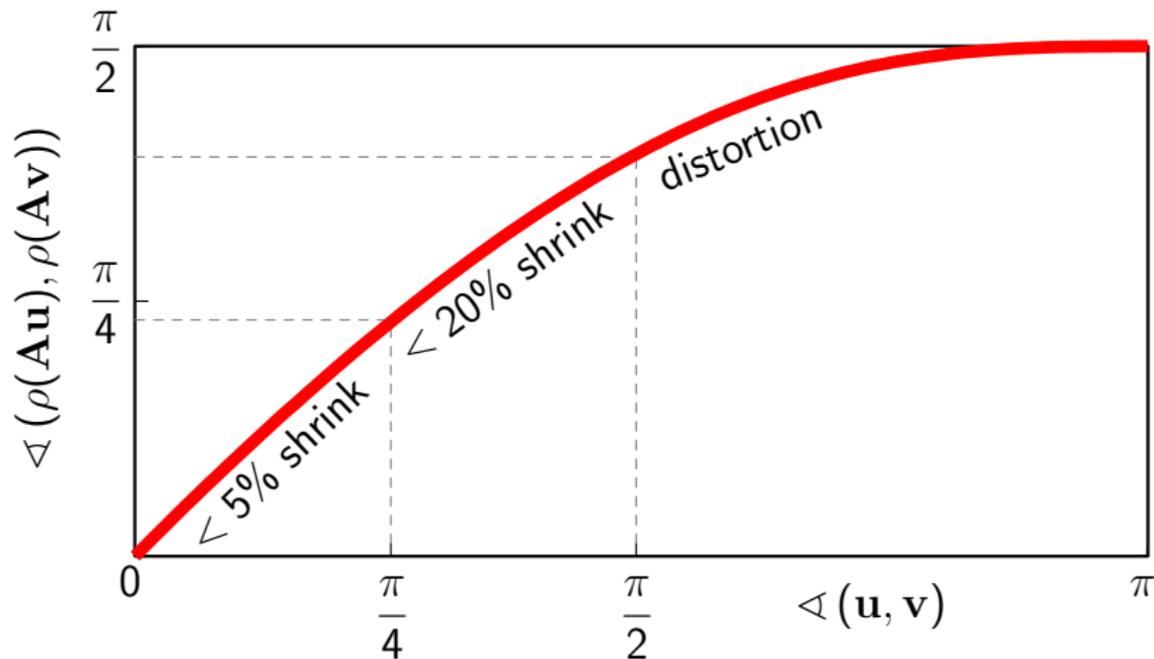
Angle distortion

$$\cos \angle(\rho(\mathbf{A}\mathbf{u}), \rho(\mathbf{A}\mathbf{v})) \approx \cos \angle(\mathbf{u}, \mathbf{v}) + \psi(\angle(\mathbf{u}, \mathbf{v}))$$



Angle distortion

$$\cos \angle(\rho(\mathbf{A}\mathbf{u}), \rho(\mathbf{A}\mathbf{v})) \approx \cos \angle(\mathbf{u}, \mathbf{v}) + \psi(\angle(\mathbf{u}, \mathbf{v}))$$



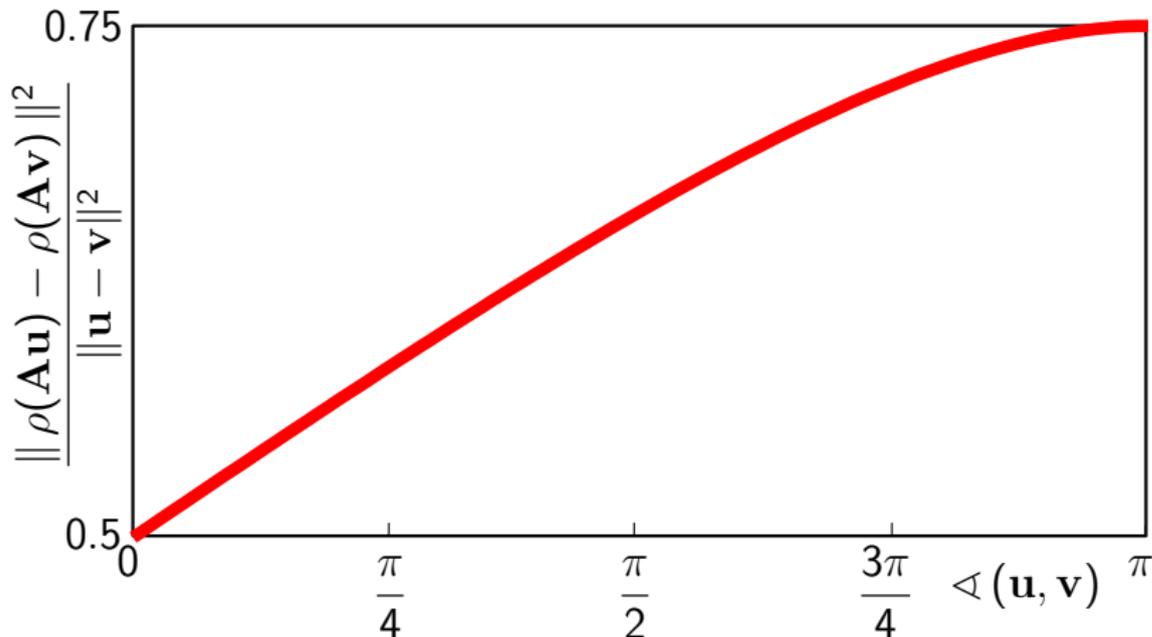
Theorem (concentration of output distance):

$$\|\rho(\mathbf{A}\mathbf{u}) - \rho(\mathbf{A}\mathbf{v})\|^2 \approx \frac{1}{2}\|\mathbf{u} - \mathbf{v}\|^2 + \|\mathbf{u}\|\|\mathbf{v}\|\psi(\angle(\mathbf{u}, \mathbf{v}))$$

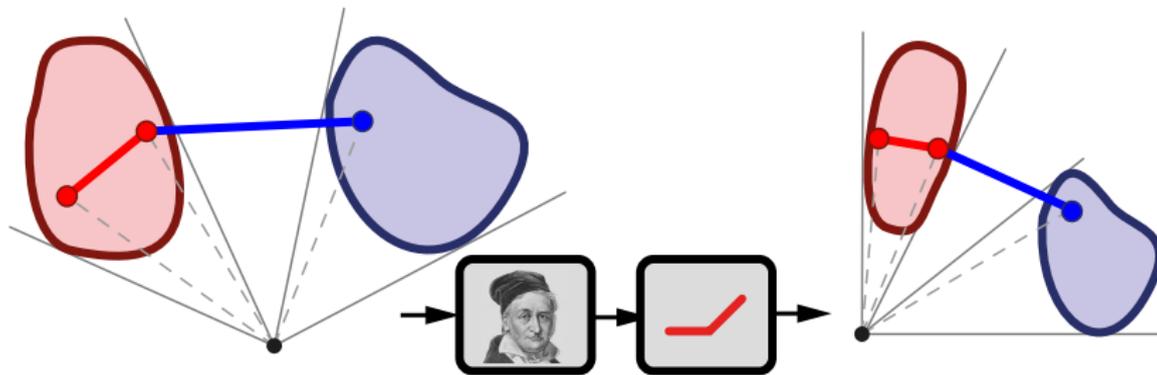
for $\mathbf{u}, \mathbf{v} \in \mathcal{K}$, where $\psi(\alpha) = \frac{1}{\pi}(\sin \alpha - \alpha \cos \alpha)$

Distance distortion

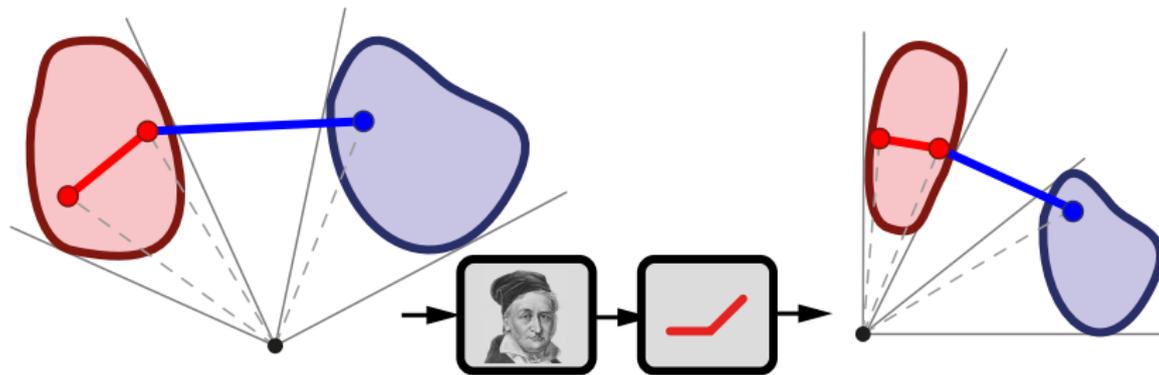
$$\|\rho(\mathbf{A}\mathbf{u}) - \rho(\mathbf{A}\mathbf{v})\|^2 \approx \frac{1}{2}\|\mathbf{u} - \mathbf{v}\|^2 + \|\mathbf{u}\|\|\mathbf{v}\|\psi(\angle(\mathbf{u}, \mathbf{v}))$$



Angle and distance distortion



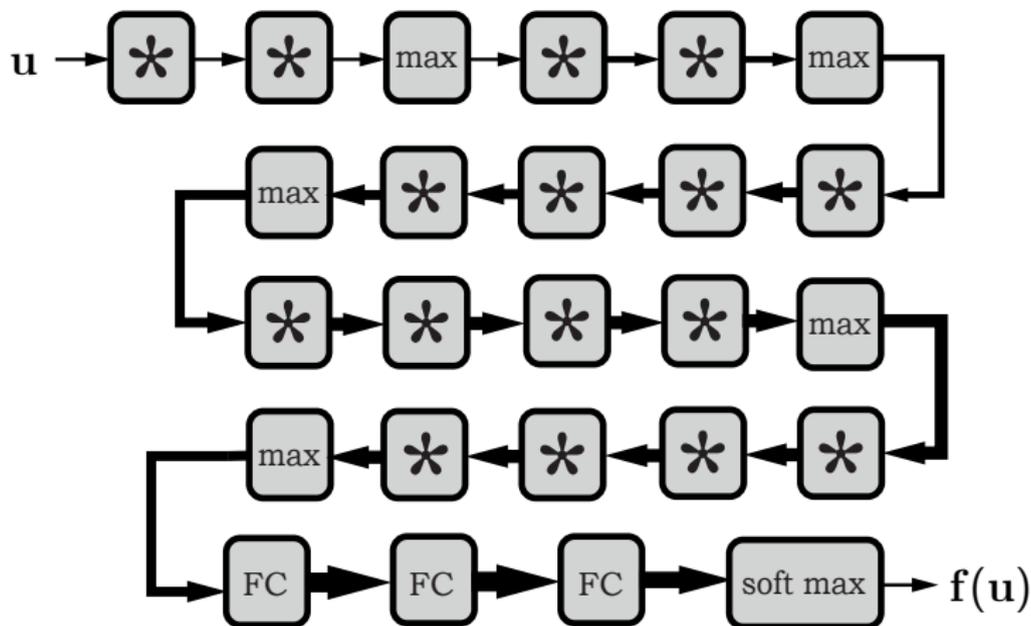
Angle and distance distortion



- Points with **small angles** between them become **closer** than points with large angles between them

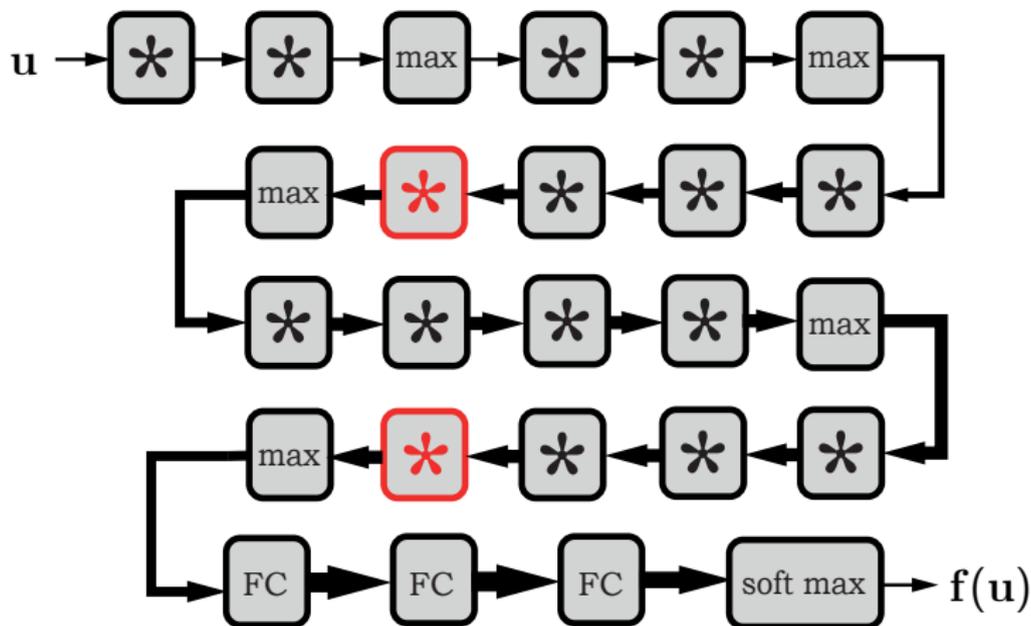
Inside a real network

State-of-the-art 19-layer CNN trained on ImageNet



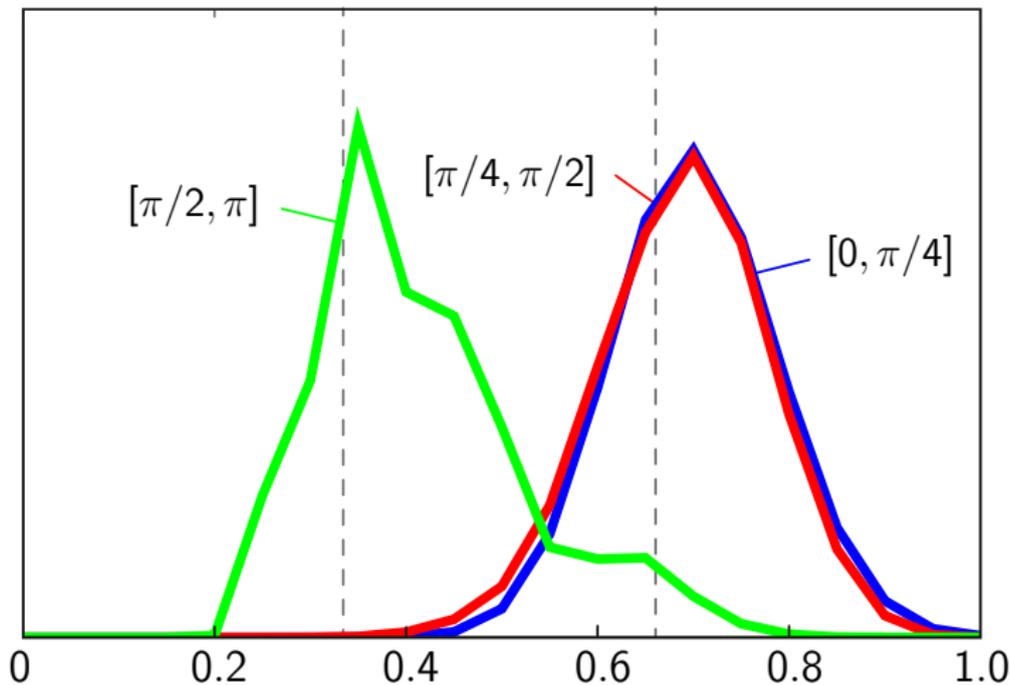
Inside a real network

State-of-the-art 19-layer CNN trained on ImageNet



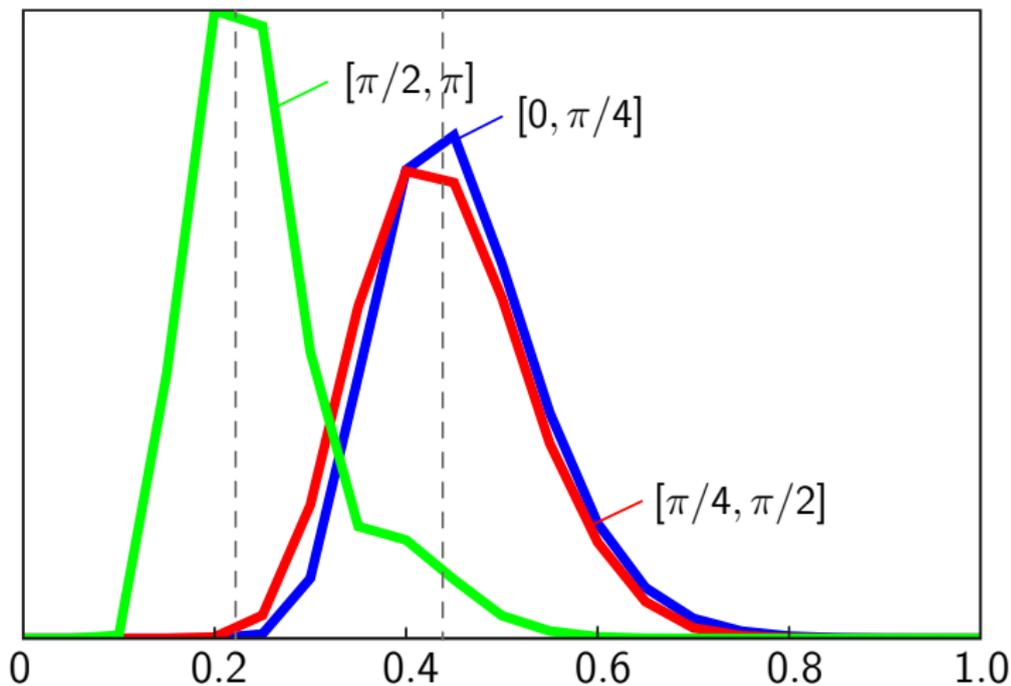
Angle distortion at 8-th layer

Distribution of $\angle(f(u), f(v)) / \angle(u, v)$



Angle distortion at 16-th layer

Distribution of $\angle(f(\mathbf{u}), f(\mathbf{v})) / \angle(\mathbf{u}, \mathbf{v})$



Training set size

- DNNs are **stable**: close points in the input are close in the output

Training set size

- DNNs are **stable**: close points in the input are close in the output
- Network performing well on an ϵ -net \mathcal{K}_ϵ in \mathcal{K} performs well on entire \mathcal{K}

Training set size

- DNNs are **stable**: close points in the input are close in the output
- Network performing well on an ϵ -net \mathcal{K}_ϵ in \mathcal{K} performs well on entire \mathcal{K}
- **Sudakov's minoration:**

$$\log |\mathcal{K}_\epsilon| \leq \frac{c \omega^2(\mathcal{K})}{\epsilon^2}$$

Training set size

- DNNs are **stable**: close points in the input are close in the output
- Network performing well on an ϵ -net \mathcal{K}_ϵ in \mathcal{K} performs well on entire \mathcal{K}
- **Sudakov's minoration:**

$$\log |\mathcal{K}_\epsilon| \leq \frac{c \omega^2(\mathcal{K})}{\epsilon^2}$$

- **Not tight!**

Training set size

- DNNs are **stable**: close points in the input are close in the output
- Network performing well on an ϵ -net \mathcal{K}_ϵ in \mathcal{K} performs well on entire \mathcal{K}
- **Sudakov's minoration:**

$$\log |\mathcal{K}_\epsilon| \leq \frac{c \omega^2(\mathcal{K})}{\epsilon^2}$$

- **Not tight!** ...but introduces Gaussian mean width $\omega(\mathcal{K})$ as the measure of data complexity

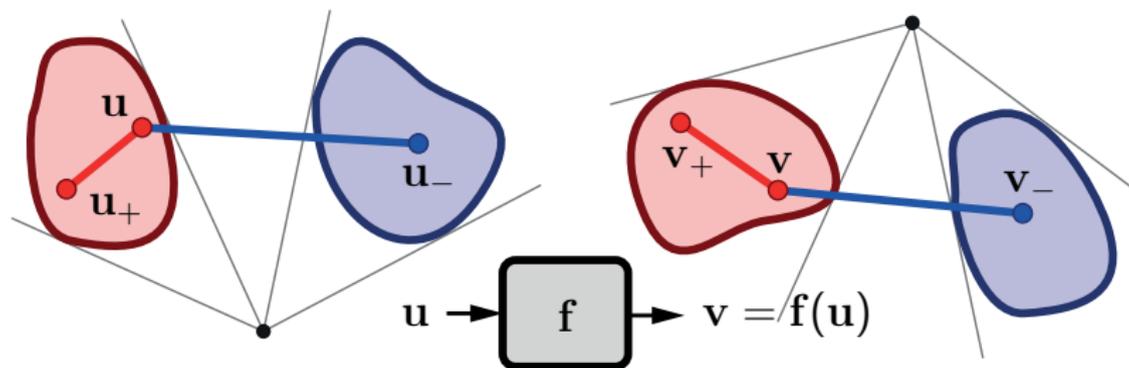
Training set size

- DNNs are **stable**: close points in the input are close in the output
- Network performing well on an ϵ -net \mathcal{K}_ϵ in \mathcal{K} performs well on entire \mathcal{K}
- **Sudakov's minoration:**

$$\log |\mathcal{K}_\epsilon| \leq \frac{c \omega^2(\mathcal{K})}{\epsilon^2}$$

- **Not tight!** ...but introduces Gaussian mean width $\omega(\mathcal{K})$ as the measure of data complexity
- Situation is **much better in practice**

Random data points



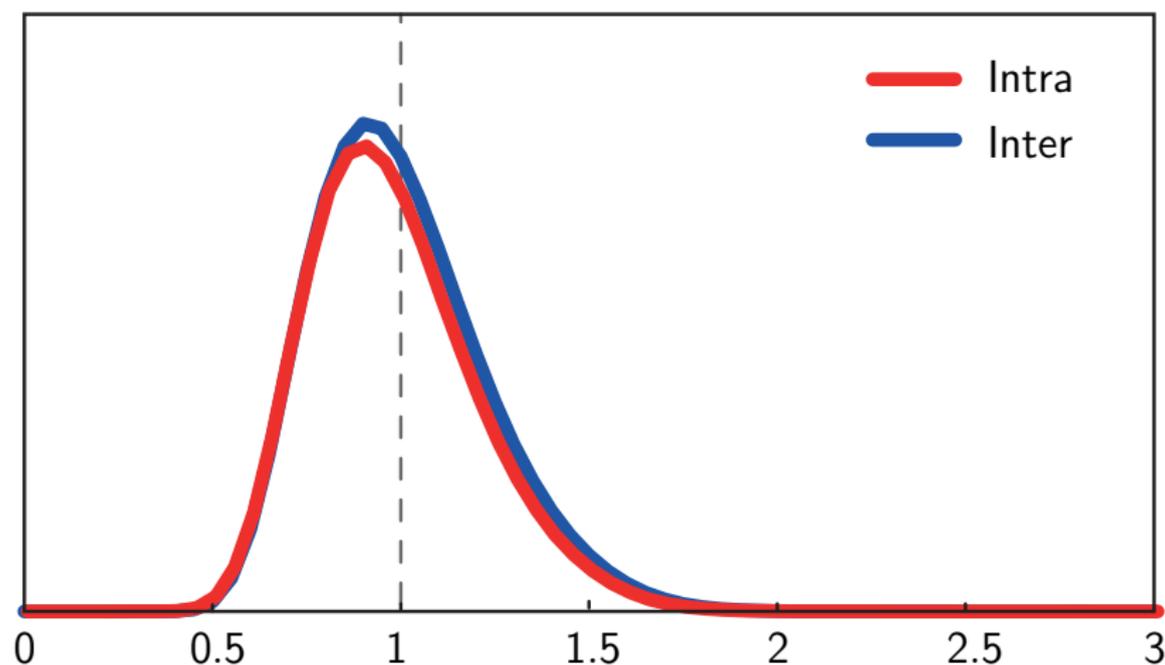
Distance ratios between random triplets (u, u_+, u_-)

$$\text{Intra-class} = \frac{\|v_+ - v\|}{\|u_+ - u\|}$$

$$\text{Inter-class} = \frac{\|v_- - v\|}{\|u_- - u\|}$$

Random data points

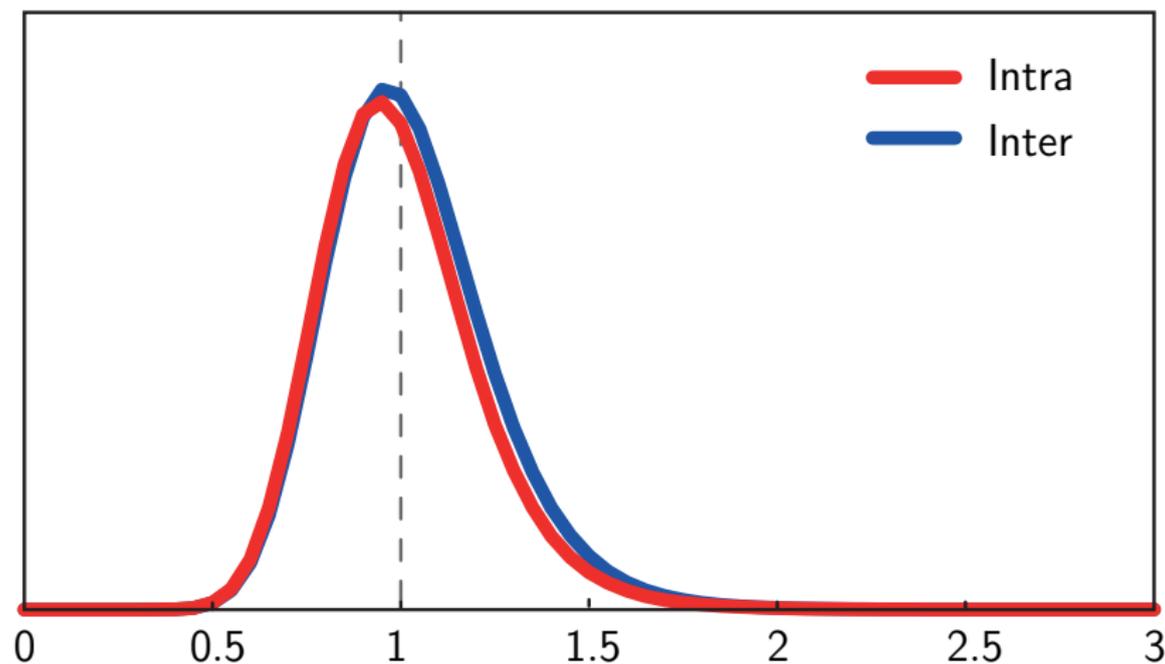
CNN on CIFAR-10 – Random weights



Giryes, Sapiro, B, 2015

Random data points

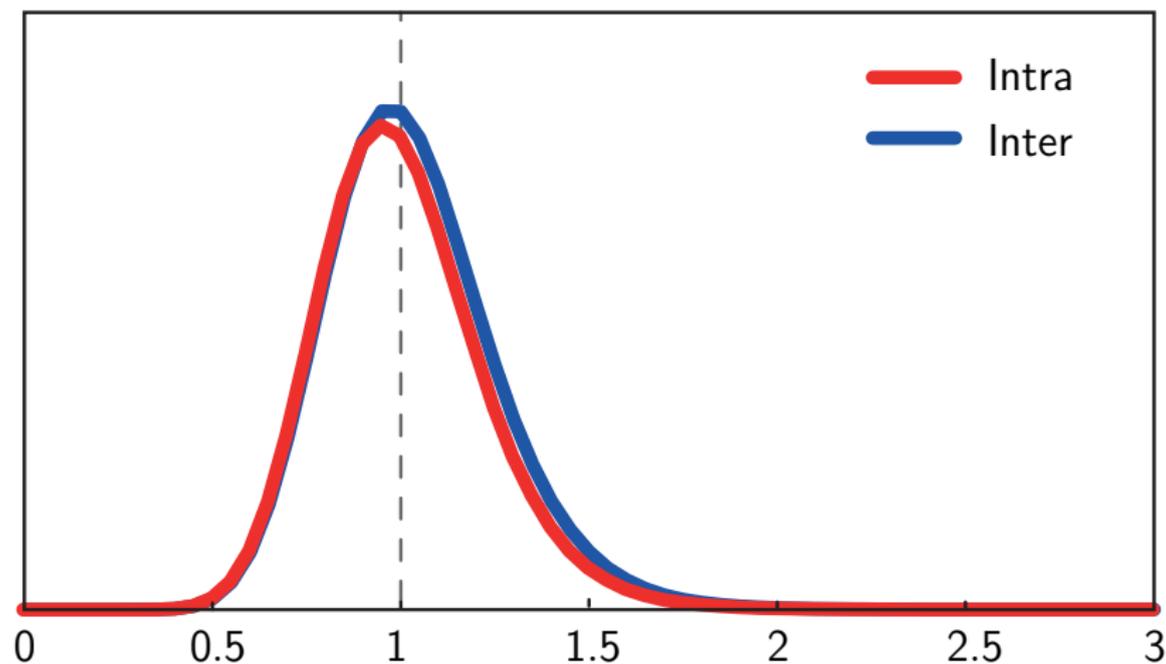
CNN on CIFAR-10 – Trained to 25% error



Giryes, Sapiro, B, 2015

Random data points

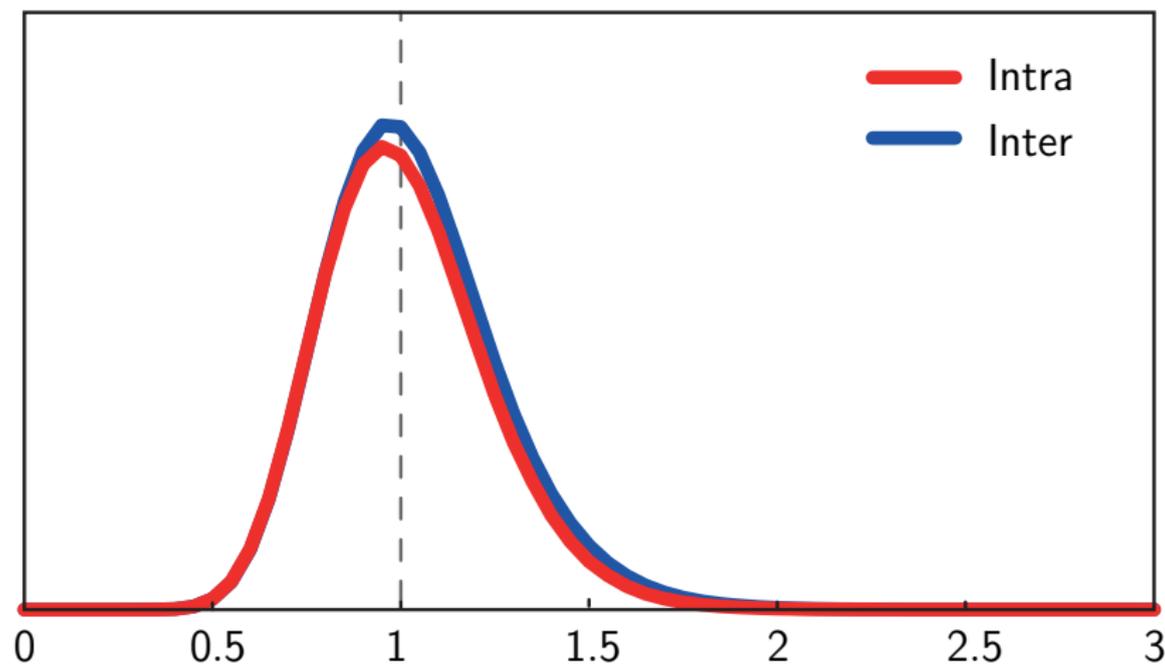
CNN on CIFAR-10 – Trained to 21% error



Giryès, Sapiro, B, 2015

Random data points

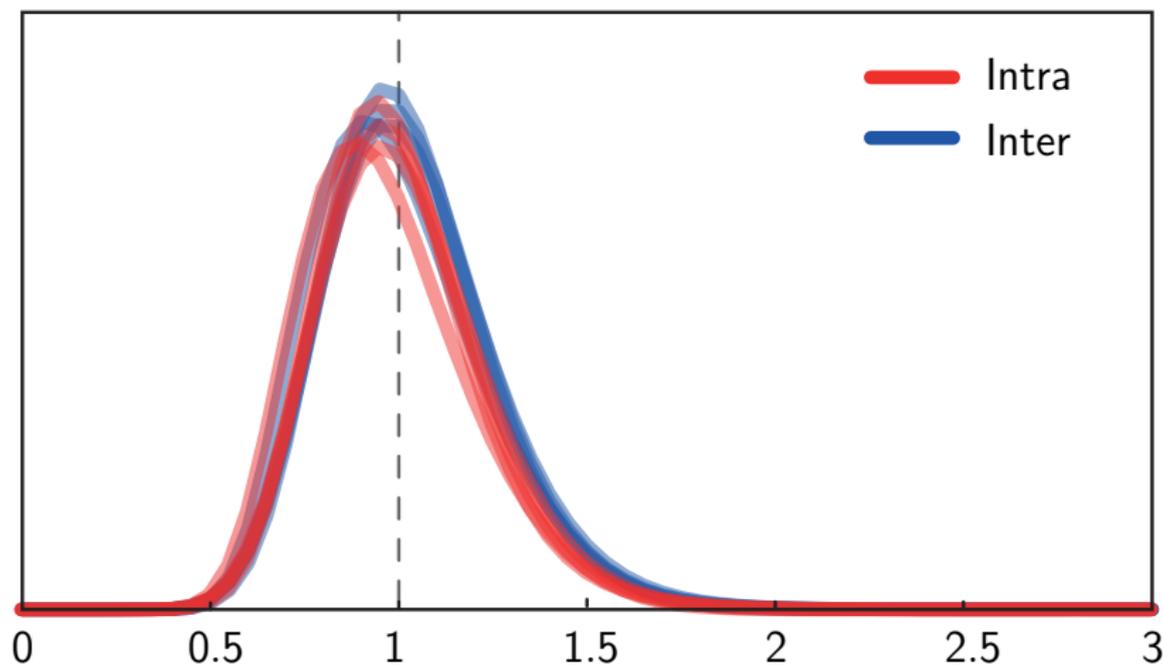
CNN on CIFAR-10 – Trained to 18% error



Giryès, Sapiro, B, 2015

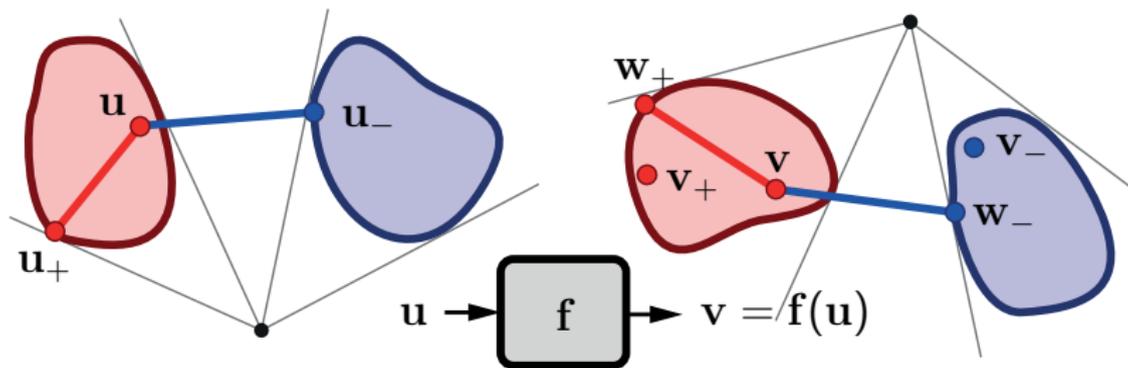
Random data points

CNN on CIFAR-10 – Random and trained



Giryes, Sapiro, B, 2015

Class boundary points



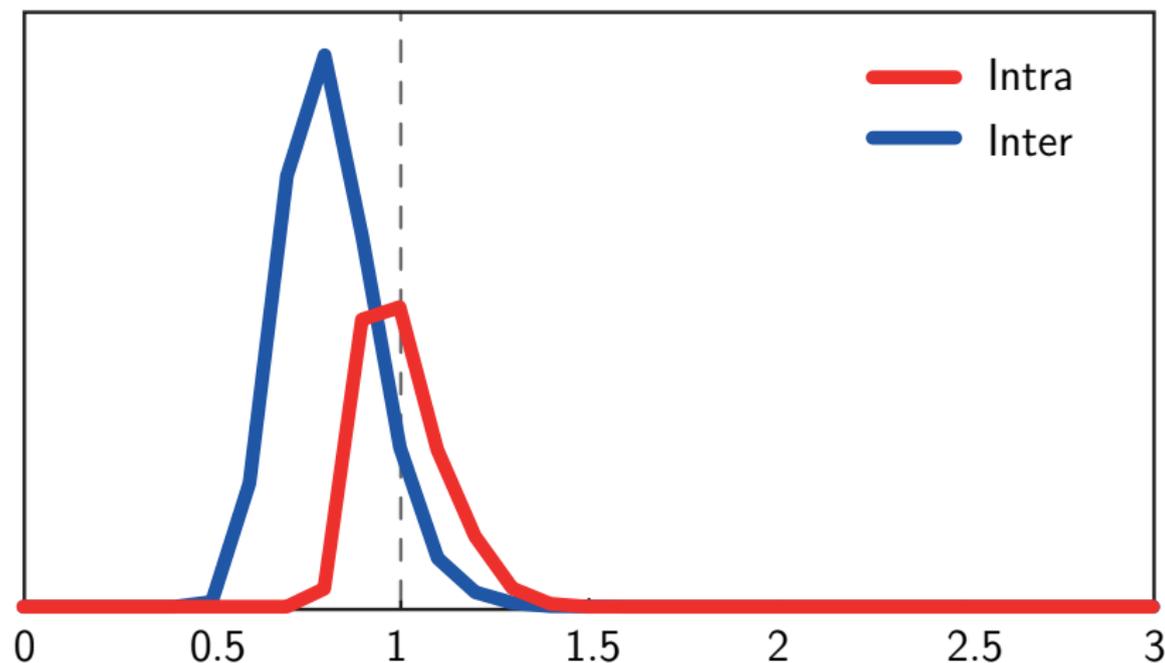
u random, u_+ farthest in class, u_- closest not in class

$$\text{Intra-class} = \frac{\|w_+ - v\|}{\|u_+ - u\|}$$

$$\text{Inter-class} = \frac{\|w_- - v\|}{\|u_- - u\|}$$

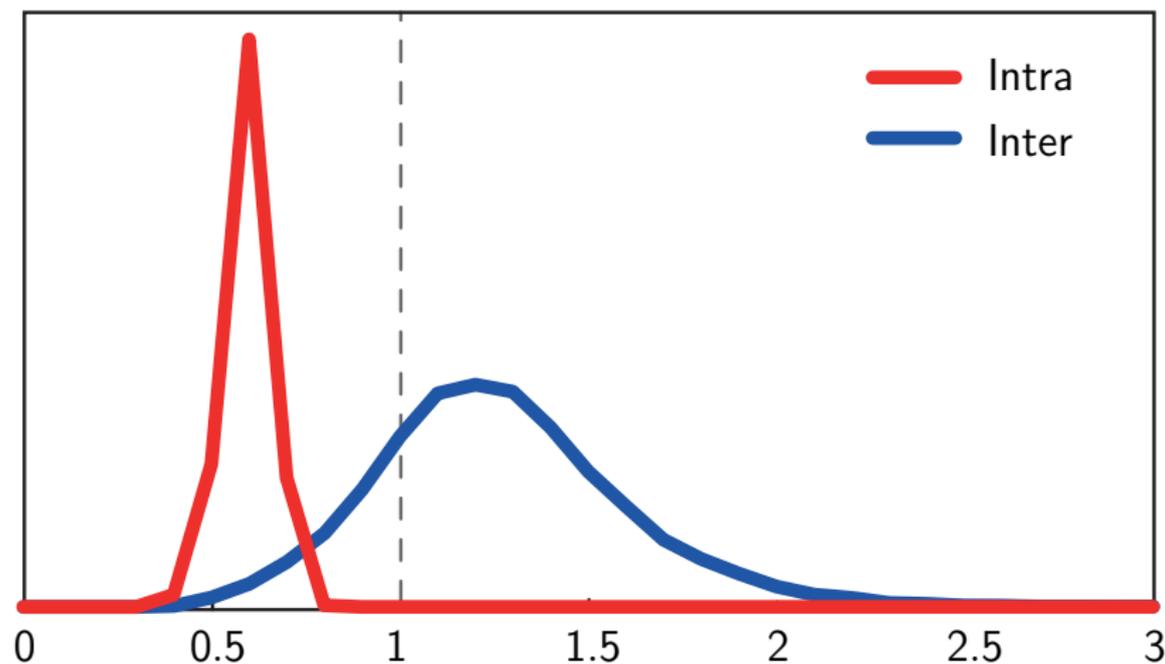
Class boundary points

CNN on CIFAR-10 – Random weights



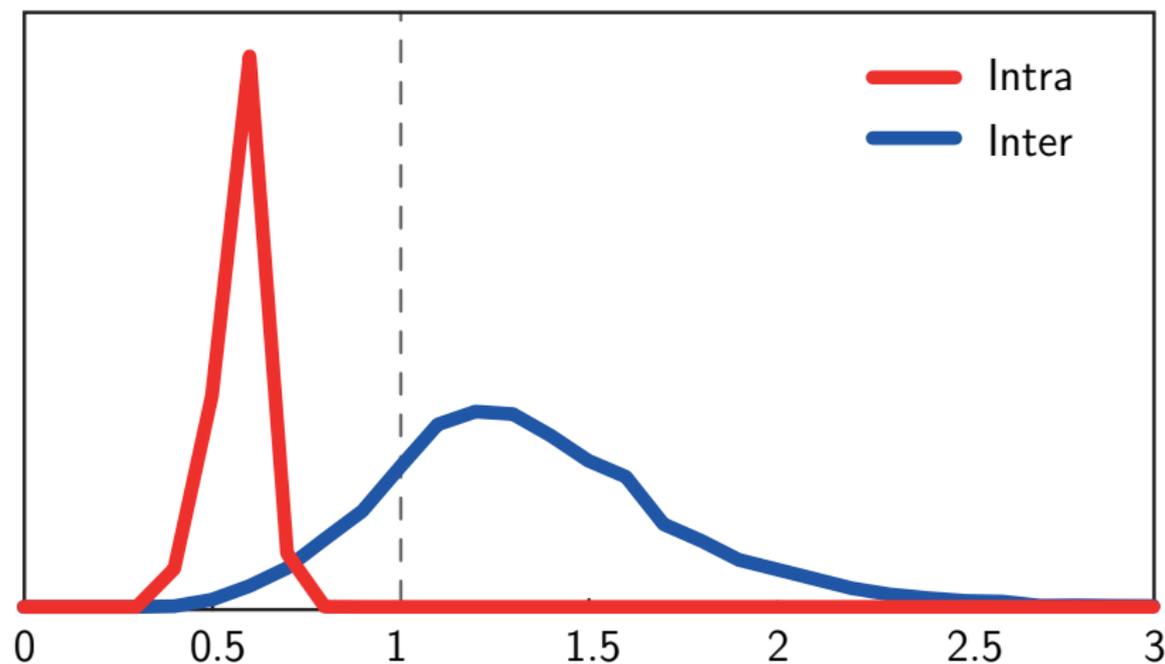
Class boundary points

CNN on CIFAR-10 – Trained to 25% error



Class boundary points

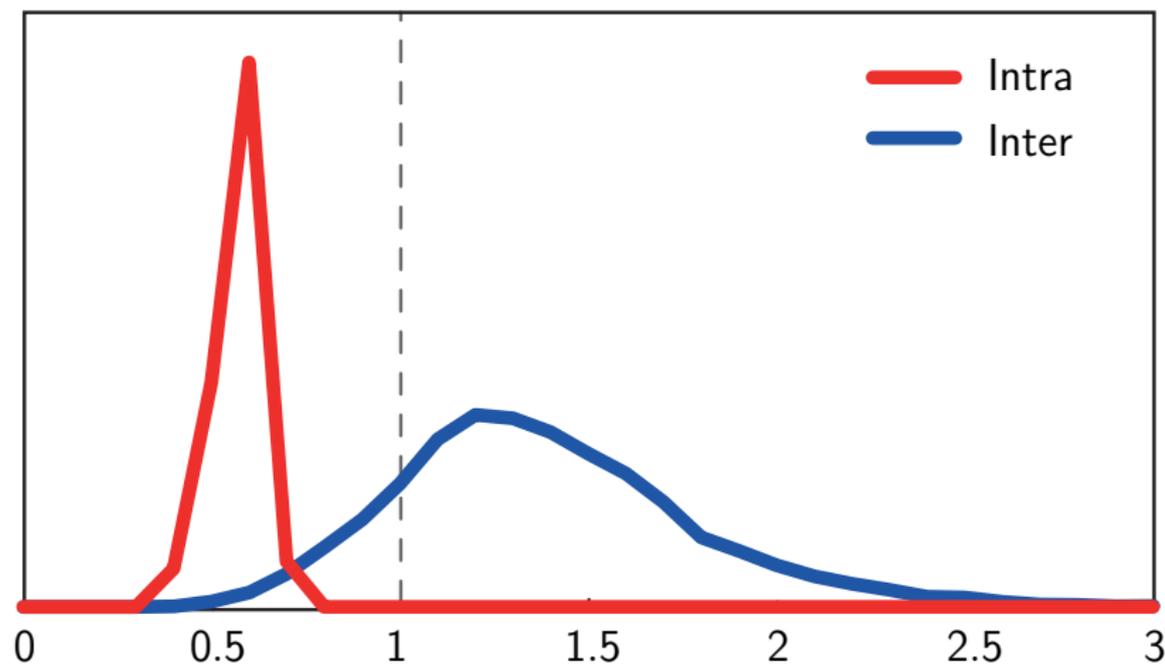
CNN on CIFAR-10 – Trained to 21% error



Giryes, Sapiro, B, 2015

Class boundary points

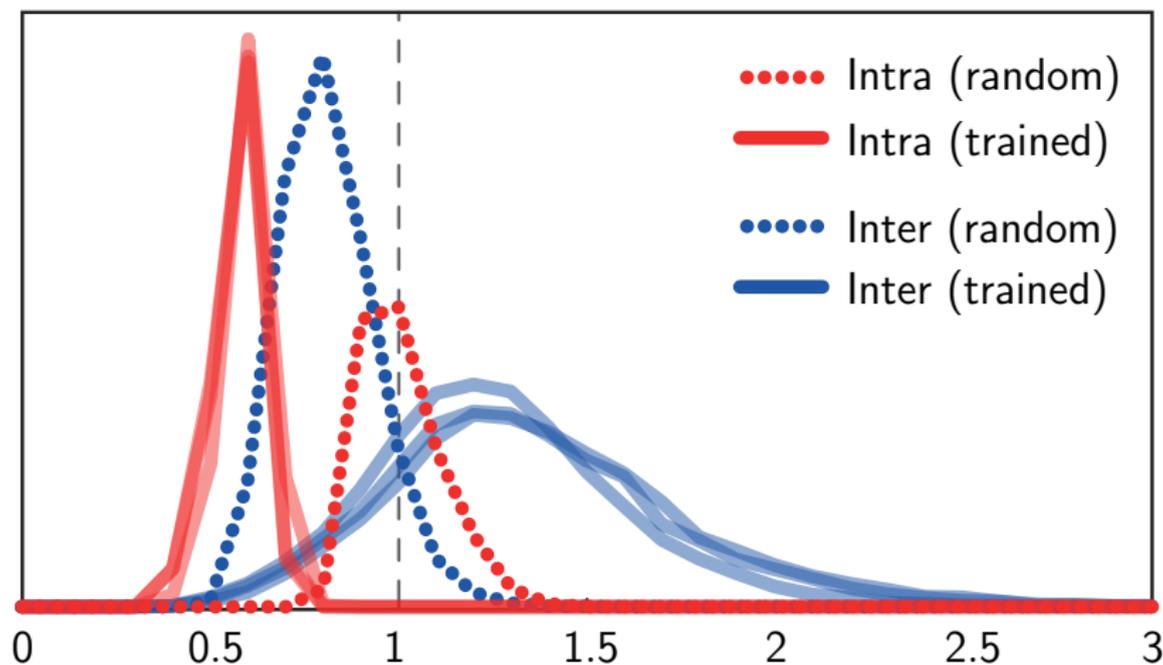
CNN on CIFAR-10 – Trained to 18% error



Giryes, Sapiro, B, 2015

Class boundary points

CNN on CIFAR-10 – Random and trained



Role of training

- Negligible effect on **random** data points

Role of training

- Negligible effect on **random** data points
- Random weights perform well **universally**

Role of training

- Negligible effect on **random** data points
- Random weights perform well **universally**
- Major effect on **class boundary** points

Role of training

- Negligible effect on **random** data points
 - Random weights perform well **universally**
 - Major effect on **class boundary** points
- Intra-class distances shrink**

Role of training

- Negligible effect on **random** data points
- Random weights perform well **universally**
- Major effect on **class boundary** points

Intra-class distances shrink

Inter-class distances grow

Role of training

- Negligible effect on **random** data points
- Random weights perform well **universally**
- Major effect on **class boundary** points
 - Intra-class distances shrink**
 - Inter-class distances grow**
- Only a small subset of \mathcal{K}_ϵ is required for training

DNNs as metric learners

- Massive supervision required for DNN training

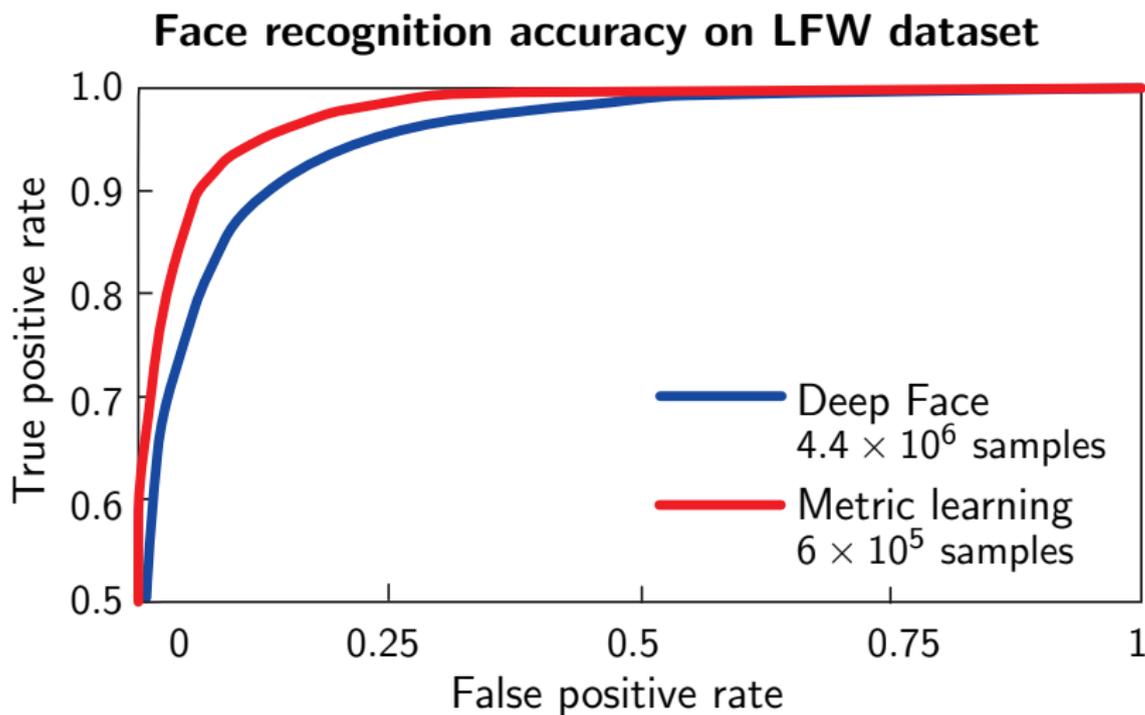
DNNs as metric learners

- Massive supervision required for DNN training
- Semi- and unsupervised training is a challenge

DNNs as metric learners

- Massive supervision required for DNN training
- Semi- and unsupervised training is a challenge
- Inject metric learning criterion into training objective to reduce the amount of labeled data

DNNs as metric learners



Compressed discrimination

Compressed sensing: reconstruct signal $\mathbf{x} \in \mathbb{R}^n$ given $m \ll n$ measurements

Compressed discrimination

Compressed sensing: reconstruct signal $\mathbf{x} \in \mathbb{R}^n$ given $m \ll n$ measurements

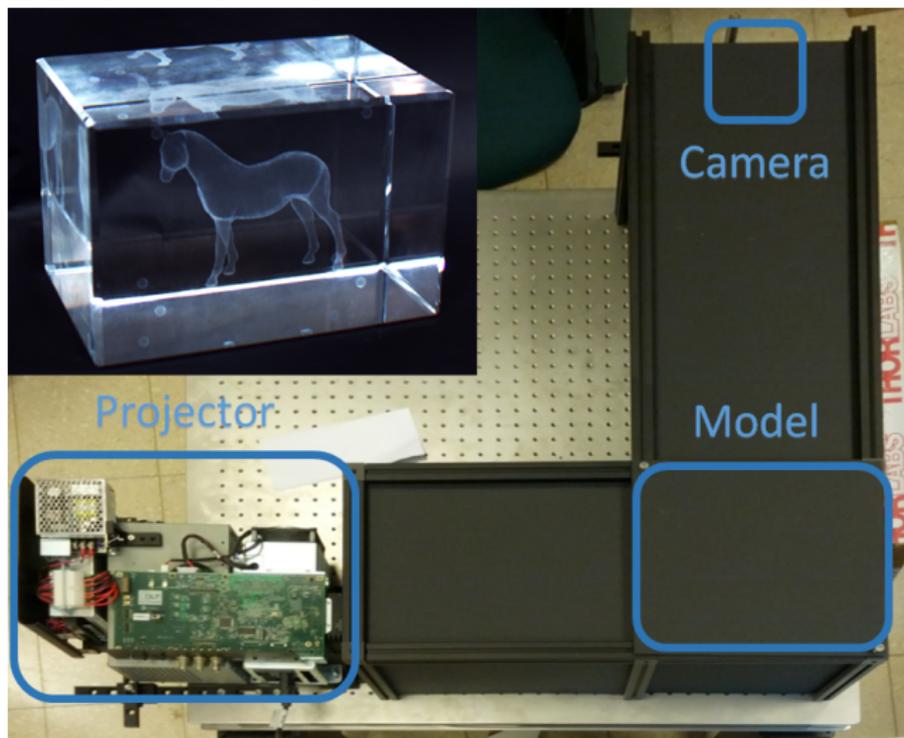
Compressed discrimination: estimate parameter $\theta \in \mathbb{R}^k$ ($k \ll n$) related to \mathbf{x} given $m' < m$ measurements

Compressed discrimination

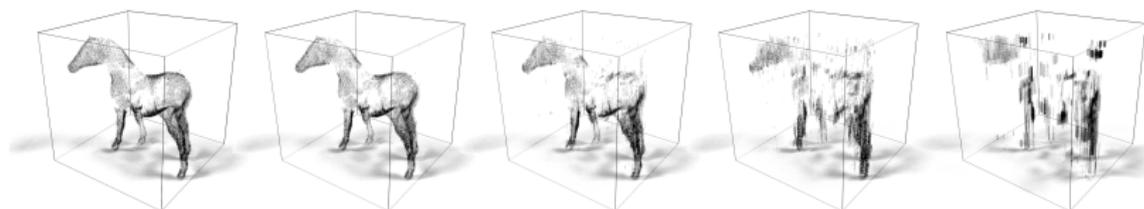
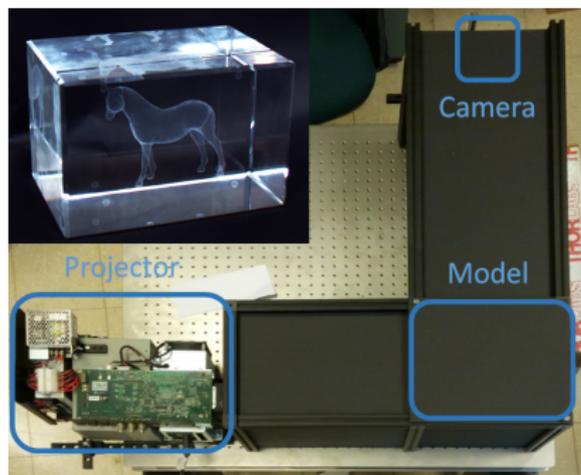
Compressed sensing: reconstruct signal $\mathbf{x} \in \mathbb{R}^n$ given $m \ll n$ measurements

Compressed discrimination: estimate parameter $\theta \in \mathbb{R}^k$ ($k \ll n$) related to \mathbf{x} given $m' < m$ measurements
 m' is insufficient to reconstruct the signal!

Compressed scattering tomography



Compressed scattering tomography



$m : n = 1 : 2$

$1 : 4$

$1 : 8$

$1 : 16$

$1 : 32$

Compressed discrimination

Compressed sensing: reconstruct signal $\mathbf{x} \in \mathbb{R}^n$ given $m \ll n$ measurements

Compressed discrimination: estimate parameter $\theta \in \mathbb{R}^k$ ($k \ll n$) related to \mathbf{x} given $m' < m$ measurements

Compressed discrimination

Compressed sensing: reconstruct signal $\mathbf{x} \in \mathbb{R}^n$ given $m \ll n$ measurements

Random projection: global & linear

Compressed discrimination: estimate parameter $\theta \in \mathbb{R}^k$ ($k \ll n$) related to \mathbf{x} given $m' < m$ measurements

Compressed discrimination

Compressed sensing: reconstruct signal $\mathbf{x} \in \mathbb{R}^n$ given $m \ll n$ measurements

Random projection: global & linear

Compressed discrimination: estimate parameter $\theta \in \mathbb{R}^k$ ($k \ll n$) related to \mathbf{x} given $m' < m$ measurements

CNN: local & non-linear

Compressed discrimination

Compressed sensing: reconstruct signal $x \in \mathbb{R}^n$ given $m \ll n$ measurements

Random projection: global & linear



Compressed discrimination: estimate parameter $\theta \in \mathbb{R}^k$ ($k \ll n$) related to x given $m' < m$ measurements

CNN: local & non-linear

Compressed discrimination

Compressed sensing: reconstruct signal $\mathbf{x} \in \mathbb{R}^n$ given $m \ll n$ measurements

Random projection: global & linear

Compressed discrimination: estimate parameter $\theta \in \mathbb{R}^k$ ($k \ll n$) related to \mathbf{x} given $m' < m$ measurements

CNN: local & non-linear

Conclusion

- Gaussian mean width as a generic data complexity measure in DNN analysis

Conclusion

- Gaussian mean width as a generic data complexity measure in DNN analysis
- DNNs keep important information of the data

Conclusion

- Gaussian mean width as a generic data complexity measure in DNN analysis
- DNNs keep important information of the data
- Random Gaussian weights are good for classifying average data points

Conclusion

- Gaussian mean width as a generic data complexity measure in DNN analysis
- DNNs keep important information of the data
- Random Gaussian weights are good for classifying average data points
- Training improves performance at class boundaries

Conclusion

- Gaussian mean width as a generic data complexity measure in DNN analysis
- DNNs keep important information of the data
- Random Gaussian weights are good for classifying average data points
- Training improves performance at class boundaries
- Deep learning can be viewed as metric learning