

# Perceptron Algorithm

- Assume for simplicity: all  $x_i$  has length 1

1. Start with the all-zeroes weight vector  $w_1 = \mathbf{0}$ , and initialize  $t$  to 1.
  2. Given example  $\mathbf{x}$ , predict positive iff  $w_t \cdot \mathbf{x} > 0$ .
  3. On a mistake, update as follows:
    - Mistake on positive:  $w_{t+1} \leftarrow w_t + \mathbf{x}$ .
    - Mistake on negative:  $w_{t+1} \leftarrow w_t - \mathbf{x}$ .
- $t \leftarrow t + 1$ .

Why is this a reasonable strategy?

# Intuition: correct the current mistake

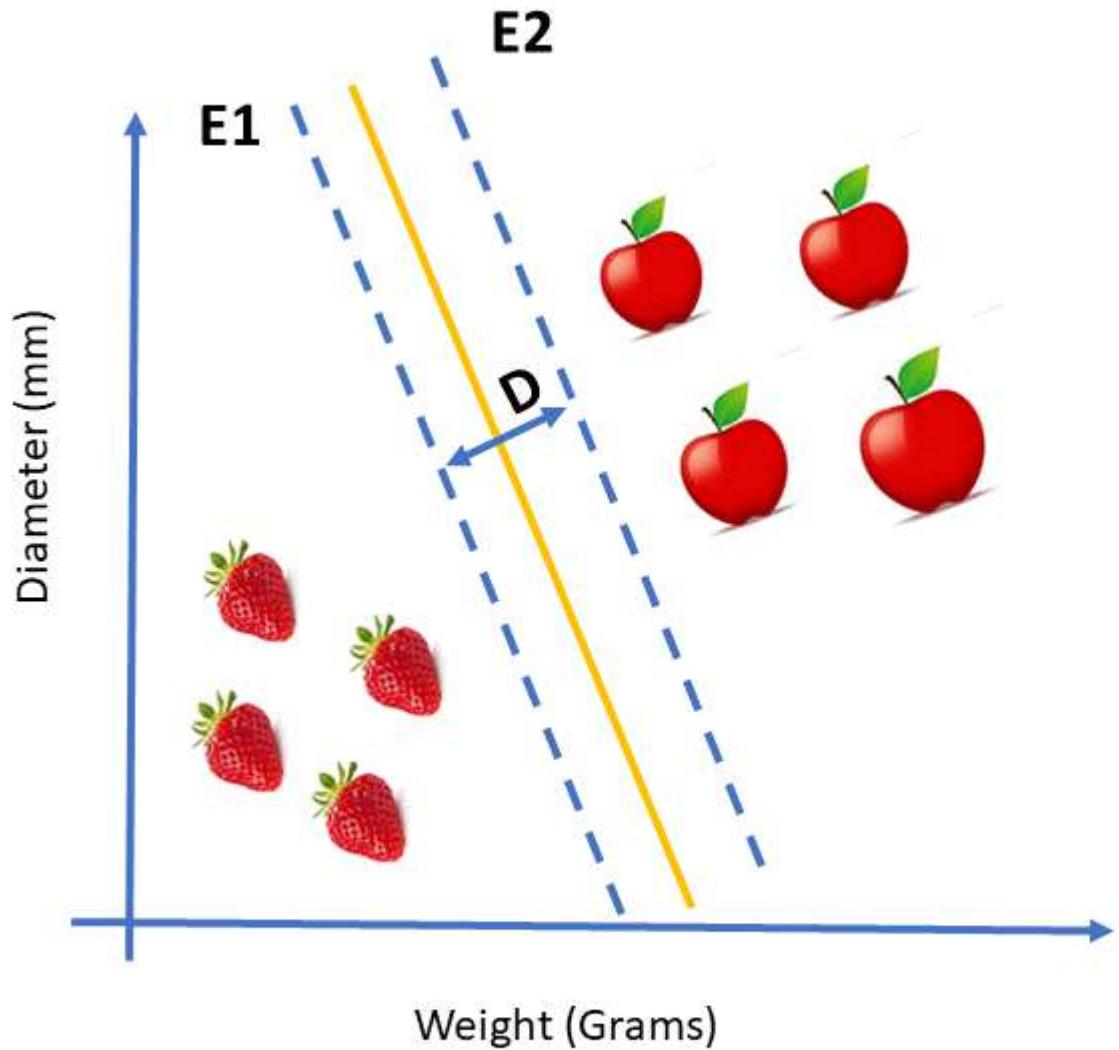
- If mistake on a positive example

$$w_{t+1}^T x = (w_t + x)^T x = w_t^T x + x^T x = w_t^T x + 1$$

- If mistake on a negative example

$$w_{t+1}^T x = (w_t - x)^T x = w_t^T x - x^T x = w_t^T x - 1$$

# Beyond the perceptron: Maximum margin



The one that maximizes the distances to the closest data points from both classes.

This is the hyperplane with maximum margin

Let's define more precisely the concept of margin

**What is the distance of a point  $\mathbf{x}$  to the hyperplane  $\mathcal{H}$ ?**

Consider some point  $\mathbf{x}$ . Let  $\mathbf{d}$  be the vector from  $\mathcal{H}$  to  $\mathbf{x}$  of minimum length. Let  $\mathbf{x}^P$  be the projection of  $\mathbf{x}$  onto  $\mathcal{H}$ . It follows then that:

$$\mathbf{x}^P = \mathbf{x} - \mathbf{d}.$$

$\mathbf{d}$  is parallel to  $\mathbf{w}$ , so  $\mathbf{d} = \alpha \mathbf{w}$  for some  $\alpha \in \mathbb{R}$ .

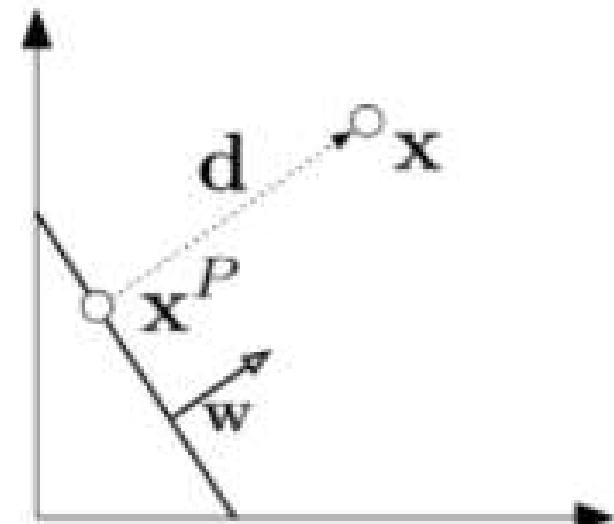
$\mathbf{x}^P \in \mathcal{H}$  which implies  $\mathbf{w}^T \mathbf{x}^P + b = 0$

therefore  $\mathbf{w}^T \mathbf{x}^P + b = \mathbf{w}^T(\mathbf{x} - \mathbf{d}) + b = \mathbf{w}^T(\mathbf{x} - \alpha \mathbf{w}) + b = 0$

$$\text{which implies } \alpha = \frac{\mathbf{w}^T \mathbf{x} + b}{\mathbf{w}^T \mathbf{w}}$$

The length of  $\mathbf{d}$ :

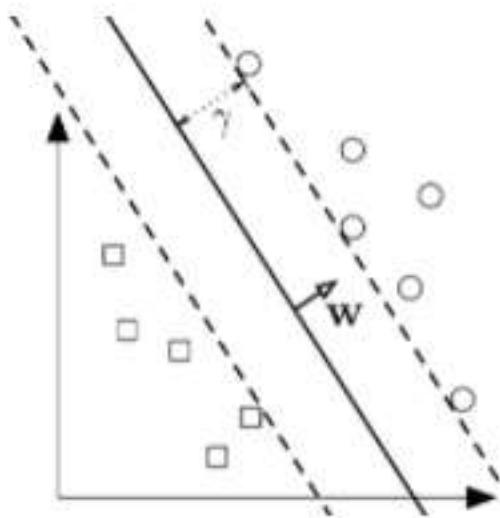
$$\|\mathbf{d}\|_2 = \sqrt{\mathbf{d}^T \mathbf{d}} = \sqrt{\alpha^2 \mathbf{w}^T \mathbf{w}} = \frac{|\mathbf{w}^T \mathbf{x} + b|}{\sqrt{\mathbf{w}^T \mathbf{w}}} = \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|_2}$$



Margin of  $\mathcal{H}$  with respect to  $D$ :  $\gamma(\mathbf{w}, b) = \min_{\mathbf{x} \in D} \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|_2}$

We can formulate our search for the maximum margin separating hyperplane as a constrained optimization problem. The objective is to maximize the margin under the constraints that all data points must lie on the correct side of the hyperplane:

$$\max_{\mathbf{w}, b} \underbrace{\gamma(\mathbf{w}, b)}_{\text{maximize margin}} \text{ such that } \underbrace{\forall i y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 0}_{\text{separating hyperplane}}$$



If we plug in the definition of  $\gamma$  we obtain:

$$\underbrace{\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \min_{\mathbf{x}_i \in D} |\mathbf{w}^T \mathbf{x}_i + b|}_{\gamma(\mathbf{w}, b)} \quad s.t. \quad \underbrace{\forall i y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 0}_{\text{separating hyperplane}}$$

*maximize margin*

Because the hyperplane is scale invariant, we can fix the scale of  $\mathbf{w}, b$  anyway we want. Let's be clever about it, and choose it such that

$$\min_{\mathbf{x} \in D} |\mathbf{w}^T \mathbf{x} + b| = 1.$$

We can add this re-scaling as an equality constraint. Then our objective becomes:

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \cdot 1 = \min_{\mathbf{w}, b} \|\mathbf{w}\|_2 = \min_{\mathbf{w}, b} \mathbf{w}^\top \mathbf{w}$$

(Where we made use of the fact  $f(z) = z^2$  is a monotonically increasing function for  $z \geq 0$  and  $\|\mathbf{w}\| \geq 0$ ; i.e. the  $\mathbf{w}$  that maximizes  $\|\mathbf{w}\|_2$  also maximizes  $\mathbf{w}^\top \mathbf{w}$ .)

The new optimization problem becomes:

$$\begin{aligned} & \min_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w} \\ \text{s.t. } & \forall i, y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 0 \\ & \min_i |\mathbf{w}^T \mathbf{x}_i + b| = 1 \end{aligned}$$

These constraints are still hard to deal with, however luckily we can show that (for the optimal solution) they are equivalent to a much simpler formulation. (Makes sure you know how to prove that the two sets of constraints are equivalent.)

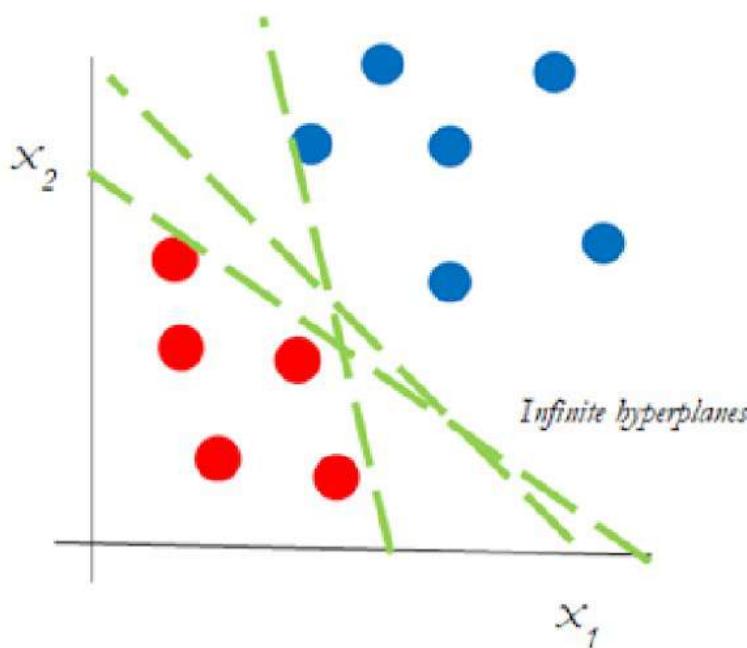
$$\begin{aligned} & \min_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w} \\ \text{s.t. } & \forall i y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \end{aligned}$$

# Support vectors

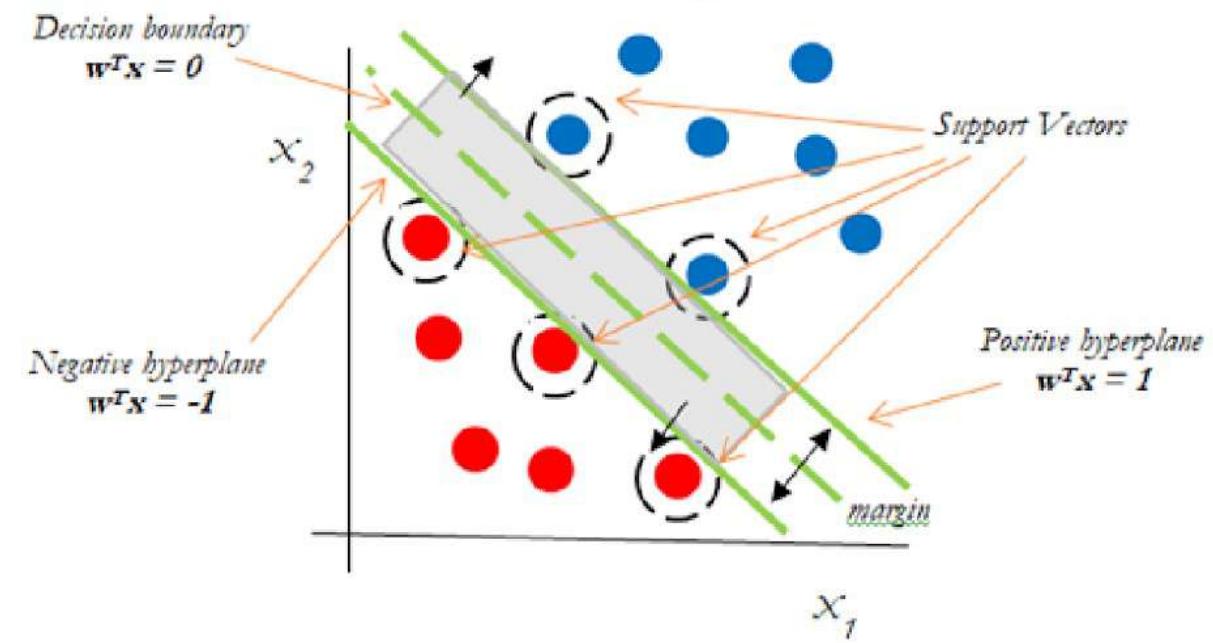
For the optimal  $\mathbf{w}, b$  pair, some training points will have tight constraints, i.e.

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1.$$

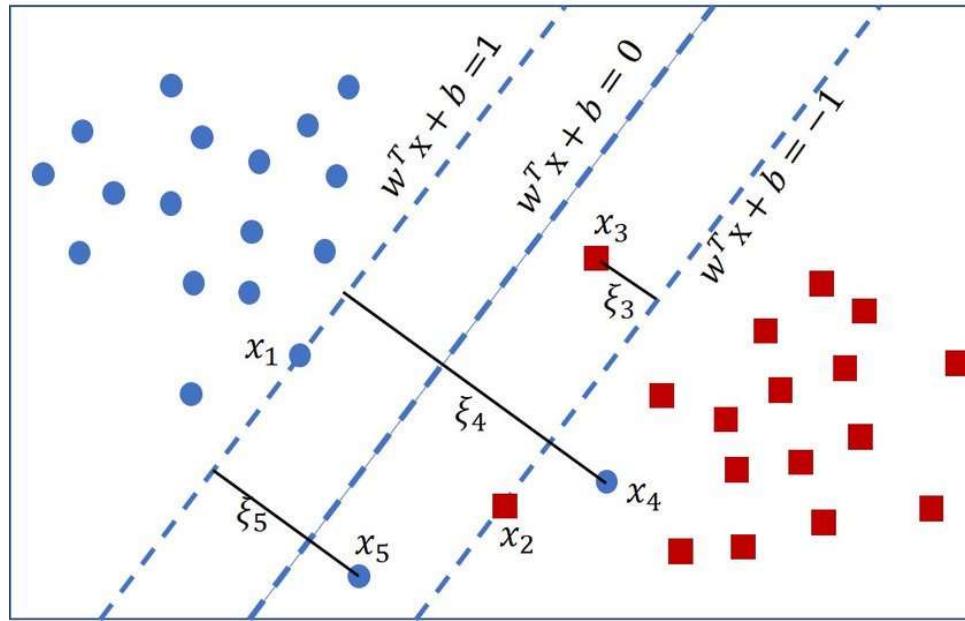
*Infinite Hyperplanes*



*Maximum Margin Classifier*



# Maximal margin classification and limitations



If the data is low dimensional it is often the case that there is no separating hyperplane between the two classes. In this case, there is no solution to the optimization problems stated above. We can fix this by allowing the constraints to be violated ever so slightly with the introduction of slack variables:

$$\begin{aligned} & \min_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i \\ \text{s.t. } & \forall i y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \\ & \forall i \xi_i \geq 0 \end{aligned}$$

The slack variable  $\xi_i$  allows the input  $\mathbf{x}_i$  to be closer to the hyperplane (or even be on the wrong side), but there is a penalty in the objective function for such "slack". If  $C$  is very large, the SVM becomes very strict and tries to get all points to be on the right side of the hyperplane. If  $C$  is very small, the SVM becomes very loose and may "sacrifice" some points to obtain a simpler (i.e. lower  $\|\mathbf{w}\|_2^2$ ) solution.

# New optimization problem

Let us consider the value of  $\xi_i$  for the case of  $C \neq 0$ . Because the objective will always try to minimize  $\xi_i$  as much as possible, the equation must hold as an *equality* and we have:

$$\xi_i = \begin{cases} 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) & \text{if } y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1 \\ 0 & \text{if } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \end{cases}$$

This is equivalent to the following closed form:

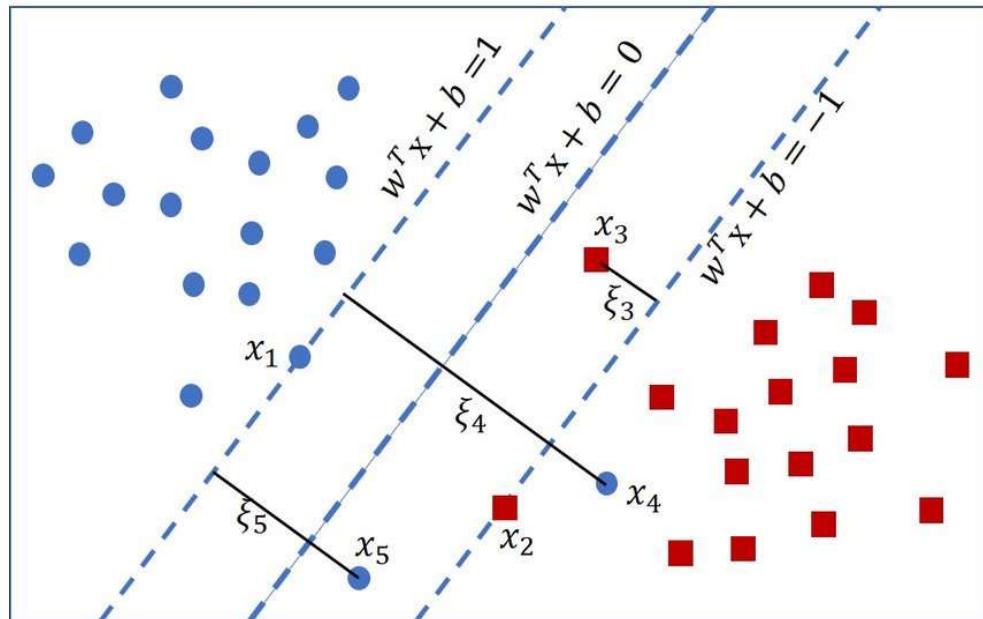
$$\xi_i = \max(1 - y_i(\mathbf{w}^T \mathbf{x}_i + b), 0).$$

If we plug this closed form into the objective of our SVM optimization problem, we obtain the following *unconstrained* version as loss function and regularizer:

$$\min_{\mathbf{w}, b} \underbrace{\frac{\mathbf{w}^T \mathbf{w}}{2}}_{l_2\text{-regularizer}} + C \sum_{i=1}^n \underbrace{\max [1 - y_i(\mathbf{w}^T \mathbf{x}_i + b), 0]}_{hinge\text{-loss}}$$

This formulation allows us to optimize the SVM parameters ( $\mathbf{w}, b$ ) just like logistic regression (e.g. through gradient descent). The only difference is that we have the **hinge-loss** instead of the **logistic loss**.

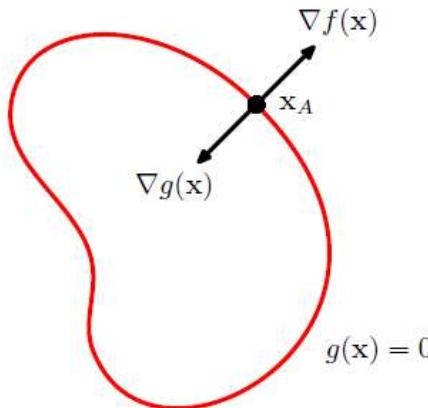
# Solving for the new optimization problem



slack penalty

$$\begin{aligned} & \min_{\substack{w \in \mathbb{R}^d \\ \xi \geq 0}} && \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i \\ & \text{s.t.} && y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i \in \{1, \dots, n\} \end{aligned}$$

# Lagrange Multipliers



Consider the problem:

$$\begin{aligned} & \max_x f(\mathbf{x}) \\ \text{s.t. } & g(\mathbf{x}) = 0 \end{aligned}$$

This is because  
on the curves  $g$  is  
constant

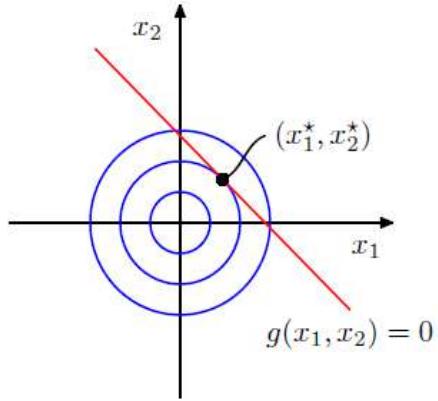
- Points on  $g(\mathbf{x}) = 0$  must have  $\nabla g(\mathbf{x})$  normal to surface
- A **stationary point** must have no change in  $f$  in the direction of the constraint surface, so  $\nabla f(\mathbf{x})$  must also be normal to the surface.
  - So there must be some  $\lambda \neq 0$  such that  $\nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$
- Define **Lagrangian**:

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

This is because they  
are the same vector up  
to a multiplicative  
constant

- Stationary points of  $L(\mathbf{x}, \lambda)$  have  $\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda) = \nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$  and  $\nabla_{\lambda} L(\mathbf{x}, \lambda) = g(\mathbf{x}) = 0$

## Lagrange Multipliers Example



- Consider the problem

$$\begin{aligned} \max_{\mathbf{x}} f(x_1, x_2) &= 1 - x_1^2 - x_2^2 \\ \text{s.t.} \quad g(x_1, x_2) &= x_1 + x_2 - 1 = 0 \end{aligned}$$

- Lagrangian:

$$L(\mathbf{x}, \lambda) = 1 - x_1^2 - x_2^2 + \lambda(x_1 + x_2 - 1)$$

- Stationary points require:

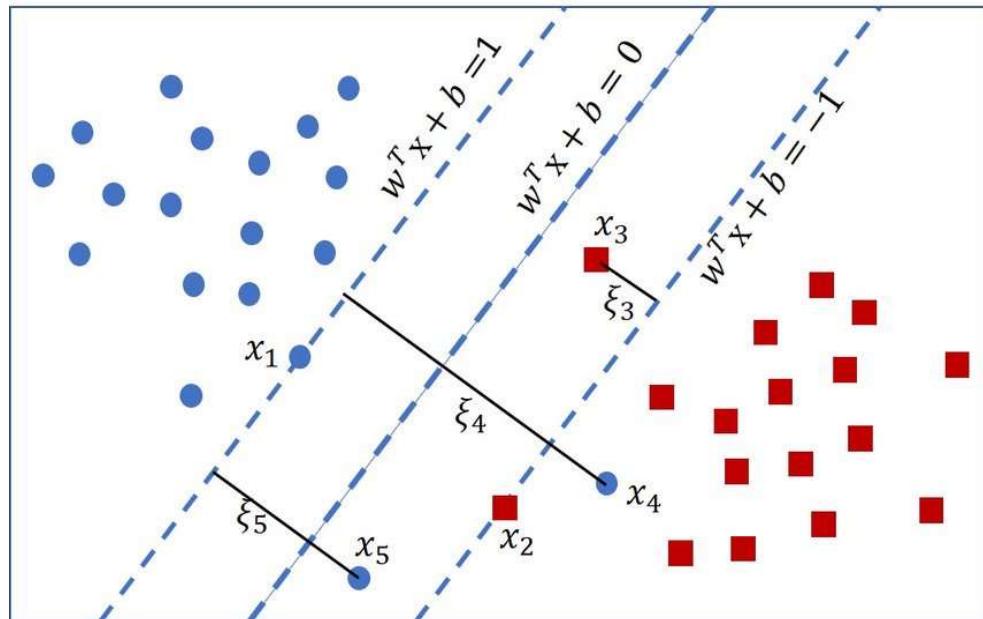
$$\frac{\partial L}{\partial x_1} = -2x_1 + \lambda = 0$$

$$\frac{\partial L}{\partial x_2} = -2x_2 + \lambda = 0$$

$$\frac{\partial L}{\partial \lambda} = x_1 + x_2 - 1 = 0$$

- So stationary point is  $(x_1^*, x_2^*) = (\frac{1}{2}, \frac{1}{2})$ ,  $\lambda = 1$

# Solving for the new optimization problem



slack penalty

$$\begin{aligned} & \min_{\substack{w \in \mathbb{R}^d \\ \xi \geq 0}} \quad \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i \in \{1, \dots, n\} \end{aligned}$$

The Lagrangian:  $L(w, b, \xi, \alpha, \lambda)$

$$= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i \left( 1 - y_i (w^\top x_i + b) - \xi_i \right) + \sum_{i=1}^n \lambda_i (-\xi_i)$$

with dual variable constraints

$$\alpha_i \geq 0, \quad \lambda_i \geq 0.$$

**Minimize wrt the primal variables  $w$ ,  $b$ , and  $\xi$ .**

Derivative wrt  $w$ :

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^n \alpha_i y_i x_i = 0 \quad w = \sum_{i=1}^n \alpha_i y_i x_i.$$

Derivative wrt  $b$ :

$$\frac{\partial L}{\partial b} = \sum_i y_i \alpha_i = 0.$$

Derivative wrt  $\xi_i$ :

$$\frac{\partial L}{\partial \xi_i} = C - \alpha_i - \lambda_i = 0 \quad \alpha_i = C - \lambda_i.$$

Noting that  $\lambda_i \geq 0$ ,

$$\alpha_i \leq C.$$

$$L(w, b, \xi, \alpha, \beta) = \frac{1}{2} w^\top w + C \sum_i \xi_i$$

$$- \sum_i \alpha_i (y_i (w^\top x_i + b) - 1 + \xi_i) - \sum_i \beta_i \xi_i$$


---

Gradients

$$\nabla_w L = w - \sum_i \alpha_i y_i x_i = 0$$

$$\nabla_b L = - \sum_i \alpha_i y_i = 0$$

$$\nabla_\xi L = C - \alpha - \beta = 0$$

Consequences

$$w = \sum_i \alpha_i y_i x_i$$

$$\sum_i \alpha_i y_i = 0$$

$$\alpha = C - \beta$$

$$\beta = C - \alpha$$

**w is a linear combination of the training points!**

Rewriting the Lagrangian...

$$\begin{aligned} L(\alpha, \lambda) &= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i \left( 1 - y_i (w^\top x_i + b) - \xi_i \right) \\ &\quad + \sum_{i=1}^n \lambda_i (-\xi_i) \\ &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^\top x_j + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^\top x_j \\ &\quad - b \underbrace{\sum_{i=1}^m \alpha_i y_i}_{0} + \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i \xi_i - \sum_{i=1}^m \underbrace{(C - \alpha_i)}_{\lambda_i} \xi_i \\ &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^\top x_j. \end{aligned}$$

$$L(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^\top x_j,$$

subject to the constraints

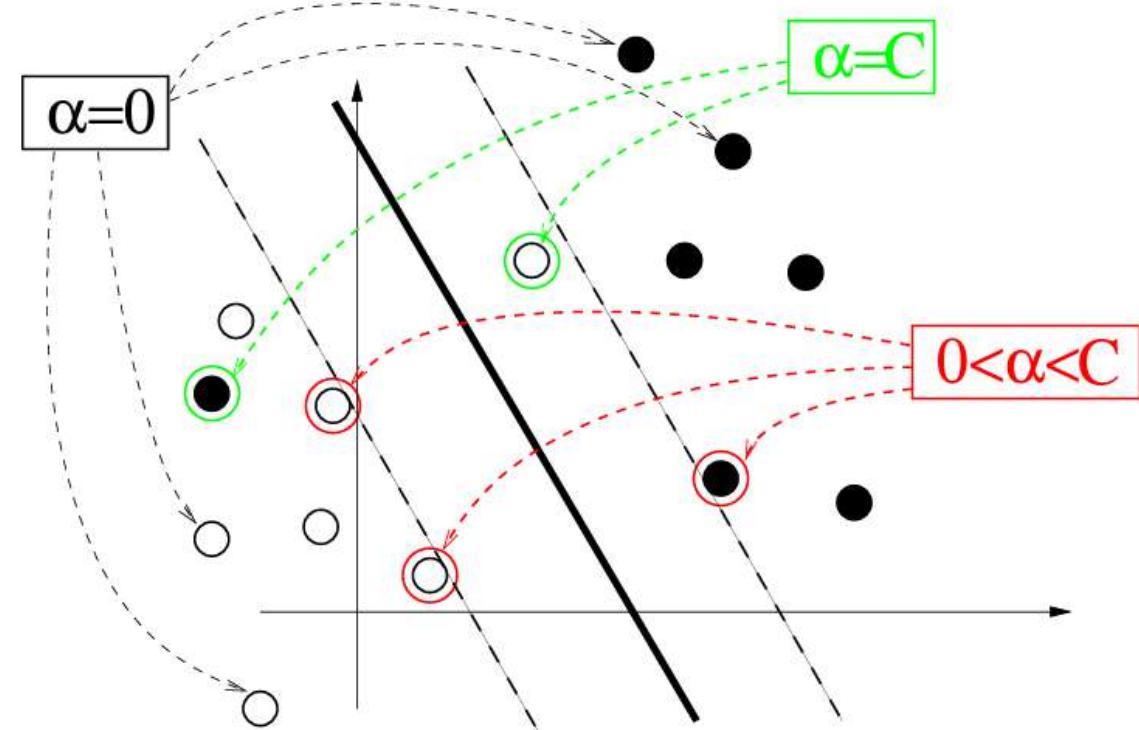
$$0 \leq \alpha_i \leq C, \quad \sum_{i=1}^n y_i \alpha_i = 0$$

# Support vectors

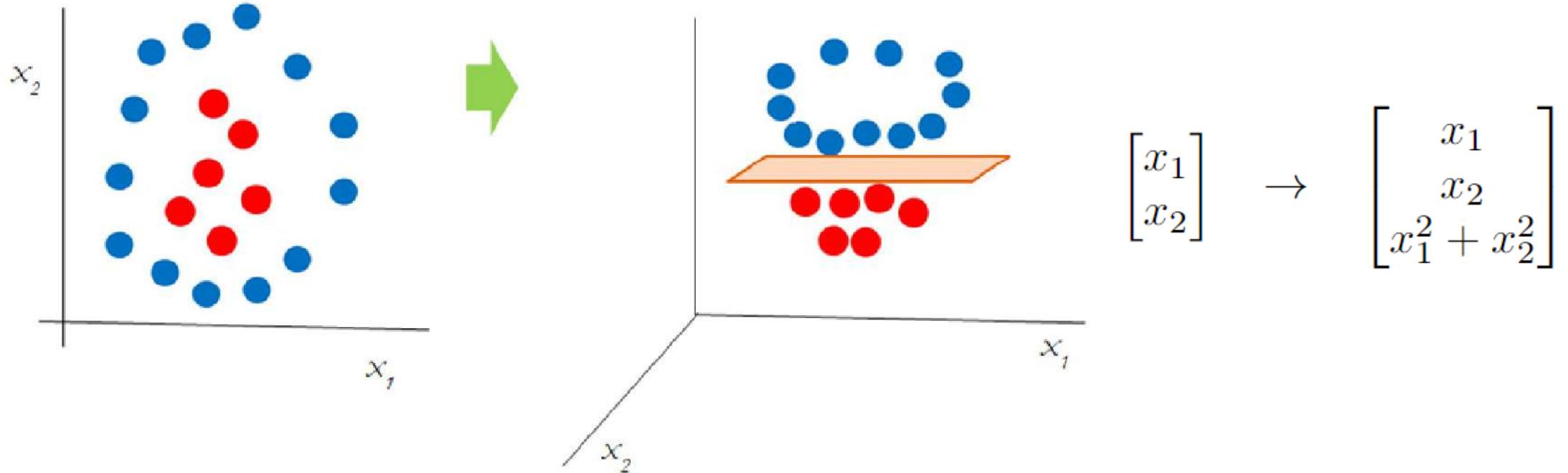
- The solution  $w$  can be written as a linear combination of the data
- Only those data for which  $0 \leq \alpha_i \leq C$ ,  $\sum_{i=1}^n y_i \alpha_i = 0$  are counted for the solution

# Support vectors

$$\begin{array}{ll}\min_{\substack{w \in \mathbb{R}^d \\ \xi \geq 0}} & \frac{1}{2} w^\top w + C \sum_{i=1}^n \xi_i \\ \text{s.t.} & y_i(w^\top x_i + b) \geq 1 - \xi_i\end{array}$$



$w$  is given by a linear combination only of the points that satisfy the inequality!



$$f(x) = w^T x \rightarrow f(x) = w^T \phi(x)$$

- A problem that is not linearly separable in 2 dim becomes separable in 3 dim.
- We know a good coordinate change; what's the problem with this?

# Polynomial kernel

For  $x, z \in [0, 1]$  and  $0 < \alpha < 1$

$$k(x, z) = \frac{1}{1 - \alpha^2 xz}$$

**Proof**

$$\frac{1}{1 - \alpha^2 xz} = \sum_{s=0}^{\infty} (\alpha^2 xz)^s = \Phi(x)^\top \Phi(z)$$

with

$$\Phi(x) = (1, \alpha x, \alpha^2 x^2, \alpha^3 x^3, \dots)$$

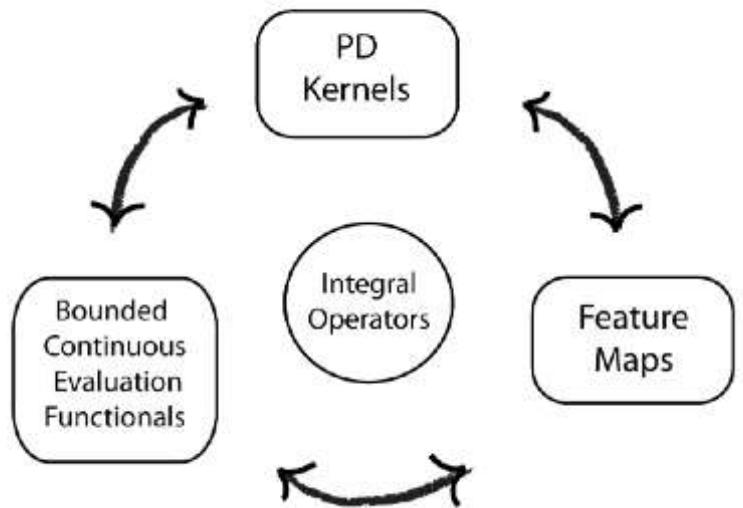
# Gaussian kernel

$$\begin{aligned} K(x_i, x_j) &= \exp\left(-\frac{1}{2\sigma^2}\|x_i - x_j\|^2\right) & \sigma = 1/\sqrt{2} \\ &= \exp(-\|x_i - x_j\|^2) \\ &= \exp(-x_i^\top x_i) \exp(-x_j^\top x_j) \exp(2x_i^\top x_j) & \exp(z) = \sum_{n=0}^{\infty} \frac{z^n}{n!} \\ &= \exp(-x_i^\top x_i) \exp(-x_j^\top x_j) \sum_{n=0}^{\infty} \frac{2^n (x_i^\top x_j)^n}{n!} & \text{order-n polynomial kernel* } \Phi^n(x) \end{aligned}$$

$$\Phi^{\text{rbf}} = \exp(-x^\top x) [\Phi^1(x)^\top, \Phi^2(x)^\top, \dots, \Phi^\infty(x)^\top]^\top$$

# Class 3: Math of kernel methods

- Aronszajn theorem (positive definite maps and features)
- Reproducing property. Reproducing kernel Hilbert spaces (RKHS)
- Pointwise continuity



Some slides from : Julien Mairal

<https://members.cbio.mines-paristech.fr/~jvert/svn/kernelcourse/slides/master2017/master2017.pdf>

# Can we start from the kernel instead of features? What defines a good kernel?

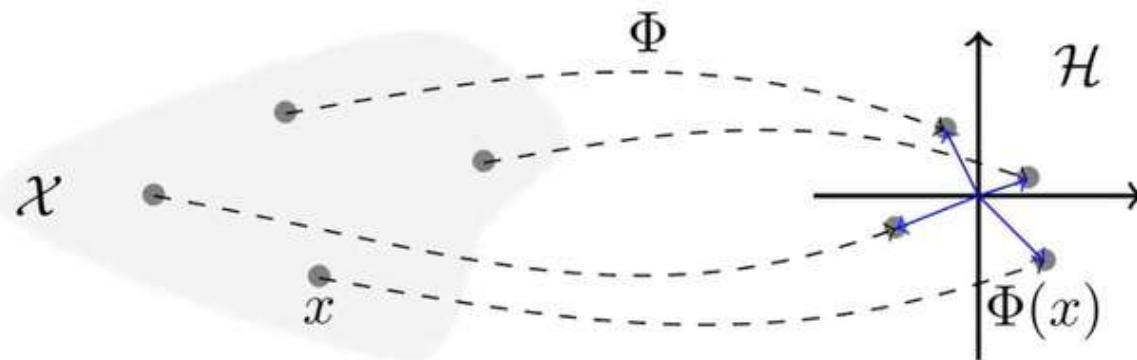
Theorem (Aronszajn, 1950)

$K$  is a p.d. kernel on the set  $\mathcal{X}$  if and only if there exists a Hilbert space  $\mathcal{H}$  and a mapping

$$\Phi : \mathcal{X} \mapsto \mathcal{H}$$

such that, for any  $\mathbf{x}, \mathbf{x}'$  in  $\mathcal{X}$ :

↳ 
$$K(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle_{\mathcal{H}} .$$



# Def1: Hilbert space

- An **inner product** on an  $\mathbb{R}$ -vector space  $\mathcal{H}$  is a mapping  $(f, g) \mapsto \langle f, g \rangle_{\mathcal{H}}$  from  $\mathcal{H}^2$  to  $\mathbb{R}$  that is **bilinear, symmetric** and such that  $\langle f, f \rangle_{\mathcal{H}} > 0$  for all  $f \in \mathcal{H} \setminus \{0\}$ .
- A vector space endowed with an inner product is called **pre-Hilbert**. It is endowed with a **norm** defined as  $\|f\|_{\mathcal{H}} = \langle f, f \rangle_{\mathcal{H}}^{\frac{1}{2}}$ .
- A **Cauchy sequence**  $(f_n)_{n \geq 0}$  is a sequence whose elements become progressively arbitrarily close to each other:

$$\lim_{N \rightarrow +\infty} \sup_{n, m \geq N} \|f_n - f_m\|_{\mathcal{H}} = 0.$$

- A **Hilbert space** is a pre-Hilbert space **complete** for the norm  $\|\cdot\|_{\mathcal{H}}$ . That is, any Cauchy sequence in  $\mathcal{H}$  converges in  $\mathcal{H}$ .

Completeness is necessary to keep “good” convergence properties of Euclidean spaces in an infinite-dimensional context.

# Def2: Positive definite kernel

Definition.  $\Omega$ : set.  $k: \Omega \times \Omega \rightarrow \mathbf{R}$  is a **positive definite kernel** if

- 1) (symmetry)  $k(x, y) = k(y, x)$
- 2) (positivity) for arbitrary  $x_1, \dots, x_n \in \Omega$

$$\begin{pmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{pmatrix}$$
 is positive semidefinite,  
(Gram matrix)

i.e.,  $\sum_{i,j=1}^n c_i c_j k(x_i, x_j) \geq 0$  for any  $c_i \in \mathbf{R}$   $c^T K c$

### Theorem (Aronszajn, 1950)

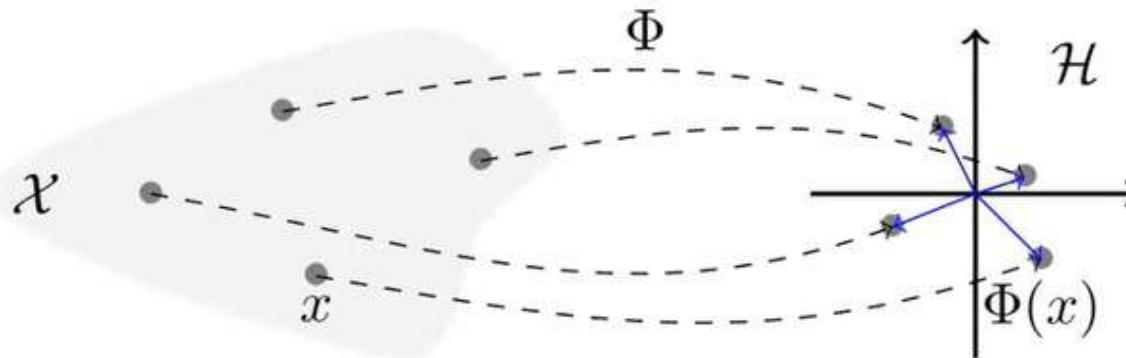
$K$  is a p.d. kernel on the set  $\mathcal{X}$  if and only if there exists a Hilbert space  $\mathcal{H}$  and a mapping

$$\Phi : \mathcal{X} \mapsto \mathcal{H}$$

such that, for any  $x, x'$  in  $\mathcal{X}$ :



$$K(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}} .$$



The mapping is not unique!

**Example 37.** Consider  $\mathcal{X} = \mathbb{R}$ , and

$$k(x, y) = xy = \begin{bmatrix} \frac{x}{\sqrt{2}} & \frac{x}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \frac{y}{\sqrt{2}} \\ \frac{y}{\sqrt{2}} \end{bmatrix},$$

where we defined the feature maps  $\phi(x) = x$  and  $\tilde{\phi}(x) = \begin{bmatrix} \frac{x}{\sqrt{2}} & \frac{x}{\sqrt{2}} \end{bmatrix}$ , and where the feature spaces are respectively,  $\mathcal{H} = \mathbb{R}$ , and  $\tilde{\mathcal{H}} = \mathbb{R}^2$ .

What is a good definition of the restriction of the  $\mathcal{H}$  space?

# The reproducing property guarantees the uniqueness

## Reproducing kernels Hilbert spaces (RKHS)

---

$$f(x) = \int f(y)K(x,y)dy$$

### Definition

Let  $\mathcal{X}$  be a set and  $\mathcal{H} \subset \mathbb{R}^{\mathcal{X}}$  be a class of functions forming a (real) Hilbert space with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ . The function  $K : \mathcal{X}^2 \mapsto \mathbb{R}$  is called a reproducing kernel (r.k.) of  $\mathcal{H}$  if

- ①  $\mathcal{H}$  contains all functions of the form

$$\forall \mathbf{x} \in \mathcal{X}, \quad K_{\mathbf{x}} : \mathbf{t} \mapsto K(\mathbf{x}, \mathbf{t}) .$$

- 
- ② For every  $\mathbf{x} \in \mathcal{X}$  and  $f \in \mathcal{H}$  the reproducing property holds:

$$f(\mathbf{x}) = \langle f, K_{\mathbf{x}} \rangle_{\mathcal{H}} .$$

If a r.k. exists, then  $\mathcal{H}$  is called a reproducing kernel Hilbert space (RKHS).

## First a definition (from matrices to operators):

Fix a symmetric function  $k : \mathcal{X}^2 \rightarrow \mathbb{R}$  on a compact set  $\mathcal{X} \subset \mathbb{R}^d$ , and consider the integral operator  $T_k : L_2(\mathcal{X}) \rightarrow L_2(\mathcal{X})$  defined as

$$T_k f(\cdot) = \int_{\mathcal{X}} k(\cdot, x) f(x) dx.$$

We say  $T_k$  is positive semidefinite if, for all  $f \in L_2(\mathcal{X})$ ,  
 $\langle f, T_k f \rangle_{L_2(\mathcal{X})} \geq 0$ , that is,

$$\int_{\mathcal{X}^2} k(u, v) f(u) f(v) du dv \geq 0.$$

## Mercer's Theorem

**Theorem:** If  $k$  is continuous and  $T_k$  is positive semidefinite, then  $T_k$  has eigenfunctions  $\psi_i \in L_2(\mathcal{X})$  (say  $\|\psi_i\|_{L_2} = 1$ ) with eigenvalues  $\lambda_i \geq 0$ , and for all  $u, v \in \mathcal{X}$ , we can write

$$k(u, v) = \sum_{i=1}^{\infty} \lambda_i \psi_i(u) \psi_i(v).$$

Furthermore, this series converges uniformly.

- A positive semidefinite operator defines a kernel
- Eigenfunctions play the role of features/basis (not normalized)

We skip the demonstration (Spectral Theorem)

## Reproducing property

Under certain conditions (Mercer's theorem and extensions), we can write

$$k(x, x') = \sum_{i=1}^{\infty} \lambda_i e_i(x) e_i(x'), \quad \int_{\mathcal{X}} e_i(x) e_j(x) d\mu(x) = \begin{cases} 1 & i = j \\ 0 & i \neq j. \end{cases}$$

where this sum is guaranteed to converge whatever the  $x$  and  $x'$ .

Infinite dimensional feature map:  $\phi(x) = \begin{bmatrix} \vdots \\ \sqrt{\lambda_i} e_i(x) \\ \vdots \end{bmatrix} \in \ell_2.$

Reproducing property  
when we have pd: more  
concrete

$$f(x) = \sum_i \alpha_i e_i(x)$$

$$K(x, x') = \sum_j \lambda_j e_j(x) e_j(x')$$

$$\begin{aligned}\langle f, K_x \rangle &= \int K(x, x') f(x') dx' \\ &= \int \left( \sum_i \lambda_i e_i(x) e_i(x') \right) f(x') dx' \\ &= \sum_i \left( \int \sqrt{\lambda_i} e_i(x') f(x') dx' \right) \sqrt{\lambda_i} e_i(x) \\ &= \sum_i \alpha_i e_i(x) = f(x)\end{aligned}$$

How does the reproducing property guarantees uniqueness?  
We divide the demonstration into steps

Kernel  $\leftrightarrow$  PD  $\leftrightarrow$  r.k.  $\leftrightarrow$  RKHS

Kernel  $\leftrightarrow$  PD  $\leftrightarrow$  r.k.  $\leftrightarrow$  RKHS

Mercer theorem

PD  $\leftrightarrow$  r.k.  $\leftrightarrow$  RKHS

A function  $k : X \times X \rightarrow \mathbb{R}$  defines a positive definite kernel if and only if there exists a Reproducing Kernel Hilbert Space (RKHS)  $\mathcal{H}$  of functions on  $X$  such that:

1. For every  $x \in X$ , the function  $k_x(\cdot) = k(x, \cdot)$  belongs to  $\mathcal{H}$ .
2. The reproducing property holds: for every  $f \in \mathcal{H}$  and every  $x \in X$ ,

$$f(x) = \langle f, k_x \rangle_{\mathcal{H}}.$$

**( $\Rightarrow$ ) If  $k$  is a positive definite kernel, then it induces an RKHS with the reproducing property.**

Construct an inner product on the space of functions generated by the kernel  $k$  and show that this construction leads to an RKHS that inherently satisfies the reproducing property.

1. **Construct the RKHS:** Given a positive definite kernel  $k : X \times X \rightarrow \mathbb{R}$ , we define the linear space  $\mathcal{H}_0$  consisting of finite linear combinations of kernel functions:

$$\mathcal{H}_0 = \left\{ f = \sum_{i=1}^n c_i k_{x_i} : c_i \in \mathbb{R}, x_i \in X, n \in \mathbb{N} \right\}.$$

Define an inner product on this space by:

$$\left\langle \sum_{i=1}^n c_i k_{x_i}, \sum_{j=1}^m d_j k_{y_j} \right\rangle_{\mathcal{H}_0} = \sum_{i=1}^n \sum_{j=1}^m c_i d_j k(x_i, y_j).$$

2. **Completion to an RKHS:** The space  $\mathcal{H}_0$  can be completed to form a Hilbert space  $\mathcal{H}$ . By construction, every function in this Hilbert space can be expressed (or approximated arbitrarily closely) as a linear combination of the kernel functions  $k_x$ .
3. **Reproducing Property:** For each  $x \in X$  and any function  $f = \sum_{i=1}^n c_i k_{x_i} \in \mathcal{H}_0$ ,

$$f(x) = \sum_{i=1}^n c_i k(x_i, x) = \left\langle \sum_{i=1}^n c_i k_{x_i}, k_x \right\rangle_{\mathcal{H}_0}.$$

Since  $k_x \in \mathcal{H}$  by construction, this property extends to the entire RKHS  $\mathcal{H}$ . Therefore,  $k$  satisfies the reproducing property.

The reproducing property emerges naturally from how the inner product was defined, due to the kernel's structure.  
This is the key of the proof: defining the space and inner product using the kernel, then completing this space to obtain the RKHS where the reproducing property holds.

**( $\Leftarrow$ ) If there exists an RKHS with the reproducing property, then  $k$  is a positive definite kernel.**

1. **Existence of the Kernel Function:** Suppose  $\mathcal{H}$  is an RKHS of functions on  $X$  with the reproducing property. For every  $x \in X$ , there exists a function  $k_x \in \mathcal{H}$  such that for any  $f \in \mathcal{H}$ ,

$$f(x) = \langle f, k_x \rangle_{\mathcal{H}}.$$

2. **Define the Kernel:** Define the function  $k : X \times X \rightarrow \mathbb{R}$  by

$$k(x, y) = \langle k_y, k_x \rangle_{\mathcal{H}}.$$

3. **Positive Definiteness:** To show that  $k$  is a positive definite kernel, take any finite set of points  $\{x_1, x_2, \dots, x_n\} \subset X$  and any real coefficients  $c_1, c_2, \dots, c_n$ . Consider the linear combination in the RKHS:

$$f = \sum_{i=1}^n c_i k_{x_i}.$$

Then, the norm squared of  $f$  in the RKHS is given by:

$$\|f\|_{\mathcal{H}}^2 = \left\langle \sum_{i=1}^n c_i k_{x_i}, \sum_{j=1}^n c_j k_{x_j} \right\rangle_{\mathcal{H}} = \sum_{i=1}^n \sum_{j=1}^n c_i c_j \langle k_{x_i}, k_{x_j} \rangle_{\mathcal{H}} = \sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j).$$

Since  $\|f\|_{\mathcal{H}}^2 \geq 0$  (as it is the squared norm of a function in a Hilbert space), it follows that

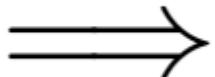
$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \geq 0. \quad c^T K c$$

Therefore,  $k$  is a positive definite kernel.

# Uniqueness of RKHS and RK

## Theorem

- If  $\mathcal{H}$  is a RKHS, then it has a unique r.k.
- Conversely, a function  $K$  can be the r.k. of at most one RKHS.



If a r.k. exists then it is unique

Let  $K$  and  $K'$  be two r.k. of a RKHS  $\mathcal{H}$ . Then for any  $\mathbf{x} \in \mathcal{X}$ :

$$\begin{aligned}\|K_{\mathbf{x}} - K'_{\mathbf{x}}\|_{\mathcal{H}}^2 &= \langle K_{\mathbf{x}} - K'_{\mathbf{x}}, K_{\mathbf{x}} - K'_{\mathbf{x}} \rangle_{\mathcal{H}} \\ &= \langle K_{\mathbf{x}} - K'_{\mathbf{x}}, K_{\mathbf{x}} \rangle_{\mathcal{H}} - \langle K_{\mathbf{x}} - K'_{\mathbf{x}}, K'_{\mathbf{x}} \rangle_{\mathcal{H}} \\ &= K_{\mathbf{x}}(\mathbf{x}) - K'_{\mathbf{x}}(\mathbf{x}) - K_{\mathbf{x}}(\mathbf{x}) + K'_{\mathbf{x}}(\mathbf{x}) \\ &= 0.\end{aligned}$$

This shows that  $K_{\mathbf{x}} = K'_{\mathbf{x}}$  as functions, i.e.,  $K_{\mathbf{x}}(\mathbf{y}) = K'_{\mathbf{x}}(\mathbf{y})$  for any  $\mathbf{y} \in \mathcal{X}$ . In other words,  $\mathbf{K}=\mathbf{K}'$ .

□

Another way of characterizing RKHS is:  
reproducing property and pointwise continuity

If  $\mathcal{H}$  is a RKHS then  $f \mapsto f(\mathbf{x})$  is continuous

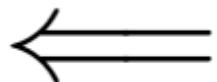
If a r.k.  $K$  exists, then for any  $(\mathbf{x}, f) \in \mathcal{X} \times \mathcal{H}$ :

$$\begin{aligned}|f(\mathbf{x})| &= |\langle f, K_{\mathbf{x}} \rangle_{\mathcal{H}}| & |\langle \mathbf{u}, \mathbf{v} \rangle|^2 &\leq \langle \mathbf{u}, \mathbf{u} \rangle \cdot \langle \mathbf{v}, \mathbf{v} \rangle \\&\leq \|f\|_{\mathcal{H}} \cdot \|K_{\mathbf{x}}\|_{\mathcal{H}} && \text{(Cauchy-Schwarz)} \\&\leq \|f\|_{\mathcal{H}} \cdot K(\mathbf{x}, \mathbf{x})^{\frac{1}{2}},\end{aligned}$$

because  $\|K_{\mathbf{x}}\|_{\mathcal{H}}^2 = \langle K_{\mathbf{x}}, K_{\mathbf{x}} \rangle_{\mathcal{H}} = K(\mathbf{x}, \mathbf{x})$ . Therefore  $f \in \mathcal{H} \mapsto f(\mathbf{x}) \in \mathbb{R}$  is a continuous linear mapping. □

Reproducing property implies the fact that is pointwise continuous

# Reproducing property and pointwise continuity



If  $f \mapsto f(\mathbf{x})$  is continuous then  $\mathcal{H}$  is a RKHS

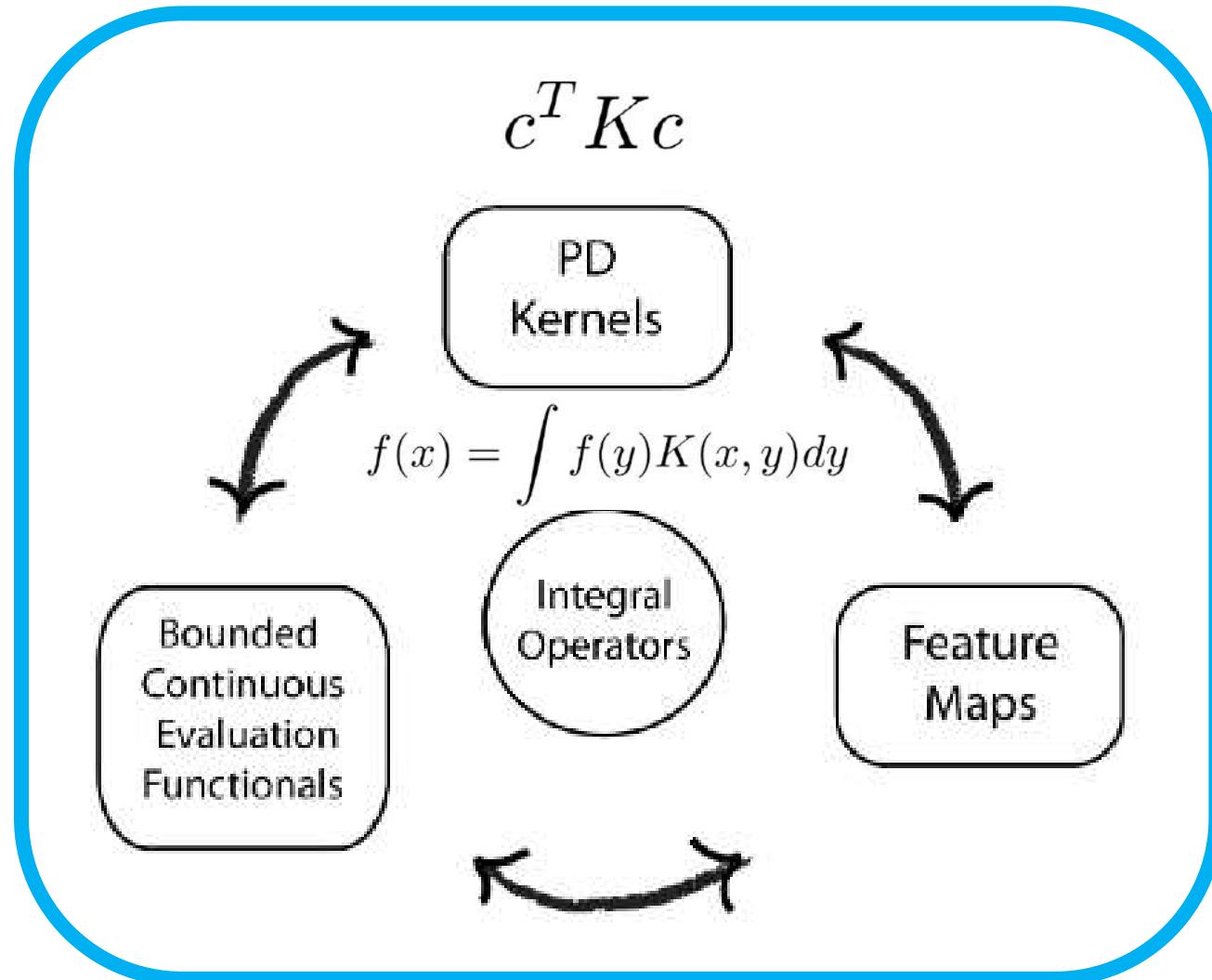
Conversely, let us assume that for any  $\mathbf{x} \in \mathcal{X}$  the linear form  $f \in \mathcal{H} \mapsto f(\mathbf{x})$  is continuous.

Then by Riesz representation theorem (general property of Hilbert spaces) there exists a unique  $g_{\mathbf{x}} \in \mathcal{H}$  such that:

$$f(\mathbf{x}) = \langle f, g_{\mathbf{x}} \rangle_{\mathcal{H}} .$$

The function  $K(\mathbf{x}, \mathbf{y}) = g_{\mathbf{x}}(\mathbf{y})$  is then a r.k. for  $\mathcal{H}$ . □

The picture is:



# Summarizing

- Feature maps define kernels and PD operators but the relation is not one to one
- Feature maps defined by the reproducing kernel Hilbert spaces property are one to one with kernels and pd
- To have the reproducing property we consider operators that satisfy the Mercer theorem
- Pointwise continuity is one to one with RKHS

# Why do we care?

---

(Representer Theorem): *The optimal solution to any vector valued function learning problem of the form.*

$f^* = \operatorname{argmin}_{f \in \mathcal{H}_K} L(f(\mathbf{x}_1) \dots f(\mathbf{x}_l)) + \lambda \|f\|_{\mathcal{H}_K}^2$ , is a sum of matrix-vector products of the form

$f^*(\mathbf{x}) = \sum_{i=1}^l K(\mathbf{x}, \mathbf{x}_i) \boldsymbol{\alpha}_i$  where  $\boldsymbol{\alpha}_i \in \mathbb{R}^n$ ,  $i = 1 \dots l$ ,  $L$  is an arbitrary loss function (which can also be an indicator function encoding arbitrary constraints) and  $\lambda > 0$  is a regularization parameter.

- **We reduced an infinite dimensional non-linear optimization to a linear finite one!!!**
- **Choosing a kernel we choose a regularization!!!**

## Motivation

- An RKHS is a space of (potentially nonlinear) functions, and  $\|f\|_{\mathcal{H}}$  measures the smoothness of  $f$ .
- Given a set of data  $(\mathbf{x}_i \in \mathcal{X}, y_i \in \mathbb{R})_{i=1,\dots,n}$ , a natural way to estimate a regression function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is to solve something like:

$$\min_{f \in \mathcal{H}} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i))}_{\text{empirical risk, data fit}} + \underbrace{\lambda \|f\|_{\mathcal{H}}^2}_{\text{regularization}}. \quad (1)$$

for a loss function  $\ell$  such as  $\ell(y, t) = (y - t)^2$ .

- How to solve in practice this problem, potentially in infinite dimension?

# Proof

**f in the subspace spanned by the kernel → The loss is minimized**

Suppose we project  $f$  onto the subspace:

$$\text{span}\{k(x_i, \cdot) : 1 \leq i \leq n\}$$

obtaining  $f_s$  (the component along the subspace) and  $f_\perp$  (the component perpendicular to the subspace).

We have:

$$f = f_s + f_\perp \Rightarrow \|f\|^2 = \|f_s\|^2 + \|f_\perp\|^2 \geq \|f_s\|^2$$

Since  $\Omega$  is non-decreasing,

$$\Omega(\|f\|_H^2) \geq \Omega(\|f_s\|_H^2)$$

implying that  $\Omega(\cdots)$  is minimized if  $f$  lies in the subspace. Furthermore, since the kernel  $k$  has the reproducing property, we have:

$$f(x_i) = \langle f, k(x_i, \cdot) \rangle = \langle f_s, k(x_i, \cdot) \rangle + \langle f_\perp, k(x_i, \cdot) \rangle = \langle f_s, k(x_i, \cdot) \rangle = f_s(x_i)$$

Implying that:

$$L(f(x_1), \dots, f(x_n)) = L(f_s(x_1), \dots, f_s(x_n))$$

Hence,  $L(\cdots)$  depends only on the component of  $f$  lying in the subspace:  $\text{span}\{k(x_i, \cdot): 1 \leq i \leq n\}$ , and  $\Omega(\cdots)$  is minimized if  $f$  lies in that subspace. Hence,  $J(f)$  is minimized if  $f$  lies in that subspace, and we can express the minimizer as:

$$f^*(\cdot) = \sum_{i=1}^n \alpha_i k(x_i, \cdot)$$

## Practical use of the representer theorem (1/2)

- When the representer theorem holds, we know that we can look for a solution of the form

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x}), \quad \text{for some } \boldsymbol{\alpha} \in \mathbb{R}^n.$$

- For any  $j = 1, \dots, n$ , we have

$$f(\mathbf{x}_j) = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x}_j) = [\mathbf{K}\boldsymbol{\alpha}]_j.$$

- Furthermore,

$$\|f\|_{\mathcal{H}}^2 = \left\| \sum_{i=1}^n \alpha_i K_{\mathbf{x}_i} \right\|_{\mathcal{H}}^2 = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha}.$$

## Practical use of the representer theorem (2/2)

- Therefore, a problem of the form

$$\min_{f \in \mathcal{H}} \Psi(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n), \|f\|_{\mathcal{H}}^2)$$

is equivalent to the following  $n$ -dimensional optimization problem:

$$\min_{\alpha \in \mathbb{R}^n} \Psi([\mathbf{K}\alpha]_1, \dots, [\mathbf{K}\alpha]_n, \alpha^\top \mathbf{K} \alpha).$$

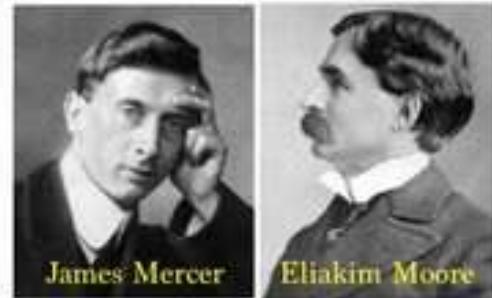
# Summary

Feature map:  $\varphi : x \in \mathbb{R}^d \mapsto \mathcal{H}$  Hilbert space.

Reproducing Hilbert space:  $f : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\|f\|_R = \min \{\|h\|_{\mathcal{H}} ; h \in \mathcal{H}, f = \langle h, \varphi(\cdot) \rangle_{\mathcal{H}}\}$$

Reproducing kernel:  $k(x, y) \stackrel{\text{def.}}{=} \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}$



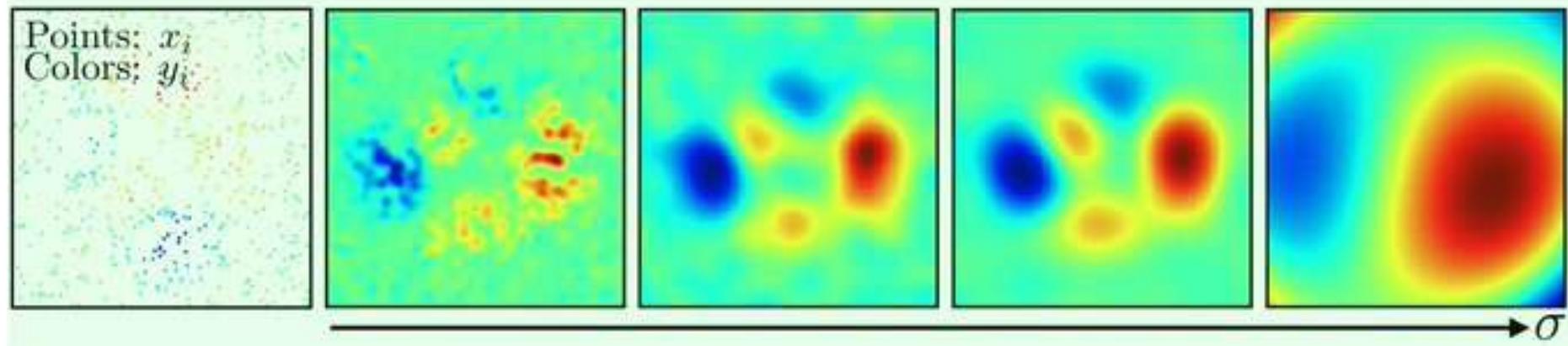
*Representer theorem:* the solution of  $\min_{f: \mathbb{R}^d \rightarrow \mathbb{R}} \mathcal{E}((f(x_i))_i) + \|f\|_R^2$

$$\text{satisfies } f(x) = \sum_i w_i K(x_i, x).$$

*Example:* regression

$$\mathcal{E}((f_i)_i) = \frac{1}{n} \sum_{i=1}^n |f_i - y_i|^2$$
$$k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$$

Points:  $x_i$   
Colors:  $y_i$



# Kernel Ridge Regression (KRR)

- Let us now consider a RKHS  $\mathcal{H}$ , associated to a p.d. kernel  $K$  on  $\mathcal{X}$ .
- KRR is obtained by **regularizing** the MSE criterion by the RKHS norm:

$$\hat{f} = \arg \min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{H}}^2 \quad (2)$$

- 1st effect = prevent overfitting by penalizing non-smooth functions.*
- By the representer theorem, any solution of (2) can be expanded as

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x}).$$

- 2nd effect = simplifying the solution.*

# Solving KRR

- Let  $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$
- Let  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)^\top \in \mathbb{R}^n$
- Let  $\mathbf{K}$  be the  $n \times n$  Gram matrix:  $\mathbf{K}_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
- We can then write:

$$(\hat{f}(\mathbf{x}_1), \dots, \hat{f}(\mathbf{x}_n))^\top = \mathbf{K}\boldsymbol{\alpha}$$

- The following holds as usual:

$$\|\hat{f}\|_{\mathcal{H}}^2 = \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha}$$

- The KRR problem (2) is therefore equivalent to:

$$\arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \frac{1}{n} (\mathbf{K}\boldsymbol{\alpha} - \mathbf{y})^\top (\mathbf{K}\boldsymbol{\alpha} - \mathbf{y}) + \lambda \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha}$$

# Solving KRR

$$\arg \min_{\alpha \in \mathbb{R}^n} \frac{1}{n} (\mathbf{K}\alpha - \mathbf{y})^\top (\mathbf{K}\alpha - \mathbf{y}) + \lambda \alpha^\top \mathbf{K}\alpha$$

- This is a convex and differentiable function of  $\alpha$ . Its minimum can therefore be found by setting the gradient in  $\alpha$  to zero:

$$\begin{aligned} 0 &= \frac{2}{n} \mathbf{K} (\mathbf{K}\alpha - \mathbf{y}) + 2\lambda \mathbf{K}\alpha \\ &= \mathbf{K} [(\mathbf{K} + \lambda n \mathbf{I}) \alpha - \mathbf{y}] \end{aligned}$$

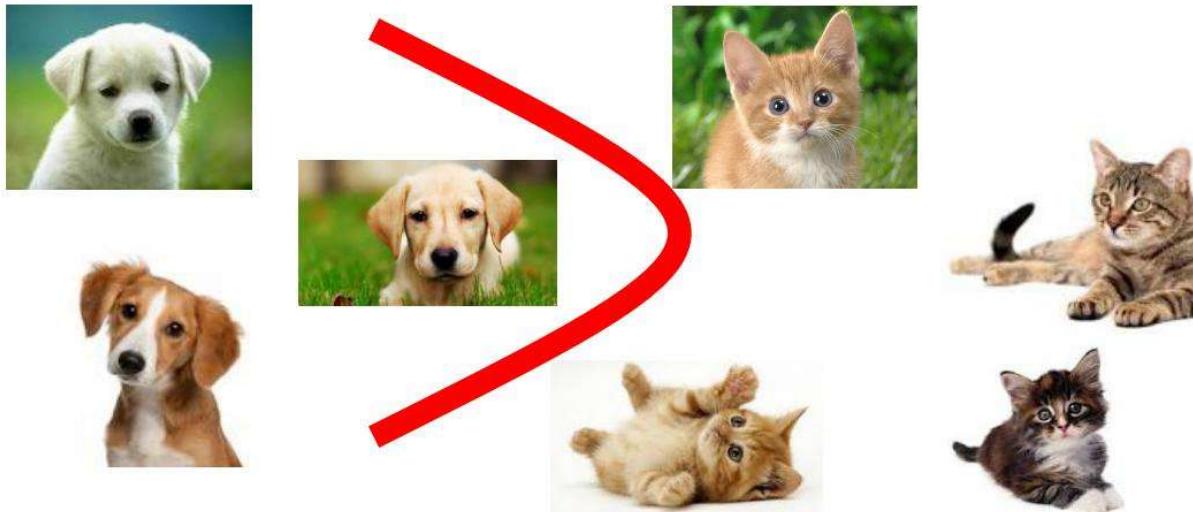
- For  $\lambda > 0$ ,  $\mathbf{K} + \lambda n \mathbf{I}$  is invertible (because  $\mathbf{K}$  is positive semidefinite) so one solution is to take:

$$\alpha = (\mathbf{K} + \lambda n \mathbf{I})^{-1} \mathbf{y}.$$

# Binary classification

## Setup

- $\mathcal{X}$  set of inputs
- $\mathcal{Y} = \{-1, 1\}$  binary outputs
- $\mathcal{S}_n = (\mathbf{x}_i, y_i)_{i=1, \dots, n} \in (\mathcal{X} \times \mathcal{Y})^n$  a training set of  $n$  pairs
- Goal = find a function  $f : \mathcal{X} \rightarrow \mathbb{R}$  to predict  $y$  by *sign*( $f(\mathbf{x})$ )



## The 0/1 loss

- The 0/1 loss measures if a prediction is correct or not:

$$\ell_{0/1}(f(\mathbf{x}), y) = \mathbf{1}(yf(\mathbf{x}) < 0) = \begin{cases} 0 & \text{if } y = \text{sign}(f(\mathbf{x})) \\ 1 & \text{otherwise.} \end{cases}$$

- It is then tempting to learn  $f$  by solving:

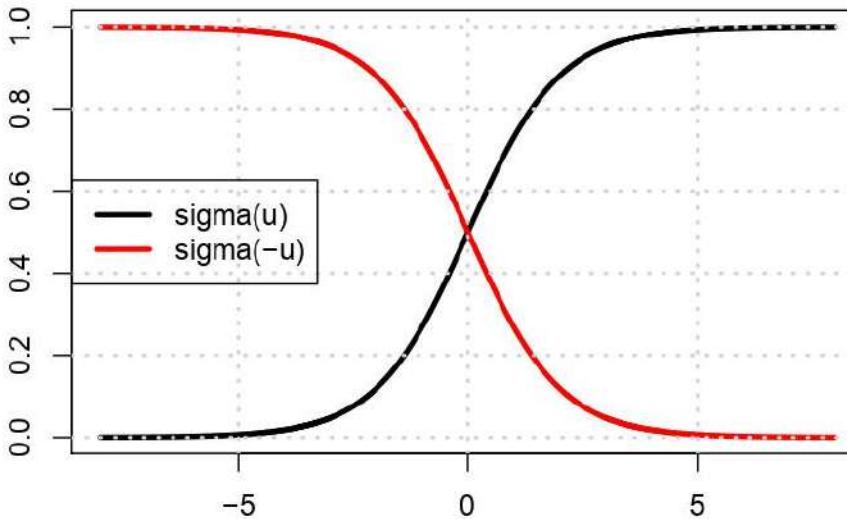
$$\min_{f \in \mathcal{H}} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell_{0/1}(f(\mathbf{x}_i), y_i)}_{\text{misclassification rate}} + \underbrace{\lambda \|f\|_{\mathcal{H}}^2}_{\text{regularization}}$$

- However:
  - The problem is non-smooth, and typically NP-hard to solve

## The logistic loss

- An alternative is to define a probabilistic model of  $y$  parametrized by  $f(\mathbf{x})$ , e.g.:

$$\forall \mathbf{y} \in \{-1, 1\}, \quad p(y | f(\mathbf{x})) = \frac{1}{1 + e^{-yf(\mathbf{x})}} = \sigma(yf(\mathbf{x}))$$



- The **logistic loss** is the negative conditional likelihood:

$$\ell_{logistic}(f(\mathbf{x}), y) = -\ln p(y | f(\mathbf{x})) = \ln(1 + e^{-yf(\mathbf{x})})$$

## Kernel logistic regression (KLR)

$$\begin{aligned}\hat{f} &= \arg \min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell_{logistic}(f(\mathbf{x}_i), y_i) + \frac{\lambda}{2} \| f \|_{\mathcal{H}}^2 \\ &= \arg \min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ln \left( 1 + e^{-y_i f(\mathbf{x}_i)} \right) + \frac{\lambda}{2} \| f \|_{\mathcal{H}}^2\end{aligned}$$

- Can be interpreted as a regularized conditional maximum likelihood estimator
- No explicit solution, but smooth convex optimization problem that can be solved numerically

## Solving KLR

- By the representer theorem, any solution of KLR can be expanded as

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x})$$

and as always we have:

$$(\hat{f}(\mathbf{x}_1), \dots, \hat{f}(\mathbf{x}_n))^T = \mathbf{K}\boldsymbol{\alpha} \quad \text{and} \quad \|\hat{f}\|_{\mathcal{H}}^2 = \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha}$$

- To find  $\boldsymbol{\alpha}$  we therefore need to solve:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n \ln \left( 1 + e^{-y_i [\mathbf{K}\boldsymbol{\alpha}]_i} \right) + \frac{\lambda}{2} \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha}$$

Given:

- $f(x_i) = \sum_{j=1}^n \alpha_j K(x_i, x_j)$ , where  $K(x_i, x_j)$  is the kernel function.
- The probability in logistic regression is  $\sigma(f(x_i)) = \frac{1}{1+e^{-\sum_{j=1}^n \alpha_j K(x_i, x_j)}}$ .

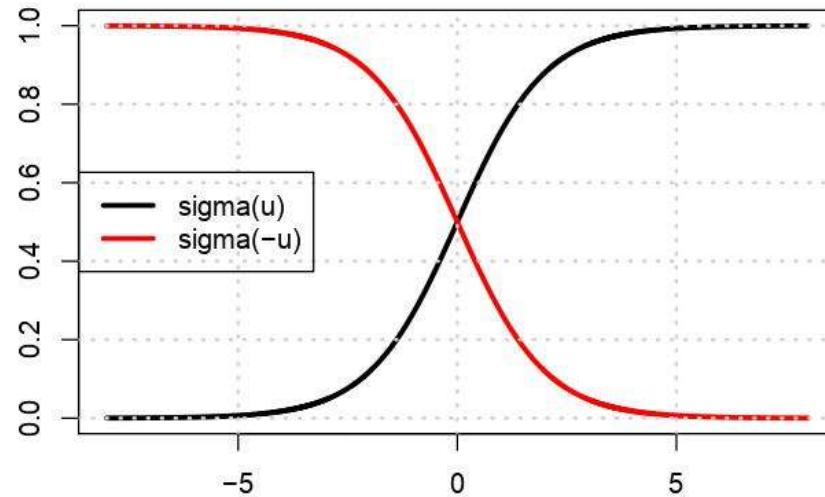
## 1. Logistic Regression Cost Function

The cost function for logistic regression is:

$$J(\alpha) = - \sum_{i=1}^n [y_i \log(\sigma(f(x_i))) + (1 - y_i) \log(1 - \sigma(f(x_i)))] + \frac{\lambda}{2} \sum_{j,k=1}^n \alpha_j \alpha_k K(x_j, x_k).$$

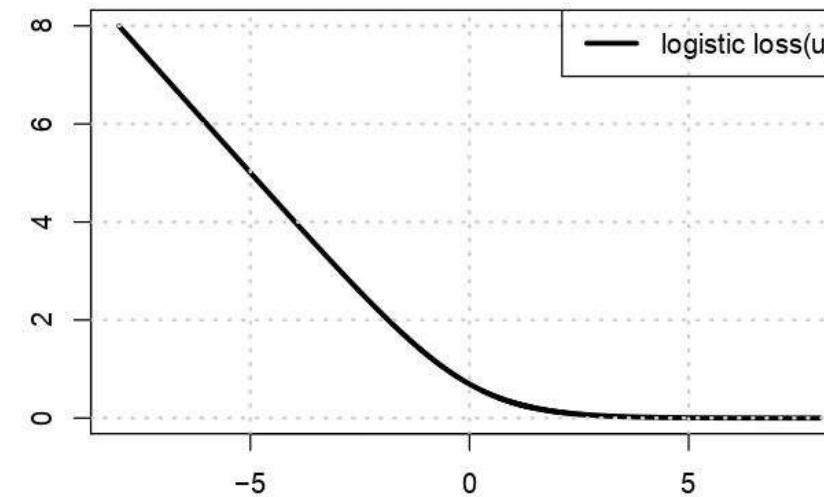
Here,  $\lambda$  is the regularization parameter.

# Technical facts



Sigmoid:

- $\sigma(u) = \frac{1}{1+e^{-u}}$
- $\sigma(-u) = 1 - \sigma(u)$
- $\sigma'(u) = \sigma(u)\sigma(-u) \geq 0$



Logistic loss:

- $\ell_{logistic}(u) = \ln(1 + e^{-u})$
- $\ell'_{logistic}(u) = -\sigma(-u)$
- $\ell''_{logistic}(u) = \sigma(u)\sigma(-u) \geq 0$

## Step 1: Derivative of the Logistic Loss

The logistic loss for a single sample  $x_i$  is:

$$L_i = -[y_i \log(\sigma(f(x_i))) + (1 - y_i) \log(1 - \sigma(f(x_i)))] .$$

Since  $f(x_i) = \sum_{j=1}^n \alpha_j K(x_i, x_j)$ , the derivative of the logistic loss with respect to  $f(x_i)$  is:

$$\frac{\partial L_i}{\partial f(x_i)} = \sigma(f(x_i)) - y_i.$$

Now, we need to express this in terms of  $\alpha_j$ . The derivative of  $f(x_i)$  with respect to  $\alpha_j$  is:

$$\frac{\partial f(x_i)}{\partial \alpha_j} = K(x_i, x_j).$$

Using the chain rule, the derivative of  $L_i$  with respect to  $\alpha_j$  is:

$$\frac{\partial L_i}{\partial \alpha_j} = \frac{\partial L_i}{\partial f(x_i)} \frac{\partial f(x_i)}{\partial \alpha_j} = (\sigma(f(x_i)) - y_i)K(x_i, x_j).$$

## Step 2: Derivative of the Regularization Term

The regularization term is:

$$\frac{\lambda}{2} \sum_{j,k=1}^n \alpha_j \alpha_k K(x_j, x_k).$$

The derivative of this term with respect to  $\alpha_j$  is:

$$\frac{\partial}{\partial \alpha_j} \left( \frac{\lambda}{2} \sum_{k=1}^n \alpha_k K(x_j, x_k) \right) = \lambda \sum_{k=1}^n \alpha_k K(x_j, x_k).$$

### 3. Combining Both Parts

Now, we combine the derivatives from both parts (logistic loss and regularization):

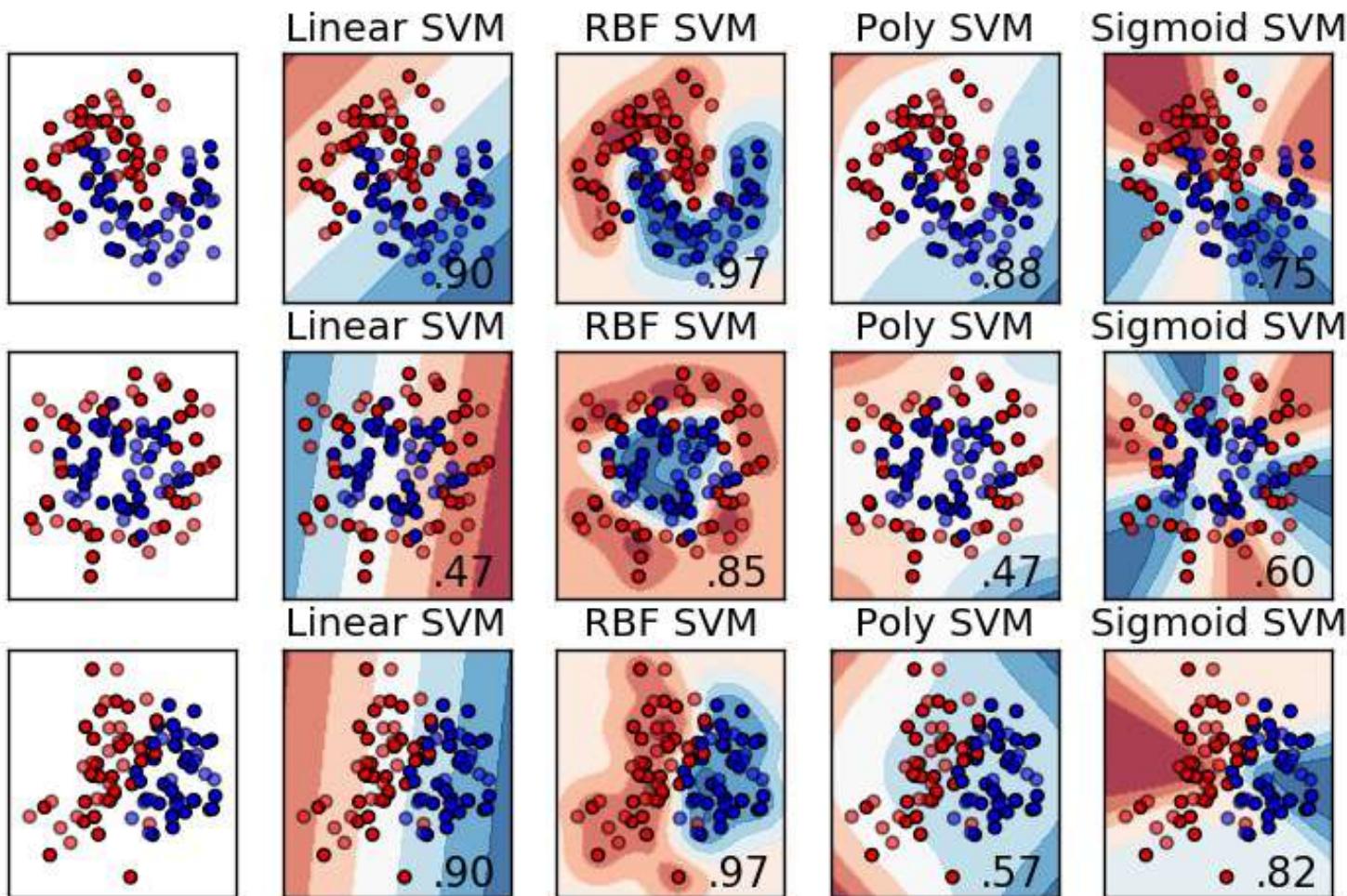
$$\frac{\partial J(\alpha)}{\partial \alpha_j} = \sum_{i=1}^n (\sigma(f(x_i)) - y_i) K(x_i, x_j) + \lambda \sum_{k=1}^n \alpha_k K(x_j, x_k).$$

### Final Expression

The derivative of the logistic regression cost function in the kernelized form with respect to the coefficient  $\alpha_j$  is:

$$\frac{\partial J(\alpha)}{\partial \alpha_j} = \sum_{i=1}^n (\sigma(\sum_{k=1}^n \alpha_k K(x_i, x_k)) - y_i) K(x_i, x_j) + \lambda \sum_{k=1}^n \alpha_k K(x_j, x_k).$$

# Kernel SVM example



# Different kernels induce different regularizations

Since  $f = \sum_{j=1}^n c_j K_{x_j}$ , then

The norm depends on the kernel

Changing the kernel we change the regularization

$$\begin{aligned}\|f\|_{\mathcal{H}}^2 &= \langle f, f \rangle_{\mathcal{H}} \\ &= \left\langle \sum_{i=1}^n c_i K_{x_i}, \sum_{j=1}^n c_j K_{x_j} \right\rangle_{\mathcal{H}} \\ &= \sum_{i=1}^n \sum_{j=1}^n c_i c_j \langle K_{x_i}, K_{x_j} \rangle_{\mathcal{H}} \\ &= \sum_{i=1}^n \sum_{j=1}^n c_i c_j K(x_i, x_j) = \mathbf{c}^T \mathbf{K} \mathbf{c}\end{aligned}$$

$$\|f\|_{\mathcal{H}} = \mathbf{c}^T \mathbf{K} \mathbf{c} = \mathbf{c}^T \mathbf{U}^T \Lambda \mathbf{U} \mathbf{c} = \tilde{\mathbf{c}}^T \Lambda \tilde{\mathbf{c}}$$

## Smoothness functional

### A simple inequality

- By Cauchy-Schwarz we have, for any function  $f \in \mathcal{H}$  and any two points  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ :

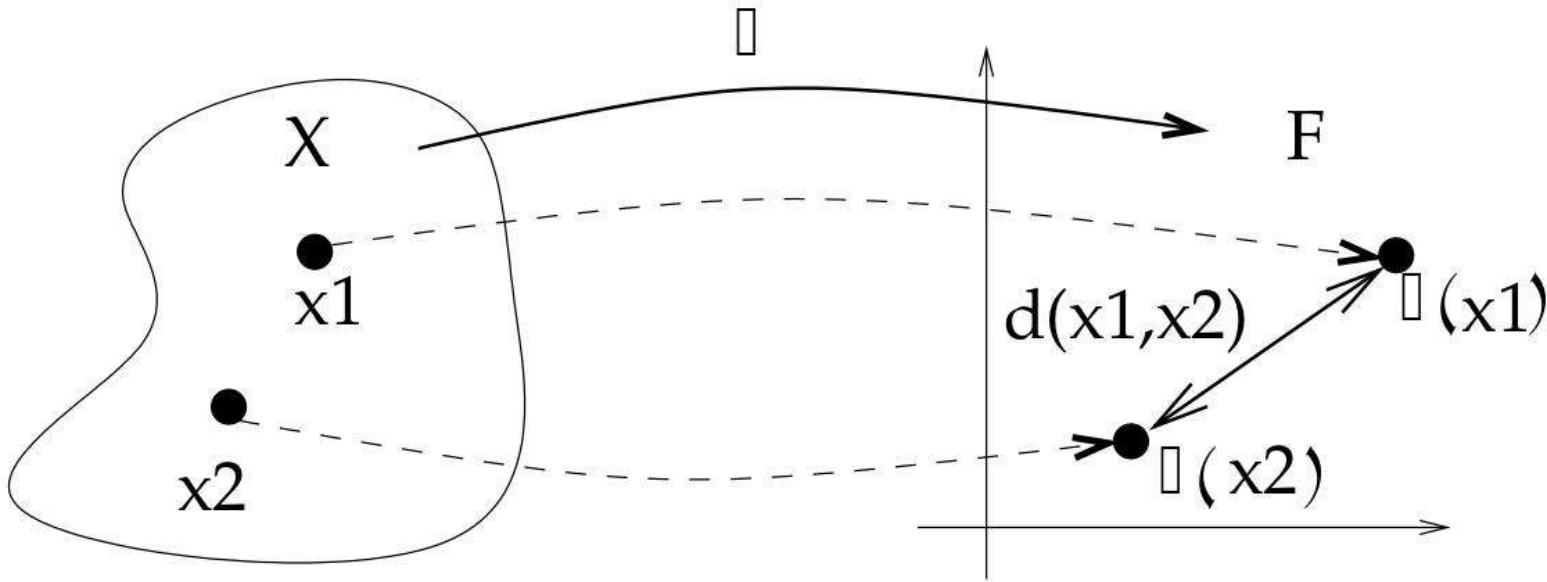
$$\begin{aligned} |f(\mathbf{x}) - f(\mathbf{x}')| &= |\langle f, K_{\mathbf{x}} - K_{\mathbf{x}'} \rangle_{\mathcal{H}}| \\ &\leq \|f\|_{\mathcal{H}} \times \|K_{\mathbf{x}} - K_{\mathbf{x}'}\|_{\mathcal{H}} \\ &= \|f\|_{\mathcal{H}} \times d_K(\mathbf{x}, \mathbf{x}'). \end{aligned}$$

- The norm of a function in the RKHS controls **how fast** the function varies over  $\mathcal{X}$  with respect to the **geometry defined by the kernel** (Lipschitz with constant  $\|f\|_{\mathcal{H}}$ ).

### Important message

Small norm  $\implies$  slow variations.

## Example 1: computing distances in the feature space

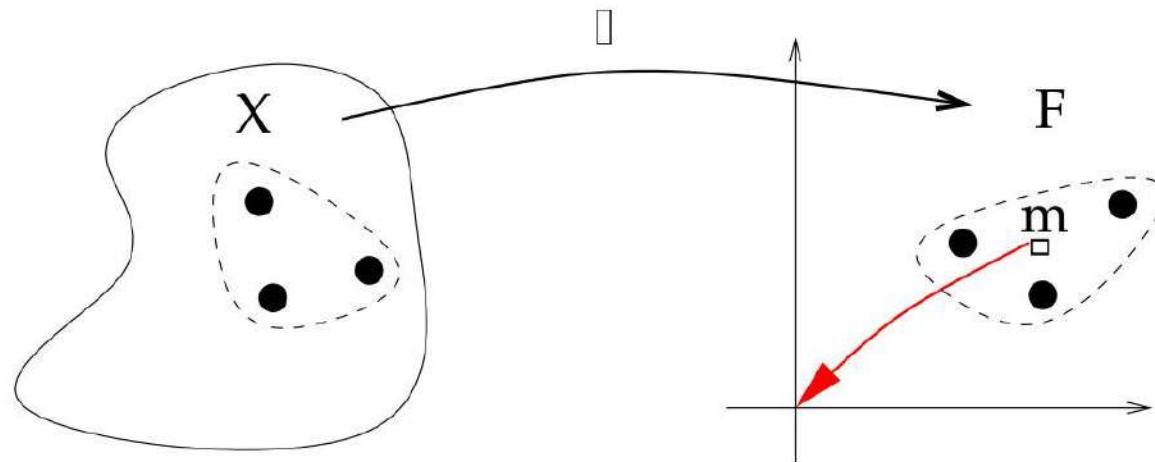


$$\begin{aligned} d_K(x_1, x_2)^2 &= \| \Phi(x_1) - \Phi(x_2) \|_{\mathcal{H}}^2 \\ &= \langle \Phi(x_1) - \Phi(x_2), \Phi(x_1) - \Phi(x_2) \rangle_{\mathcal{H}} \\ &= \langle \Phi(x_1), \Phi(x_1) \rangle_{\mathcal{H}} + \langle \Phi(x_2), \Phi(x_2) \rangle_{\mathcal{H}} - 2 \langle \Phi(x_1), \Phi(x_2) \rangle_{\mathcal{H}} \\ d_K(x_1, x_2)^2 &= K(x_1, x_1) + K(x_2, x_2) - 2K(x_1, x_2) \end{aligned}$$

## Example 3: Centering data in the feature space

### Problem

- Let  $\mathcal{S} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  be a finite set of points in  $\mathcal{X}$  endowed with a p.d. kernel  $K$ . Let  $\mathbf{K}$  be their  $n \times n$  Gram matrix:  $[\mathbf{K}]_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ .
- Let  $\boldsymbol{\mu} = 1/n \sum_{i=1}^n \Phi(\mathbf{x}_i)$  their barycenter, and  $\mathbf{u}_i = \Phi(\mathbf{x}_i) - \boldsymbol{\mu}$  for  $i = 1, \dots, n$  be centered data in  $\mathcal{H}$ .
- How to compute the centered Gram matrix  $[\mathbf{K}^c]_{i,j} = \langle \mathbf{u}_i, \mathbf{u}_j \rangle_{\mathcal{H}}$ ?



## Computation

- A direct computation gives, for  $0 \leq i, j \leq n$ :

$$\begin{aligned}\mathbf{K}_{i,j}^c &= \langle \Phi(\mathbf{x}_i) - \boldsymbol{\mu}, \Phi(\mathbf{x}_j) - \boldsymbol{\mu} \rangle_{\mathcal{H}} \\ &= \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle_{\mathcal{H}} - \langle \boldsymbol{\mu}, \Phi(\mathbf{x}_i) + \Phi(\mathbf{x}_j) \rangle_{\mathcal{H}} + \langle \boldsymbol{\mu}, \boldsymbol{\mu} \rangle_{\mathcal{H}} \\ &= \mathbf{K}_{i,j} - \frac{1}{n} \sum_{k=1}^n (\mathbf{K}_{i,k} + \mathbf{K}_{j,k}) + \frac{1}{n^2} \sum_{k,l=1}^n \mathbf{K}_{k,l}.\end{aligned}$$

- This can be rewritten in matricial form:

$$\mathbf{K}^c = \mathbf{K} - \mathbf{U}\mathbf{K} - \mathbf{K}\mathbf{U} + \mathbf{U}\mathbf{K}\mathbf{U} = (\mathbf{I} - \mathbf{U})\mathbf{K}(\mathbf{I} - \mathbf{U}),$$

where  $\mathbf{U}_{i,j} = 1/n$  for  $1 \leq i, j \leq n$ .

# Kernel PCA

We want to extend the PCA in the features space.

- Calculate the covariance matrix in the feats space:

$$C_F = \frac{1}{N} \sum_i^N \phi(x_i) \phi^T(x_i)$$

- We want to solve  $C_F v = \lambda v$

Eigenvectors can be expressed as a linear combination of the features i.e.

$$v = \sum_{i=1}^N \alpha_i \phi(x_i)$$

- We have

$$C_F v = \lambda v = \frac{1}{N} \sum_i^N \phi(x_i) \phi^T(x_i) v$$

- Thus

$$v = \frac{1}{N\lambda} \sum_i^N \phi(x_i) \phi^T(x_i) v = \sum_i^N \phi(x_i) \frac{\beta_i}{N\lambda} = \sum_{i=1}^N \alpha_i \phi(x_i)$$

Let's multiply both sides by  $\phi^T(x_i)$

- Expanding

$$\lambda\phi^T(x_k)v = \phi^T(x_k)C_Fv$$

$$\lambda \sum_i \alpha_i \phi^T(x_k) \phi(x_i) = \frac{1}{N} \sum_i \alpha_i (\phi^T(x_k) \sum_j \phi(x_j) (\phi^T(x_j) \phi(x_i)))$$

- Define  $K_{ij} = \phi^T(x_j) \phi(x_i)$  we have

$$N\lambda K\alpha = K^2\alpha \rightarrow N\lambda\alpha = K\alpha$$

# Normalization

$$\hat{\phi}(x_i) = \phi(x_i) - \frac{1}{N} \sum_{j=1}^N \phi(x_k)$$

$$\begin{aligned}\hat{K}_{ij} &= \left( \phi(x_i) - \frac{1}{N} \sum_{k=1}^N \phi(x_k) \right) \left( \phi(x_j) - \frac{1}{N} \sum_{l=1}^N \phi(x_l) \right) \\ &= K_{ij} - \frac{1}{N} \sum_{k=1}^N \phi^T(x_i) \phi(x_k) - \frac{1}{N} \sum_{l=1}^N \phi^T(x_j) \phi(x_l) + \frac{1}{N^2} \sum_{kl=1}^N \phi^T(x_k) \phi(x_l) \\ &= K_{ij} - \frac{1}{N} \sum_{k=1}^N K_{i,k} - \frac{1}{N} \sum_{k=1}^N K_{j,k} + \frac{1}{N^2} \sum_{k,l=1}^N K_{k,l}\end{aligned}$$

# Different kernels different regularizations

Since  $f = \sum_{j=1}^n c_j K_{x_j}$ , then

$$\begin{aligned}\|f\|_{\mathcal{H}}^2 &= \langle f, f \rangle_{\mathcal{H}} \\ &= \left\langle \sum_{i=1}^n c_i K_{x_i}, \sum_{j=1}^n c_j K_{x_j} \right\rangle_{\mathcal{H}} \\ &= \sum_{i=1}^n \sum_{j=1}^n c_i c_j \langle K_{x_i}, K_{x_j} \rangle_{\mathcal{H}} \\ &= \sum_{i=1}^n \sum_{j=1}^n c_i c_j K(x_i, x_j) = \mathbf{c}^T \mathbf{K} \mathbf{c}\end{aligned}$$

The norm depends on the kernel

Changing the kernel we change the regularization

$$\|f\|_{\mathcal{H}} = \mathbf{c}^T \mathbf{K} \mathbf{c} = \mathbf{c}^T \mathbf{U}^T \Lambda \mathbf{U} \mathbf{c} = \tilde{\mathbf{c}}^T \Lambda \tilde{\mathbf{c}}$$

# Regularization due to kernel choice

Gaussian kernel

$$K(x, y) = e^{-\frac{(x-y)^2}{2\sigma^2}}$$

corresponds to:

$$\varphi(t) = e^{-\frac{t^2}{2\sigma^2}}$$

$$\hat{\varphi}(\omega) = e^{-\frac{\sigma^2 \omega^2}{2}}$$

and

$$\mathcal{H} = \left\{ f : \int \left| \hat{f}(\omega) \right|^2 e^{\frac{\sigma^2 \omega^2}{2}} d\omega < \infty \right\}.$$

In particular, all functions in  $\mathcal{H}$  are **infinitely differentiable** with all derivatives in  $L^2$ .

## Laplace kernel

$$K(x, y) = \frac{1}{2} e^{-\gamma|x-y|}$$

corresponds to:

$$\varphi(t) = \frac{1}{2} e^{-\gamma|t|}$$

$$\hat{\varphi}(\omega) = \frac{\gamma}{\gamma^2 + \omega^2}$$

and

$$\mathcal{H} = \left\{ f : \int \left| \hat{f}(\omega) \right|^2 \frac{(\gamma^2 + \omega^2)}{\gamma} d\omega < \infty \right\},$$

the set of functions  $L^2$  differentiable with derivatives in  $L^2$  (Sobolev norm).

## Compare with Gaussian Kernel

## Low-frequency filter

$$K(x, y) = \frac{\sin(\Omega(x - y))}{\pi(x - y)}$$

corresponds to:

$$\varphi(t) = \frac{\sin(\Omega t)}{\pi t}$$

$$\hat{\varphi}(\omega) = 1_{[-\Omega, \Omega]}(\omega)$$

and

$$\mathcal{H} = \left\{ f : \int_{|\omega| > \Omega} |\hat{f}(\omega)|^2 d\omega = 0 \right\},$$

the set of functions whose spectrum is included in  $[-\Omega, \Omega]$ .

Compare with Gaussian Kernel

# Hypothesis space

- Linear separable

$$f(x) = w^T x$$

We learn the separator on data

- Non-linear separable

$$f(x) = w^T \phi(x)$$

We learn the separator on a  
priori hand-crafted features

# Models

- ▶ Linear

$$f(x) = w^\top x.$$

- ▶ Features

$$f(x) = w^\top \Phi(x) = \sum_{j=1}^p w^j \varphi_j(x).$$

- ▶ Kernels

$$f(x) = \sum_{i=1}^n K(x, x_i) c_i.$$

- Can we learn the features together with the separator?
- Which features? We choose this parametrization:

$$f(x) = w^T \sigma(W^T x)$$

$$w \in \mathbb{R}^k, \quad W \in \mathbb{R}^{d \times k}, \quad \sigma : \mathbb{R} \rightarrow \mathbb{R}$$

# (1):Linear filtering

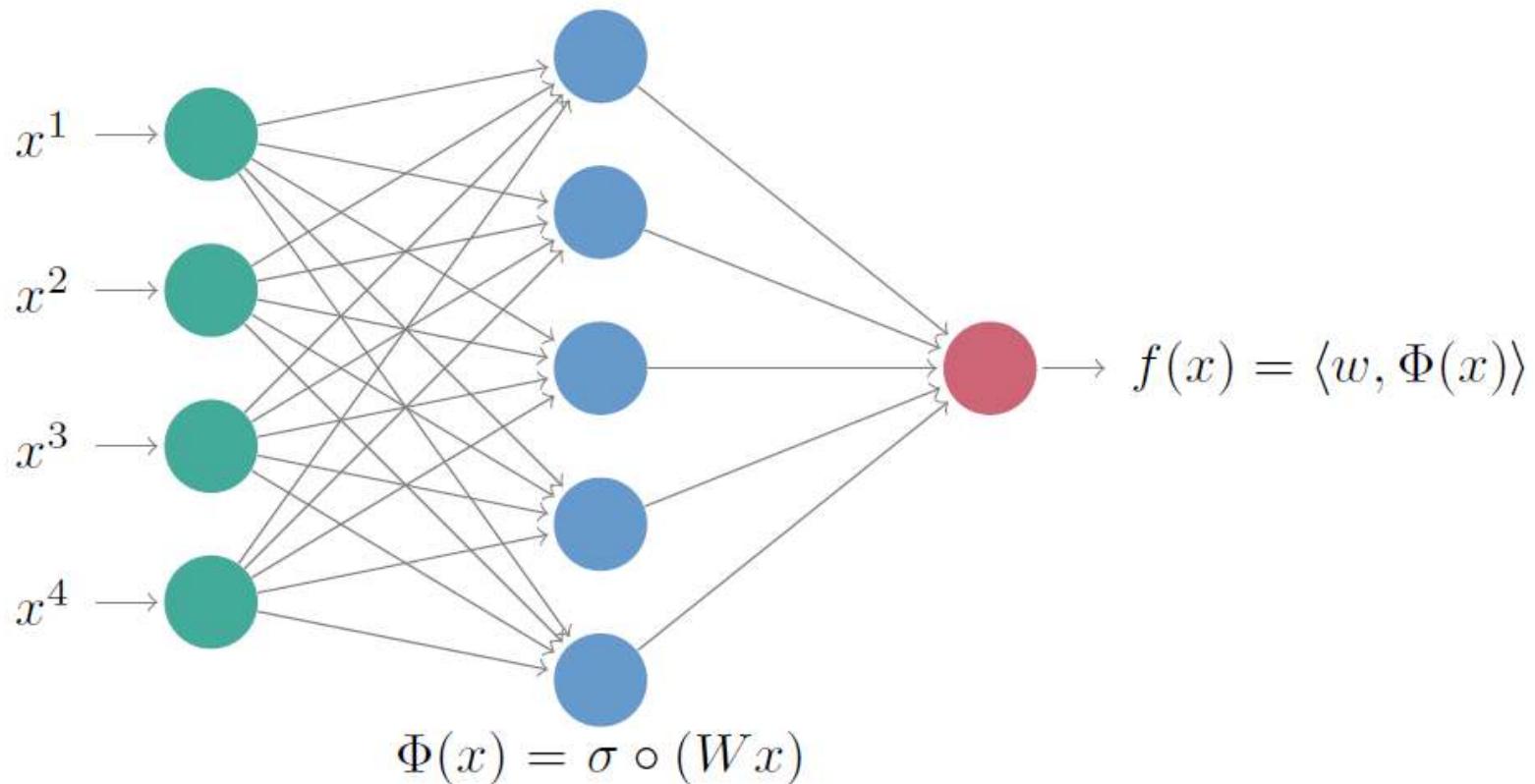
$$y = W^T x + b = \begin{bmatrix} w_1^T x + b_1 \\ w_2^T x + b_2 \\ \vdots \\ w_M^T x + b_M \end{bmatrix}$$

## (2):Non-Linear pointwise filtering

$$y = \sigma(W^T x + b) = \begin{bmatrix} \sigma(w_1^T x + b_1) \\ \sigma(w_2^T x + b_2) \\ \vdots \\ \sigma(w_M^T x + b_M) \end{bmatrix}$$

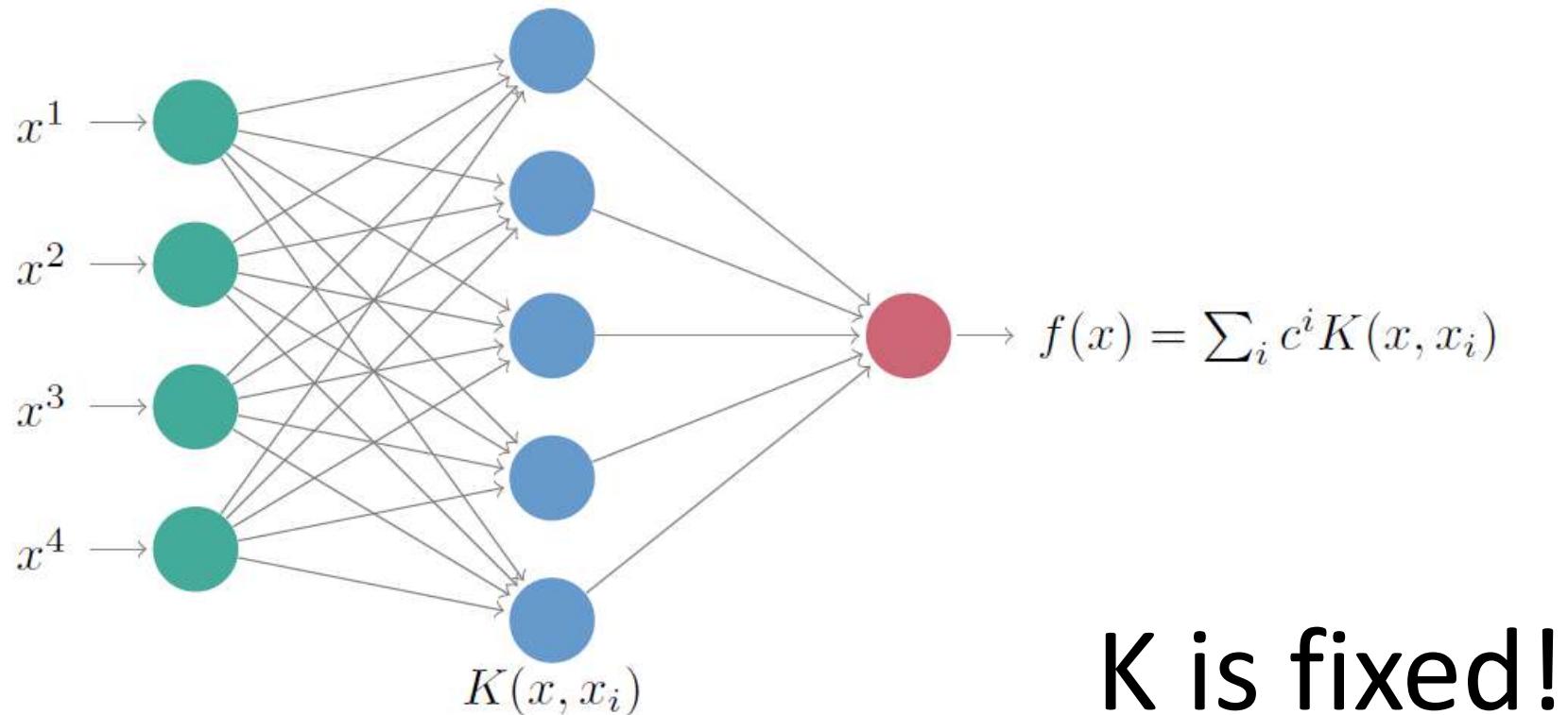
## Shallow Neural Network Representations

$l_0$ , input       $l_1$ , hidden layer     $l_L$ , output layer



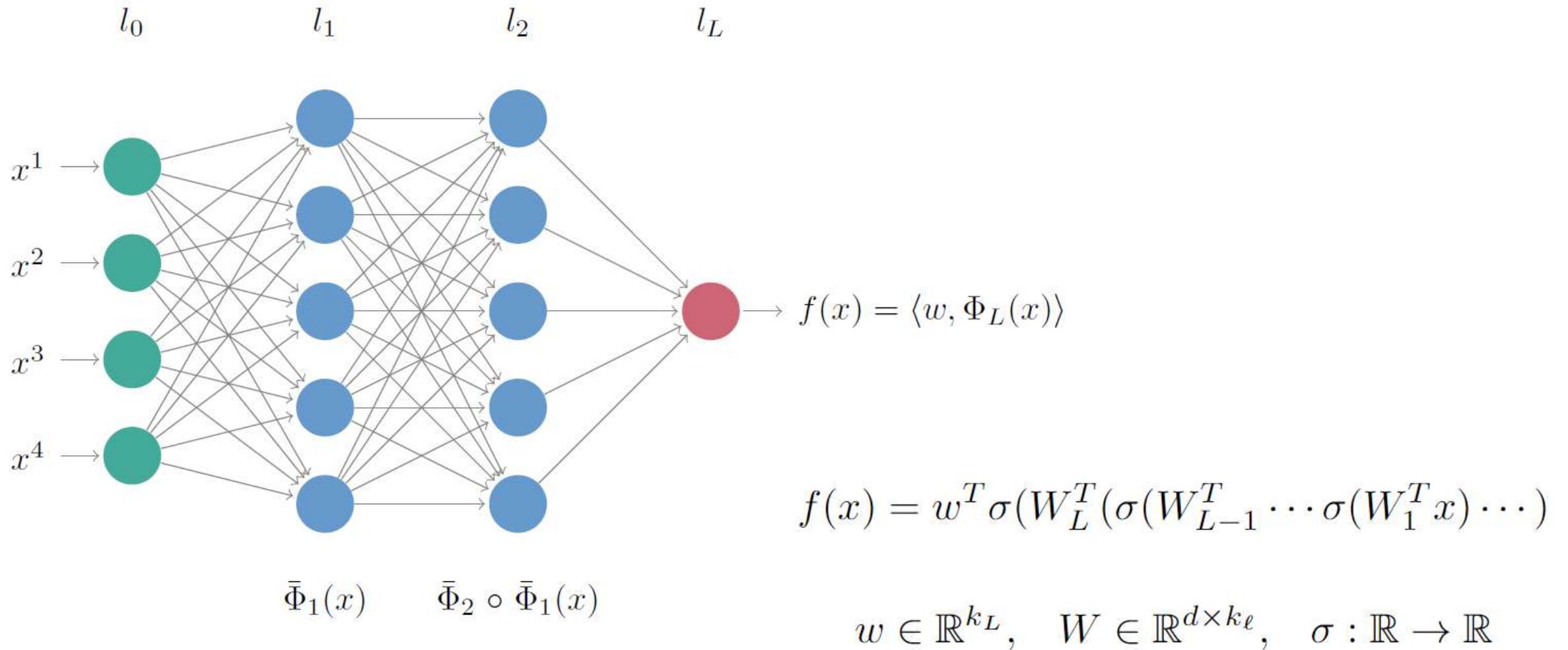
## Kernel Machine Representations (e.g. RBF)

$l_0$ , input       $l_1$ , hidden layer     $l_L$ , output layer



**K is fixed!**

## Deep Neural Network Representations



# Link with Kernels: Neural Tangent Kernel

---

---

## **Neural Tangent Kernel: Convergence and Generalization in Neural Networks**

---

**Arthur Jacot**

École Polytechnique Fédérale de Lausanne  
[arthur.jacot@netopera.net](mailto:arthur.jacot@netopera.net)

**Franck Gabriel**

Imperial College London and École Polytechnique Fédérale de Lausanne  
[franckrgabriel@gmail.com](mailto:franckrgabriel@gmail.com)

**Clément Hongler**

École Polytechnique Fédérale de Lausanne  
[clement.hongler@gmail.com](mailto:clement.hongler@gmail.com)

Slides by Nolan Dey

<https://arxiv.org/abs/1806.07572>

# Setup

- We have some neural network:  $f(x, w)$ 
  - Input data:  $x \in \mathbb{R}^{n \times d}$
  - Network parameters:  $w \in \mathbb{R}^p$
- Loss:  $L(f(x, w), y) = \frac{1}{2}(f(x, w) - y)^2$
- Optimize using full-batch gradient descent

# What is the Neural Tangent Kernel (NTK)?

- Infinitely wide neural networks can fit any function 
- Infinite width networks trained to convergence can be described by the NTK
- NTK describes training dynamics: 
$$\frac{df(x, w)}{dt} = -NTK(w_0)(f(x, w) - y)$$

Universal  
approximation  
theorem

Is the kernel associated to the network features

# Universal approximation theorem

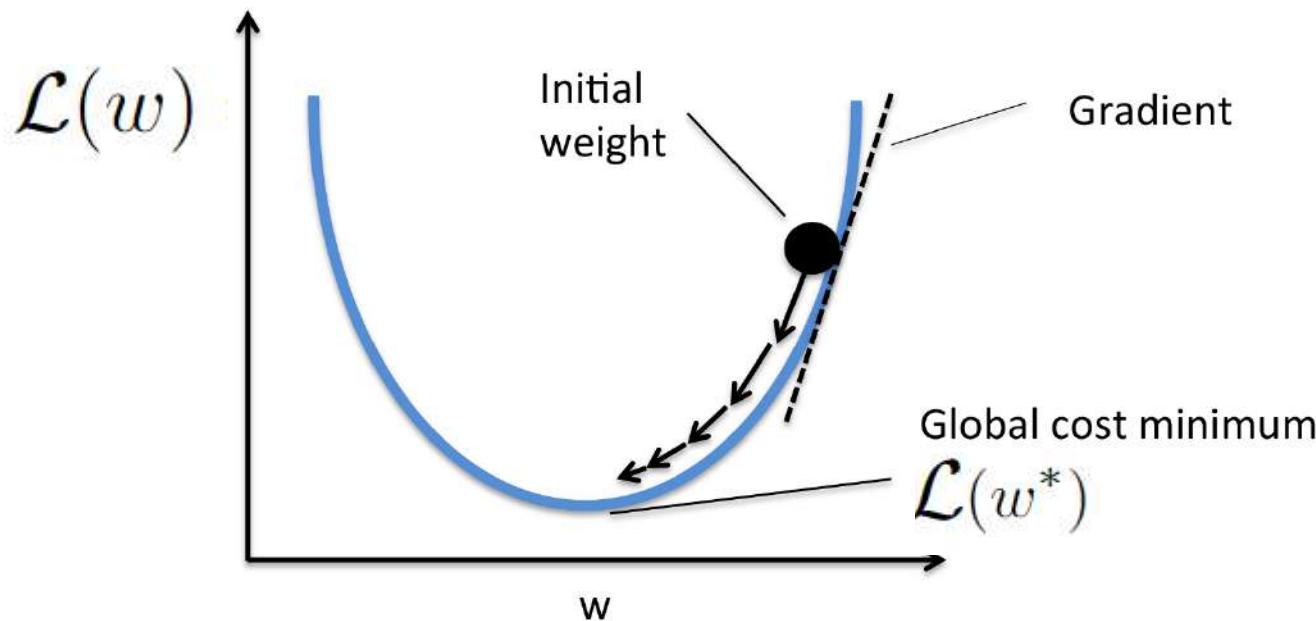
Any **continuous function** defined in the  
**n-dimensional unit hypercube** may be  
approximated by a **finite sum** of the type:

$$\sum_{j=1}^N v_j \varphi \left( \vec{\omega}^{(j)} \cdot \vec{x} + b_j \right),$$

wherein  $v_j, b_j \in \mathbb{R}$ ,  $\vec{\omega}^{(j)} \in \mathbb{R}^n$ , and  $\varphi$  is a  
**continuous discriminatory function**.

How big is N?

# Finding minima with gradient descent



$$w_{t+1} = w_t - \gamma \frac{d\mathcal{L}(w)}{dw}$$

# Taylor expansion of neural network

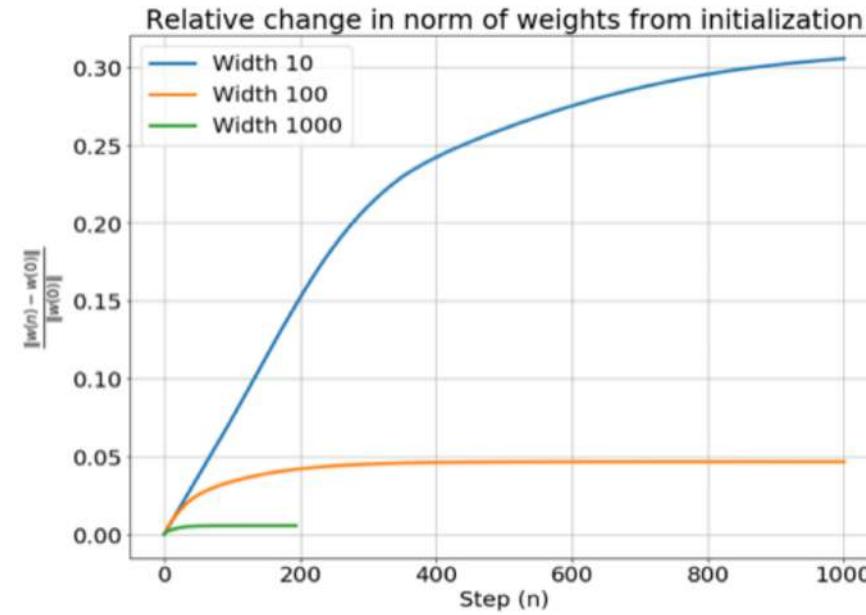
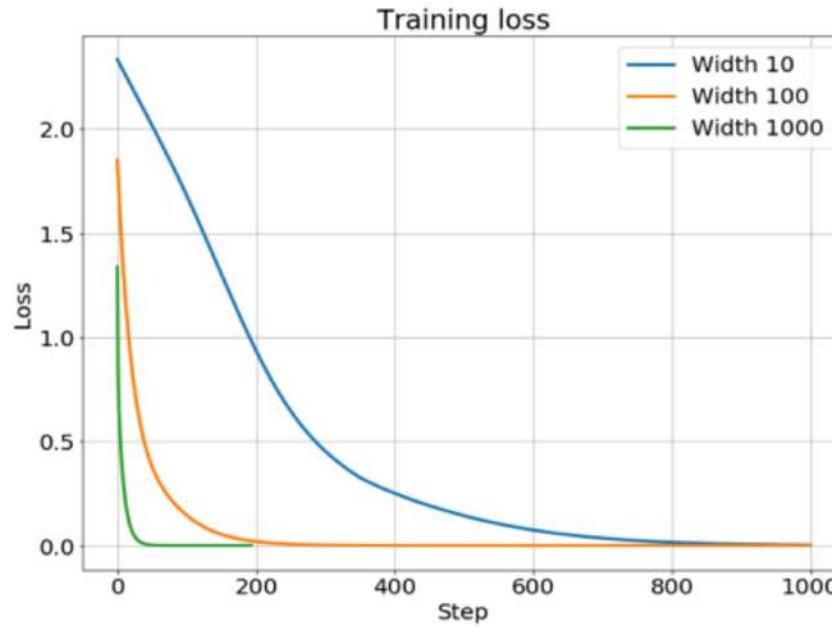
- Taylor expansion of neural network with respect to initialization weights  $w_0$ 
  - $f(x, w) \approx f(x, w_0) + \nabla_w f(x, w_0)^T(w - w_0)$
- This Taylor expansion is only accurate when weights remain close to initialization



**Small learning!**

# When is Taylor expansion accurate?

- Taylor expansion is only accurate when the weights don't change much during training
  - Weights don't change much when network is sufficiently wide
  - Lazy training = weights don't need to change



This is true when the width of the network is very large. Why?

# Gradient Flow

- Gradient descent:  $w_{t+1} = w_t - \eta \nabla_w L(w_t)$

- Rewrite as:  $\frac{w_{t+1} - w_t}{\eta} = - \nabla_w L(w_t)$

- Resembles finite difference ^

- Take infinitesimally small learning rate  $\eta$

- Gradient flow!  $\rightarrow \frac{dw(t)}{dt} = - \nabla_w L(w(t))$

$$\begin{aligned}
\frac{dw(t)}{dt} &= -\nabla_w \mathcal{L}(w(t), x) \\
&= \nabla_w \frac{1}{N} \sum_i (f(w(t), x_i) - y_i)^2 \\
&= -\frac{2}{N} \sum_i (f(w(t), x_i) - y_i) \nabla_w f(w(t), x_i)
\end{aligned}$$



This is not static!!!!

$$\begin{aligned}
\frac{df(w(t), x)}{dt} &= \nabla_w f(w(t), x) \frac{dw(t)}{dt} \\
&= \nabla_w^T f(w(t), x) \frac{-2}{N} \sum_i (f(w(t), x_i) - y_i) \nabla_w f(w(t), x_i) \\
&= -\frac{2}{N} \sum_i (f(w(t), x_i) - y_i) \nabla_w^T f(w(t), x) \nabla_w f(w(t), x_i) \\
&= -\frac{2}{N} \sum_i (f(w(t), x_i) - y_i) K(x, x_i)
\end{aligned}$$

↑  
Supposing quasi-static

↑ Feature  $\Phi(x_i)$

If the width is infinite the change in w is very small and K is the kernel associated to the dynamic of the output

Calculate the kernel for

$$f(w, x) = \sum_i v_i \sigma(w_i^T x)$$

# What does this mean?

- NTK is a useful tool for studying the dynamics of infinitely wide neural networks
- Successful nets in practice DO NOT OPERATE IN THE NTK REGIME
  - Finite nets still outperform their exactly computed infinite width counterparts
  - SGD vs full-batch gradient descent

This analysis is useful to show:

- Kernel models are not good enough to capture NNs
- Implicit Bias/Regularization (next)

# Implicit regularization: NTK solution

$$\begin{aligned}\frac{df(\theta)}{dt} &= -\eta \nabla_{\theta} f(\theta)^{\top} \nabla_{\theta} f(\theta) \nabla_f \mathcal{L} \\ &= -\eta \nabla_{\theta} f(\theta)^{\top} \nabla_{\theta} f(\theta) \nabla_f \mathcal{L} \\ &= -\eta K(\theta) \nabla_f \mathcal{L} \\ &= \textcolor{red}{-\eta K_{\infty} \nabla_f \mathcal{L}}\end{aligned}$$

# NTK solution

$$\begin{aligned}\frac{df(\theta)}{dt} &= -\eta K_\infty(f(\theta) - \mathcal{Y}) \\ \frac{dg(\theta)}{dt} &= -\eta K_\infty g(\theta) \quad ; \text{ let } g(\theta) = f(\theta) - \mathcal{Y} \\ \int \frac{dg(\theta)}{g(\theta)} &= -\eta \int K_\infty dt \\ g(\theta) &= Ce^{-\eta K_\infty t}\end{aligned}$$

When  $t = 0$ , we have  $C = f(\theta(0)) - \mathcal{Y}$  and therefore,

$$f(\theta) = (f(\theta(0)) - \mathcal{Y})e^{-\eta K_\infty t} + \mathcal{Y} = f(\theta(0))e^{-K_\infty t} + (I - e^{-\eta K_\infty t})\mathcal{Y}$$

The expression dictates how the error changes in time!

# Implicit regularization

Solving for  $f$  and using the fact that in the large-width limit (1-layer shallow network with many hidden units) the kernel is approximately stationary:

$$f(X, W(t)) \approx (I - e^{-Kt}) \cdot Y.$$

Now, since the kernel  $K$  is positive semi-definite, we can take its spectral decomposition  $K = Q\Lambda Q^T$  where  $Q$  is an orthogonal matrix whose  $i$ -th column is the eigenvector  $q_i$  of  $K$  and  $\Lambda$  is a diagonal matrix whose diagonal entries  $\lambda_i$  are the corresponding eigenvalues. Since  $e^{-Kt} = Qe^{-\Lambda t}Q^T$  we have

$$Q^T (f(X, W(t)) - Y) = -e^{-\Lambda t} Q^T Y,$$

i.e.,

$$\begin{bmatrix} q_1^T \\ q_2^T \\ \vdots \\ q_N^T \end{bmatrix} (f(X, W(t)) - Y) = \begin{bmatrix} e^{-\lambda_1 t} & & & \\ & e^{-\lambda_2 t} & \cdots & \\ & & \ddots & \\ & & & e^{-\lambda_N t} \end{bmatrix} \begin{bmatrix} q_1^T \\ q_2^T \\ \vdots \\ q_N^T \end{bmatrix} Y.$$

# Implicit regularization

The above equation shows that the convergence rate of  $q_i^T(f(X, W(t)) - Y)$ , *the error in the X predictions*, is determined by the  $i$ -th eigenvalue  $\lambda_i$ . Moreover, we can decompose the training error into the eigenspace of the NTK as

$$\begin{aligned} f(X, W(t)) - Y &= \sum_{i=1}^N (f(X, W(t)) - Y, q_i) q_i \\ &= \sum_{i=1}^N q_i^T (f(X, W(t)) - Y) q_i \\ &= \sum_{i=1}^N (e^{-\lambda_i t} q_i^T Y) q_i. \end{aligned}$$

Clearly, the network is biased to first learn the target function along the eigen-directions of the neural tangent kernel with larger eigenvalues, and then the rest components corresponding to smaller eigenvalues.

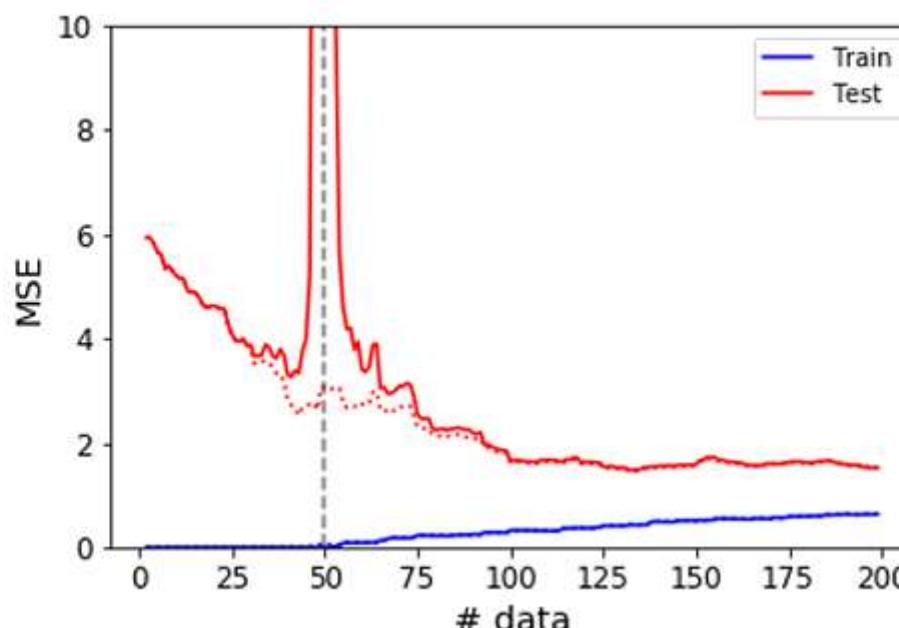
# Double Descent Phenomenon



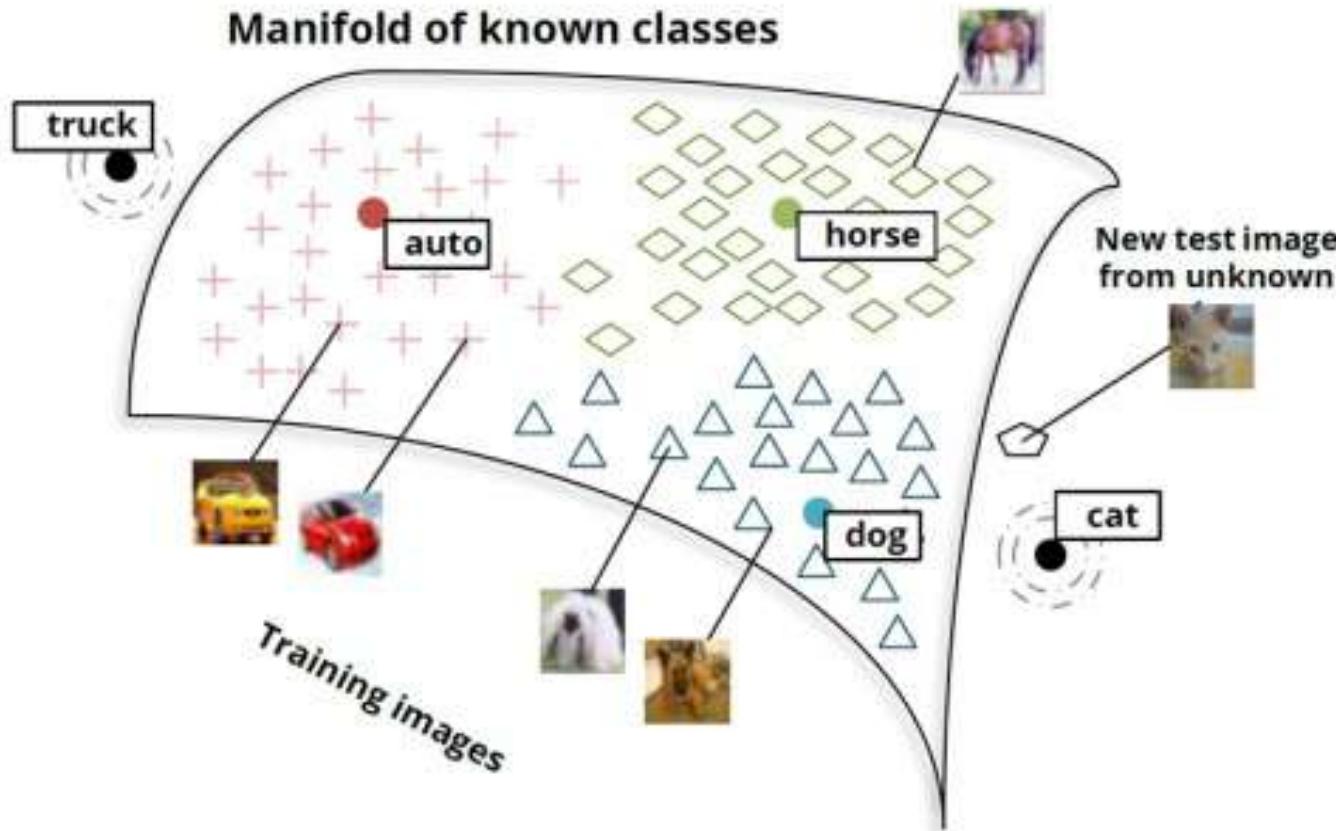
<https://openai.com/blog/deep-double-descent/>

## Intuition:

- **Case 1:**  $N \gg D$ . There's way more than enough data to pin down the optimal parameters, so it generalizes well.
- **Case 2:**  $N \approx D$ . It can memorize the training set, but just barely. It might need a large  $\|\mathbf{w}\|$  to do so.
- **Case 3:**  $N \ll D$ . It can fit the training set easily. The implicit regularization of gradient descent makes it do so with a small  $\|\mathbf{w}\|$ .



# Intuition: data manifold hypothesis

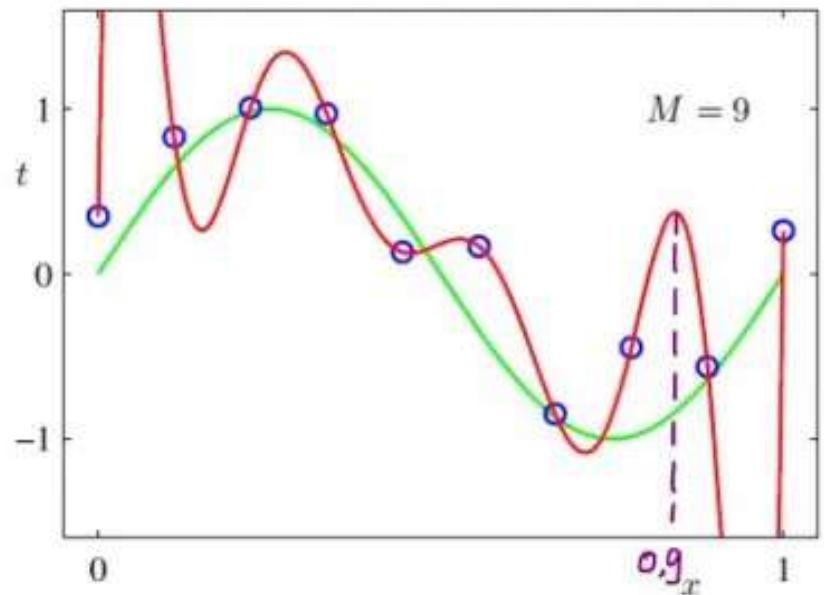


How can we understand if the network separator makes sense?

The network is capable to grasp the geometry of the manifold without overfitting.  
How is this possible?

Ansuini et al.

# Explicit regularization: $\| \cdot \|_p$ norm



$$h_a(x) = a_0 + a_1 x + a_2 x^2 + \cdots + a_m x^m$$

$$\mathcal{L}(w) = \sum_{i=1}^N (y_i - h(x_i))^2 + \lambda \|a\|_p$$

What is the effect of the regularization on the loss surface?

# Weight decay and hessian interpretation

- We often refer to  $\mathbf{H}$  as the **curvature** of a function.
- Suppose you move along a line defined by  $\theta + t\mathbf{v}$  for some vector  $\mathbf{v}$ .
- Second-order Taylor approximation:

$$\mathcal{J}(\theta + t\mathbf{v}) \approx \mathcal{J}(\theta) + t\nabla\mathcal{J}(\theta)^\top\mathbf{v} + \frac{t^2}{2}\mathbf{v}^\top\mathbf{H}(\theta)\mathbf{v}$$

- Hence, in a direction where  $\mathbf{v}^\top\mathbf{H}\mathbf{v} > 0$ , the cost function curves upwards, i.e. has **positive curvature**. Where  $\mathbf{v}^\top\mathbf{H}\mathbf{v} < 0$ , it has **negative curvature**.

# Hessian matrix

- The **Hessian matrix**, denoted  $\mathbf{H}$ , or  $\nabla^2 \mathcal{J}$  is the matrix of second derivatives:

$$\mathbf{H} = \nabla^2 \mathcal{J} = \begin{pmatrix} \frac{\partial^2 \mathcal{J}}{\partial \theta_1^2} & \frac{\partial^2 \mathcal{J}}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 \mathcal{J}}{\partial \theta_1 \partial \theta_D} \\ \frac{\partial^2 \mathcal{J}}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 \mathcal{J}}{\partial \theta_2^2} & \cdots & \frac{\partial^2 \mathcal{J}}{\partial \theta_2 \partial \theta_D} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 \mathcal{J}}{\partial \theta_D \partial \theta_1} & \frac{\partial^2 \mathcal{J}}{\partial \theta_D \partial \theta_2} & \cdots & \frac{\partial^2 \mathcal{J}}{\partial \theta_D^2} \end{pmatrix}$$

- It's a symmetric matrix because  $\frac{\partial^2 \mathcal{J}}{\partial \theta_i \partial \theta_j} = \frac{\partial^2 \mathcal{J}}{\partial \theta_j \partial \theta_i}$ .

Assuming no bias parameter, the objective function is:

$$\tilde{J}(w; X, y) = \frac{\alpha}{2} w^\top w + J(w; X, y)$$

with the gradient:

$$\nabla_w \tilde{J}(w; X, y) = \alpha w + \nabla_w J(w; X, y).$$

A single gradient update step is:

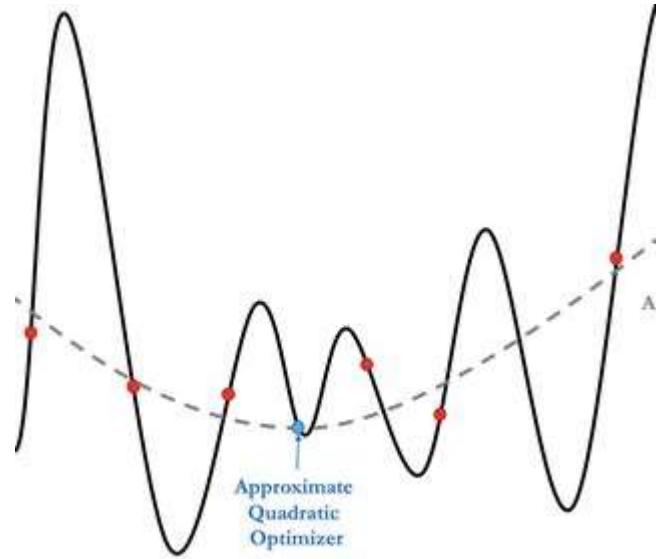
$$w \leftarrow w - \epsilon(\alpha w + \nabla_w J(w; X, y)),$$

or equivalently:

$$w \leftarrow (1 - \epsilon\alpha)w - \epsilon \nabla_w J(w; X, y).$$

This shows that weight decay modifies the learning rule by shrinking the weights by a factor of  $(1 - \epsilon\alpha)$  before the gradient update. Further analysis considers a quadratic approximation around the unregularized minimum  $w^*$ .

# A simplification



The approximation of the objective function  $\tilde{J}$  is given by:

$$\tilde{J}(\theta) = J(w^*) + \frac{1}{2}(w - w^*)^\top H(w - w^*),$$

where  $H$  is the Hessian matrix of  $J$  with respect to  $w$ , evaluated at  $w^*$ . Since  $w^*$  is defined to be a minimum, the first-order term vanishes. The Hessian  $H$  is positive semidefinite because  $w^*$  is a minimum.

The minimum of  $\tilde{J}$  occurs where its gradient:

$$\nabla_w \tilde{J}(w) = H(w - w^*) \quad *$$

is equal to zero.

To study the effect of weight decay, we modify equation (\*) by adding the weight decay gradient and solve for the minimum of the regularized objective function. The location of the minimum is represented by  $\tilde{w}$ . (the new minimum!)

Starting from:

$$\alpha\tilde{w} + H(\tilde{w} - w^*) = 0,$$

we can rewrite this as:

$$(H + \alpha I)\tilde{w} = Hw^*,$$

leading to:

$$\tilde{w} = (H + \alpha I)^{-1}Hw^*.$$

As  $\alpha \rightarrow 0$ ,  $\tilde{w}$  approaches  $w^*$ . To explore what happens as  $\alpha$  increases, we decompose  $H$  into  $Q\Lambda Q^\top$ , where  $Q$  is an orthonormal matrix of eigenvectors and  $\Lambda$  is a diagonal matrix of eigenvalues. Applying this decomposition:

$$\tilde{w} = (Q\Lambda Q^\top + \alpha I)^{-1}Q\Lambda Q^\top w^*,$$

simplifying to:

$$\tilde{w} = Q(\Lambda + \alpha I)^{-1}\Lambda Q^\top w^*.$$

$$\tilde{w} = Q(\Lambda + \alpha I)^{-1} \Lambda Q^\top w^*$$

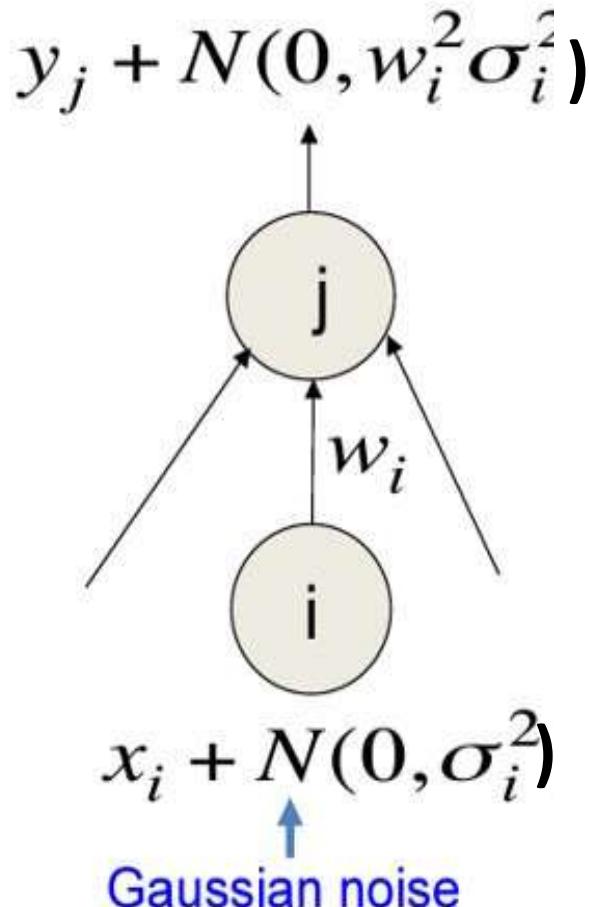
This shows that weight decay scales  $w^*$  along the eigenvector directions of  $H$ . Each component along the  $i$ -th eigenvector of  $H$  is rescaled by a factor of  $\frac{\lambda_i}{\lambda_i + \alpha}$ .

- If  $\lambda_i \gg \alpha$ , regularization has little effect.
- If  $\lambda_i \ll \alpha$ , components are shrunk significantly, approaching zero.

This explains how regularization influences directions with small eigenvalues, shrinking their magnitudes more aggressively.

# L2 regularization via noisy inputs

- Suppose we add Gaussian noise to the inputs.
  - The variance of the noise is amplified by the squared weight before going into the next layer.
- In a simple net with a linear output unit directly connected to the inputs, the amplified noise gets added to the output.
- This makes an additive contribution to the squared error.
  - So minimizing the squared error tends to minimize the squared weights when the inputs are noisy.



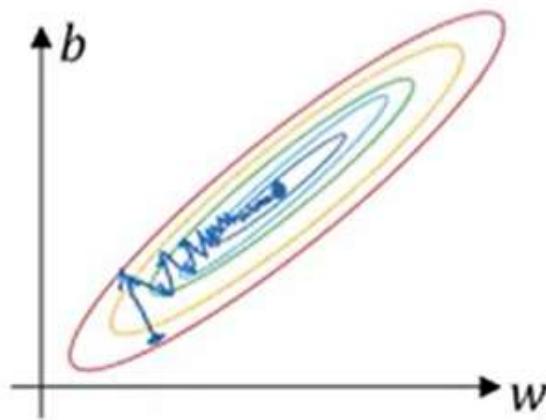
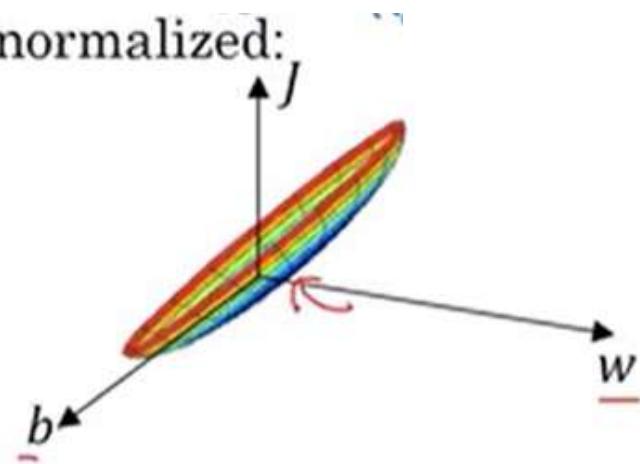
output on one case  $\rightarrow y^{noisy} = \sum_i w_i x_i + \sum_i w_i \varepsilon_i$  where  $\varepsilon_i$  is sampled from  $N(0, \sigma_i^2)$

$$\begin{aligned} E[(y^{noisy} - t)^2] &= E\left[\left(y + \sum_i w_i \varepsilon_i - t\right)^2\right] = E\left[\left((y - t) + \sum_i w_i \varepsilon_i\right)^2\right] \\ &= (y - t)^2 + E\left[2(y - t) \sum_i w_i \varepsilon_i\right] + E\left[\left(\sum_i w_i \varepsilon_i\right)^2\right] \\ &= (y - t)^2 + E\left[\sum_i w_i^2 \varepsilon_i^2\right] \quad \text{because } \varepsilon_i \text{ is independent of } \varepsilon_j \\ &\quad \text{and } \varepsilon_i \text{ is independent of } (y - t) \\ &= (y - t)^2 + \sum_i w_i^2 \sigma_i^2 \quad \text{So } \sigma_i^2 \text{ is equivalent to an L2 penalty} \end{aligned}$$

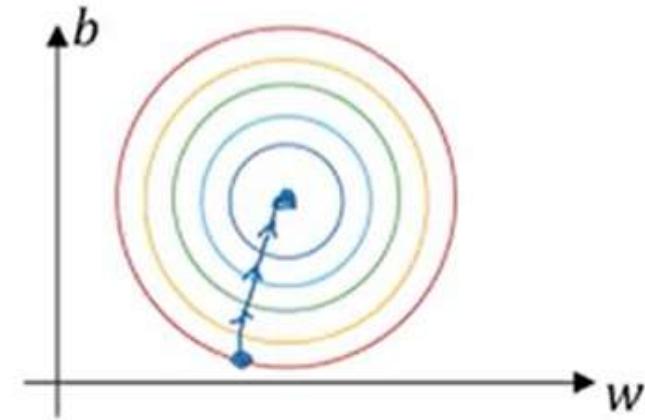
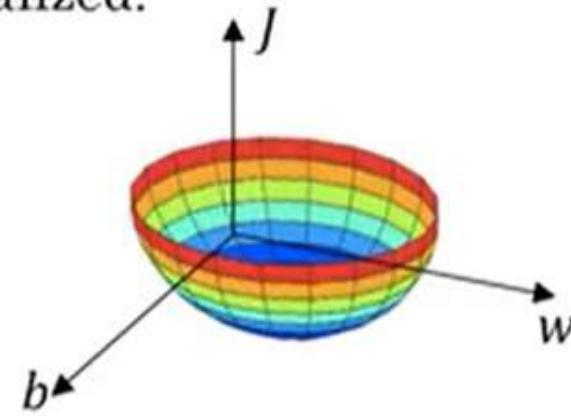
Is this true for more complex nets?

# Data normalization changes the geometry of the loss

Unnormalized:

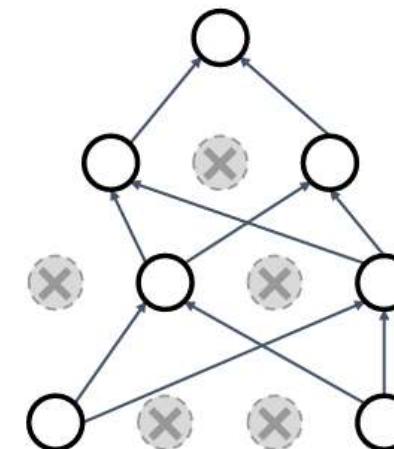
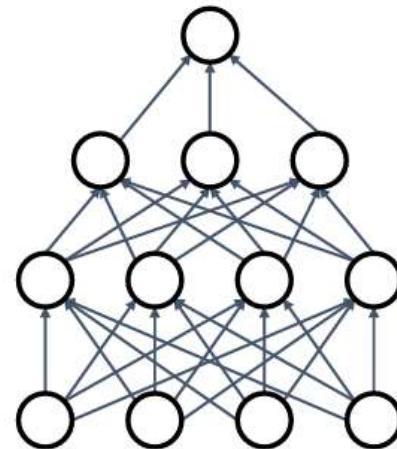


Normalized:

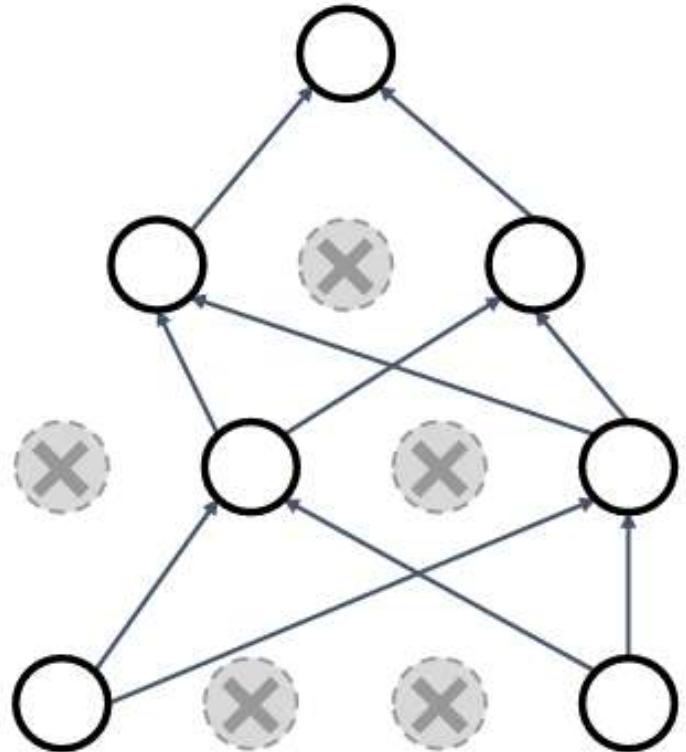


# Another example of explicit NNs regularization: Dropout

In each forward pass, randomly set some neurons to zero  
Probability of dropping is a hyperparameter; 0.5 is common



# Regularization: Dropout



Another interpretation:

Dropout is training a large **ensemble** of models (that share parameters).

Each binary mask is one model

An FC layer with 4096 units has  $2^{4096} \sim 10^{1233}$  possible masks!  
Only  $\sim 10^{82}$  atoms in the universe...

# Example with **regression**: minimum norm solution

$$A \in \mathbb{R}^{m \times n}, \quad x \in \mathbb{R}^n, \quad m \gg n$$
$$Ax = y$$

Parameters are more than data

one particular solution is

$$x_{\text{ln}} = A^T(AA^T)^{-1}y$$

Moore-Penrose pseudo-inverse  
(a possible algorithm)

( $AA^T$  is invertible since  $A$  full rank)

in fact,  $x_{\text{ln}}$  is the solution of  $y = Ax$  that minimizes  $\|x\|$

i.e.,  $x_{\text{ln}}$  is solution of optimization problem

$$\begin{array}{ll}\text{minimize} & \|x\| \\ \text{subject to} & Ax = y\end{array}$$

(with variable  $x \in \mathbb{R}^n$ )

# Example with regression: close form is minimum norm solution

suppose  $Ax = y$ , so  $A(x - x_{\text{ln}}) = 0$  and

$$\begin{aligned}(x - x_{\text{ln}})^T x_{\text{ln}} &= (x - x_{\text{ln}})^T A^T (AA^T)^{-1}y \\&= (A(x - x_{\text{ln}}))^T (AA^T)^{-1}y \\&= 0\end{aligned}$$

i.e.,  $(x - x_{\text{ln}}) \perp x_{\text{ln}}$ , so

$$\|x\|^2 = \|x_{\text{ln}} + x - x_{\text{ln}}\|^2 = \|x_{\text{ln}}\|^2 + \|x - x_{\text{ln}}\|^2 \geq \|x_{\text{ln}}\|^2$$

i.e.,  $x_{\text{ln}}$  has smallest norm of any solution

# Regression: minimum norm solution with GD

Let cost function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be defined by

$$f(x) := \frac{1}{2} \|Ax - b\|_2^2 \quad A \in \mathbb{R}^{m \times n}, \quad x \in \mathbb{R}^n, \quad m \gg n$$

whose gradient is

$$\nabla f(x) = A^\top (Ax - b)$$

Using gradient descent with step  $\mu > 0$ ,

$$\begin{aligned} x_{k+1} &= x_k - \mu \nabla f(x_k) \\ &= (I - \mu A^\top A)x_k + \mu A^\top b \end{aligned}$$

Hence,

$$x_k = (I - \mu A^\top A)^k x_0 + \mu \sum_{\ell=0}^{k-1} (I - \mu A^\top A)^\ell A^\top b$$

$$A = U\Sigma V^\top = U \begin{bmatrix} \Sigma_1 & 0 \end{bmatrix} \begin{bmatrix} V_1^\top \\ V_2^\top \end{bmatrix} = U\Sigma_1 V_1^\top$$

Letting  $y := V^\top x$ , we rewrite

$$\begin{aligned} y_k &= (I - \mu \Sigma^\top \Sigma)^k y_0 + \mu \sum_{\ell=0}^{k-1} (I - \mu \Sigma^\top \Sigma)^\ell \Sigma^\top U^\top b \\ &= \begin{bmatrix} (I - \mu \Sigma_1^2)^k & 0 \\ 0 & I \end{bmatrix} y_0 + \mu \sum_{\ell=0}^{k-1} \begin{bmatrix} (I - \mu \Sigma_1^2)^\ell & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} \Sigma_1 \\ 0 \end{bmatrix} U^\top b \\ &= \begin{bmatrix} (I - \mu \Sigma_1^2)^k & 0 \\ 0 & I \end{bmatrix} y_0 + \mu \sum_{\ell=0}^{k-1} \begin{bmatrix} (I - \mu \Sigma_1^2)^\ell \Sigma_1 \\ 0 \end{bmatrix} U^\top b \end{aligned}$$

# Example with regression: minimum norm solution with GD

Choosing  $\mu > 0$  such that all eigenvalues of  $\mathbf{I} - \mu \Sigma_1^2$  are strictly inside the unit circle, then  $\mathbf{y}_k \rightarrow \mathbf{y}_\infty$ , where

$$\mathbf{y}_\infty = \begin{bmatrix} \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} \end{bmatrix} \mathbf{y}_0 + \mu \sum_{\ell=0}^{\infty} \begin{bmatrix} (\mathbf{I} - \mu \Sigma_1^2)^\ell \Sigma_1 \\ \mathbf{O} \end{bmatrix} \mathbf{U}^\top \mathbf{b}$$

where

$$\mu \sum_{\ell=0}^{\infty} (\mathbf{I} - \mu \Sigma_1^2)^\ell \Sigma_1 = \mu (\mathbf{I} - \mathbf{I} + \mu \Sigma_1^2)^{-1} \Sigma_1 = \Sigma_1^{-1}$$

and, thus,

$$\mathbf{y}_\infty = \begin{bmatrix} \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} \end{bmatrix} \mathbf{y}_0 + \begin{bmatrix} \Sigma_1^{-1} \\ \mathbf{O} \end{bmatrix} \mathbf{U}^\top \mathbf{b}$$

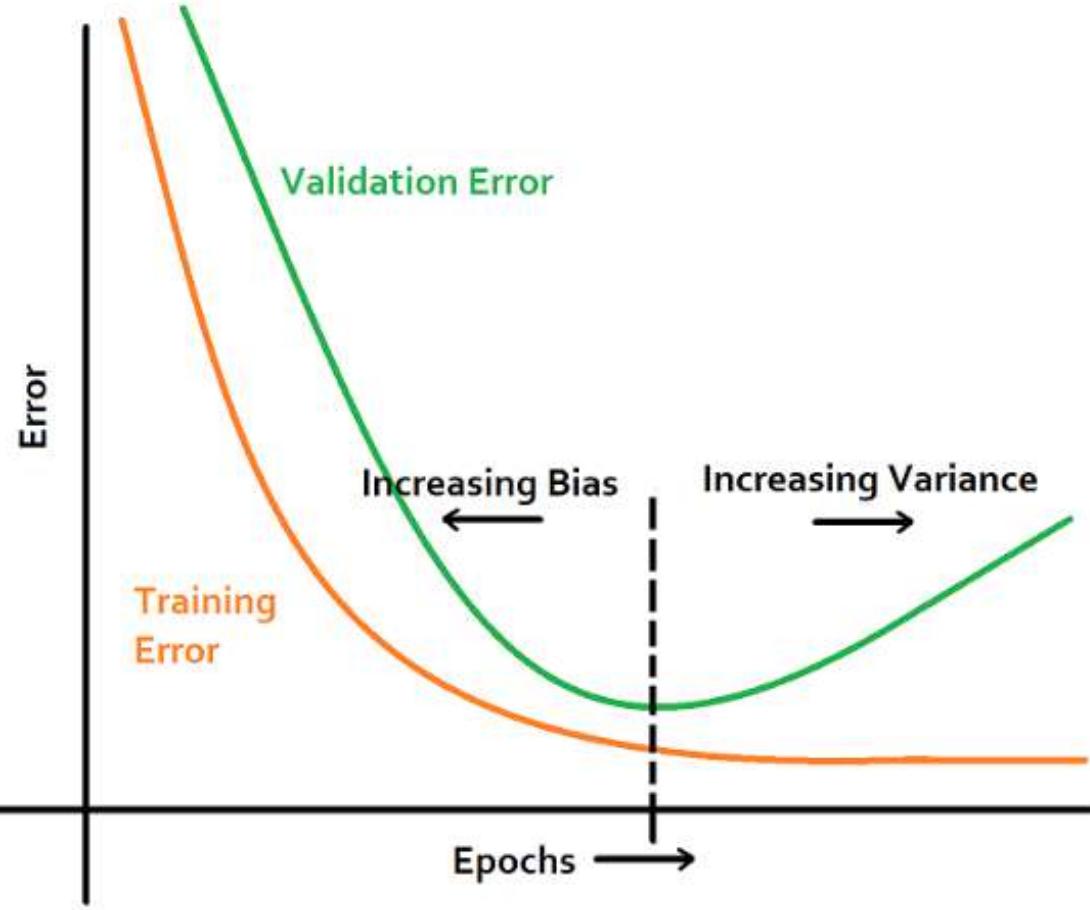
Since  $\mathbf{x} := \mathbf{V}\mathbf{y}$ ,

This is the pseudo-inverse

$$\mathbf{x}_\infty = \mathbf{V}_2 \mathbf{V}_2^\top \mathbf{x}_0 + \underbrace{\mathbf{V}_1 \Sigma_1^{-1} \mathbf{U}^\top \mathbf{b}}_{=\mathbf{x}_{LN}}$$

Therefore, we conclude that if  $\mathbf{x}_0$  is orthogonal to the null space of  $\mathbf{A}$ , then gradient descent will converge to the least-norm solution.

# Early stopping, L2 norm and learning rate



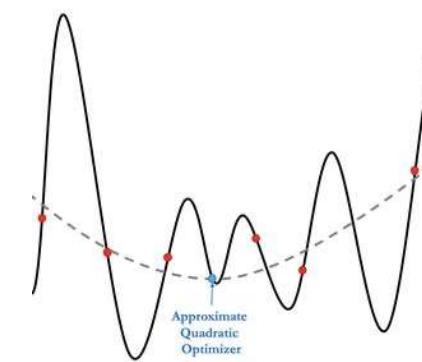
From Taylor series expansion, we will make a quadratic approximation around the empirically optimal value of the weights  $\mathbf{w}^*$ . The matrix  $\mathbf{H}$  is the hessian matrix.

$$\hat{J}(\boldsymbol{\theta}) = J(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}^*),$$

$$\nabla_{\mathbf{w}} \hat{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w} - \mathbf{w}^*).$$

Since  $\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w})$

$$\nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y}) \Big|_{\mathbf{w}^*} = 0$$



The gradient of this approximation comes out to be a linear function of  $H$  and  $w$ .  
 Now, we can calculate the weights at each step of the gradient descent.

$$\begin{aligned}\mathbf{w}^{(\tau)} &= \mathbf{w}^{(\tau-1)} - \epsilon \nabla_{\mathbf{w}} \hat{J}(\mathbf{w}^{(\tau-1)}) \\ &= \mathbf{w}^{(\tau-1)} - \epsilon \mathbf{H}(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*) \\ \mathbf{w}^{(\tau)} - \mathbf{w}^* &= (\mathbf{I} - \epsilon \mathbf{H})(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*).\end{aligned}$$

Since the Hessian matrix is real, symmetric and positive semi-definite. By eigenvalue decomposition. matrix  $H$  can be written as,

$$\begin{aligned}\mathbf{H} &= \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^\top \\ \mathbf{w}^{(\tau)} - \mathbf{w}^* &= (\mathbf{I} - \epsilon \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^\top)(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*) \\ \mathbf{Q}^\top(\mathbf{w}^{(\tau)} - \mathbf{w}^*) &= (\mathbf{I} - \epsilon \boldsymbol{\Lambda}) \mathbf{Q}^\top (\mathbf{w}^{(\tau-1)} - \mathbf{w}^*)\end{aligned}$$

Assuming  $\epsilon$  is chosen to be small and  $w(0) = 0$ . The weights after  $\tau$  iterations are:

$$\mathbf{Q}^\top \mathbf{w}^{(\tau)} = [\mathbf{I} - (\mathbf{I} - \epsilon \boldsymbol{\Lambda})^\tau] \mathbf{Q}^\top \mathbf{w}^* \quad (1)$$

The optimal weights obtained on imposing L2 norm penalty on the objective function is given by ([refer](#)):

$$\begin{aligned} \mathbf{Q}^\top \tilde{\mathbf{w}} &= (\Lambda + \alpha \mathbf{I})^{-1} \Lambda \mathbf{Q}^\top \mathbf{w}^* \\ \mathbf{Q}^\top \tilde{\mathbf{w}} &= [\mathbf{I} - (\Lambda + \alpha \mathbf{I})^{-1} \alpha] \mathbf{Q}^\top \mathbf{w}^* \end{aligned} \quad (2)$$

$$1 - \frac{\alpha}{\Lambda + \alpha I} = \frac{\Lambda + \alpha I - \alpha I}{\Lambda + \alpha I} = \frac{\Lambda}{\Lambda + \alpha I}$$

$$\log(1 - x) \approx x$$

Comparing Eq(1) and Eq(2), we derive the following relation between  $\tau$ ,  $\mathbf{w}$ , and  $\epsilon$ . The equivalence can be seen between the L2 regularization and early stopping.

$$(\mathbf{I} - \epsilon \Lambda)^\tau = (\Lambda + \alpha \mathbf{I})^{-1} \alpha,$$

Taking log both sides and using the series approximation of  $\log(1+x)$ , we can conclude that if all  $\lambda_i$  are small (that is,  $\epsilon \lambda_i \ll 1$  and  $\lambda_i/\alpha \ll 1$ ) then the following equation holds.

$$\begin{aligned} \tau &\approx \frac{1}{\epsilon \alpha}, \\ \alpha &\approx \frac{1}{\tau \epsilon}. \end{aligned}$$

$$\tau \log(I - \epsilon \Lambda) = \log(\alpha) - \log(\Lambda + \alpha I)$$

$$\tau \log(I - \epsilon \Lambda) = \log(\alpha) - \log(\alpha(\frac{\Lambda}{\alpha} + I))$$

$$\tau \log(I - \epsilon \Lambda) = -\log(\frac{\Lambda}{\alpha} + I)$$

$$\tau \epsilon \Lambda = \frac{\lambda}{\alpha}$$

Here,  $\alpha$  is the regularization constant,  $\tau$  is no. of iterations, and  $\epsilon$  is the learning rate.

Increasing no. of epochs/iterations  $\tau$  is equivalent to reducing the regularization constant. Similarly, early stopping the model, i.e. reducing the no. of iterations is similar to L2 regularization with large  $\alpha$ . Thus, we can say that early stopping regularizes the model.

**References:**

Deep Learning (Adaptive Computation and Machine Learning series) — By Ian Goodfellow, Yoshua Bengio, Aaron Courville

# Weights initialization

Weight Initialization: Xavier Initialization

"Xavier" initialization:  
std = 1/sqrt(Din)

**Derivation:** Variance of output = Variance of input

$$y = Wx$$

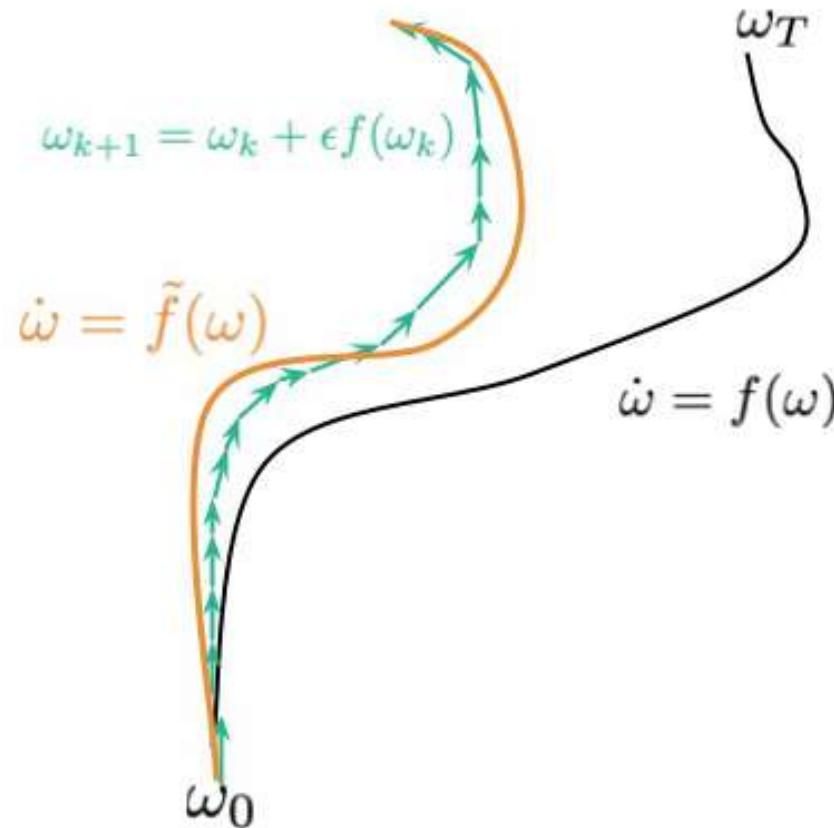
$$y_i = \sum_{j=1}^{Din} x_j w_j$$

$$\begin{aligned} \text{Var}(y_i) &= \text{Din} * \text{Var}(x_i w_i) && [\text{Assume } x, w \text{ are iid}] \\ &= \text{Din} * (\text{E}[x_i^2] \text{E}[w_i^2] - \text{E}[x_i]^2 \text{E}[w_i]^2) && [\text{Assume } x, w \text{ independent}] \\ &= \text{Din} * \text{Var}(x_i) * \text{Var}(w_i) && [\text{Assume } x, w \text{ are zero-mean}] \end{aligned}$$

If  $\text{Var}(w_i) = 1/\text{Din}$  then  $\text{Var}(y_i) = \text{Var}(x_i)$

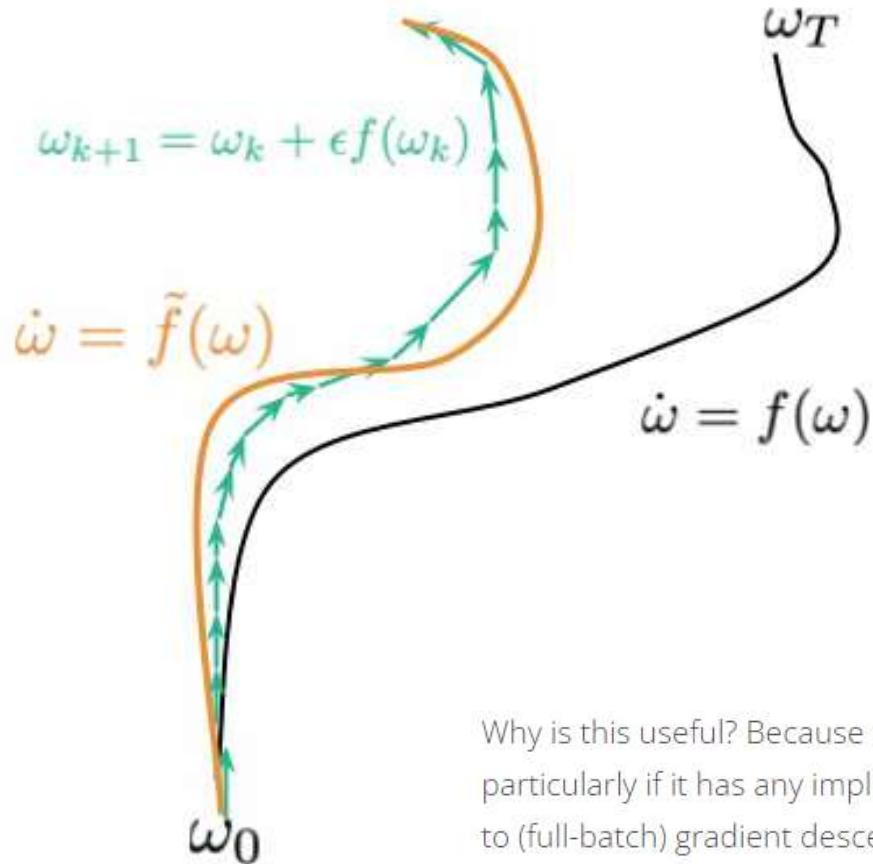
**Effectively reduce the complexity of the model!!**

# Discretization in GD and implicit regularization



Let's say we have a differential equation  $\dot{\omega} = f(\omega)$ . The solution to this ODE with initial condition  $\omega_0$  is a continuous trajectory  $\omega_t$ , shown in the image in black. We usually can't compute this solution in closed form, and instead simulate the ODE using the Euler's method,  $\omega_{k+1} = \omega_k + \epsilon f(\omega_k)$ . This results in a discrete trajectory shown in teal. Due to discretization error, for finite stepsize  $\epsilon$ , this discrete path may not lie exactly where the continuous black path lies. Errors accumulate over time, as shown in this illustration. The goal of backward error analysis is to find a different ODE,  $\dot{\omega} = \tilde{f}(\omega)$  such that the approximate discrete path we got from Euler's method lies near the the continuous path which solves this new ODE. Our goal is to reverse engineer a modified  $\tilde{f}$  such that the discrete iteration can be well-modelled by an ODE.

See also: <https://arxiv.org/pdf/2101.12176.pdf> , <https://www.inference.vc/notes-on-the-origin-of-implicit-regularization-in-stochastic-gradient-descent/>



### IMPLICIT GRADIE REGULARIZATION DUE TO DISCRETIZATION

Why is this useful? Because the form  $\tilde{f}$  takes can reveal interesting aspects of the behaviour of the discrete algorithm, particularly if it has any implicit bias towards moving into different areas of the space. When the authors apply this technique to (full-batch) gradient descent, it already suggests the kind of implicit regularization bias gradient descent has.

In Gradient descent with a cost function  $C$ , the original ODE is  $f(\omega) = -\nabla C(\omega)$ . The modified ODE which corresponds to a finite stepsize  $\epsilon$  takes the form  $\dot{\omega} = -\nabla \tilde{C}_{GD}(\omega)$  where

$$\tilde{C}_{GD}(\omega) = C(\omega) + \frac{\epsilon}{4} \|\nabla C(\omega)\|^2$$

## Problem Setup

Given a continuous flow represented by an ordinary differential equation (ODE) of the form:

$$\dot{\omega} = f(\omega)$$

we are interested in finding a modified loss function that accounts for the discretization error when integrating this ODE using a discrete method, like Euler's method.

## Discretization Error

For simplicity, let's assume that the ODE is solved using a discrete update rule, such as the explicit Euler method:

$$\omega_{i+1} = \omega_i + \epsilon f(\omega_i)$$

where  $\epsilon$  is the step size (e.g., the learning rate in gradient descent). This discrete step introduces an approximation error because the true solution of the continuous flow is not exactly reached by the discrete step.

## Backward Error Analysis

The goal is to find a modified flow (and hence, a modified loss function) such that the continuous flow with this modified loss matches the discrete flow to a given order in  $\epsilon$ . This means we want to express the continuous flow as:

$$\dot{\omega} = f_{\text{mod}}(\omega) = f(\omega) + \epsilon f_1(\omega) + \epsilon^2 f_2(\omega) + \dots$$

where the correction terms  $f_1(\omega), f_2(\omega), \dots$  account for the discretization error at different orders of  $\epsilon$ .

## Derivation of the First Correction Term

We start by expanding the discrete step  $\omega_{i+1}$  in a Taylor series about  $\omega_i$ :

$$\omega_{i+1} = \omega_i + \epsilon \dot{\omega}_i + \frac{\epsilon^2}{2} \ddot{\omega}_i + \mathcal{O}(\epsilon^3)$$

where  $\dot{\omega}_i = f(\omega_i)$  and  $\ddot{\omega}_i$  is the time derivative of  $\dot{\omega}_i$ . Using the chain rule,  $\ddot{\omega}_i$  is:

$$\ddot{\omega}_i = \frac{d}{dt} f(\omega) = \nabla f(\omega) \cdot \dot{\omega} = \nabla f(\omega) f(\omega)$$

Now we match this Taylor expansion to the discrete update  $\omega_{i+1} = \omega_i + \epsilon f(\omega_i)$ , which only has terms of order  $\epsilon$ . To ensure consistency between the continuous and discrete flows, we need to add correction terms to the continuous flow at higher orders in  $\epsilon$ .

At order  $\epsilon^2$ , the correction term  $f_1(\omega)$  is derived by requiring that the Taylor expansion of the continuous flow matches the discrete update. This leads to the condition:

$$f_1(\omega) + \frac{1}{2} \nabla f(\omega) f(\omega) = 0$$

Solving for  $f_1(\omega)$ , we obtain:

$$f_1(\omega) = -\frac{1}{2} \nabla f(\omega) f(\omega)$$

## General Correction to the Loss

To interpret this result in terms of a correction to the **loss function**, we note that the flow  $\dot{\omega} = f(\omega)$  is typically derived as the gradient of a loss function  $C(\omega)$ :

$$f(\omega) = -\nabla C(\omega)$$

Substituting this into the expression for  $f_1(\omega)$ , we find that the correction term becomes:

$$f_1(\omega) = -\frac{1}{2}\nabla\nabla C(\omega)\nabla C(\omega)$$

This is exactly the gradient of  $\frac{1}{4}||\nabla C(\omega)||^2$ , meaning the corrected flow is:

$$\dot{\omega} = -\nabla C(\omega) - \frac{\epsilon}{4}\nabla||\nabla C(\omega)||^2$$

Thus, the **corrected loss function** to match the discrete flow to the continuous flow up to order  $\epsilon$  is:

$$C_{\text{mod}}(\omega) = C(\omega) + \frac{\epsilon}{4}||\nabla C(\omega)||^2$$

# Discretization in GD and implicit regularization

Intuition:

$$\begin{aligned}\nabla C(\omega + (\epsilon/2)\nabla C) &\approx \nabla C + (\epsilon/2)\nabla \nabla C \nabla C \\ &= \nabla C + (\epsilon/4)\nabla (||\nabla C||^2)\end{aligned}$$

On average, the gradient  
is half a step stale

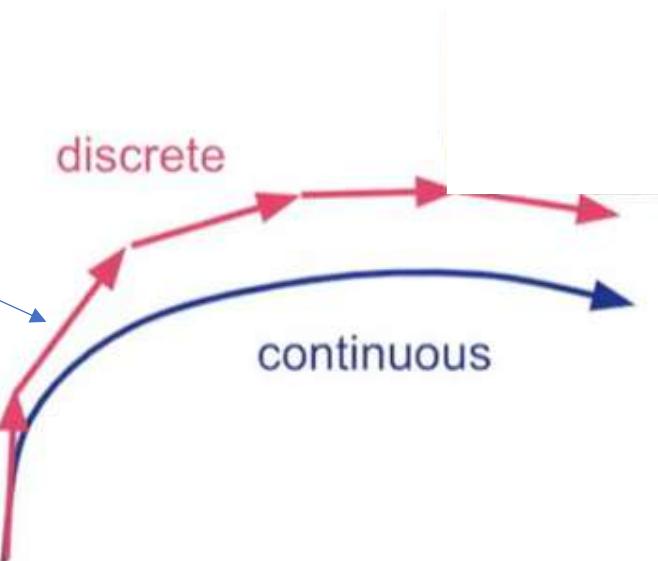
Taylor expansion

For "small finite  
learning rates":

$$\begin{aligned}\dot{\omega} &= -\nabla \tilde{C}_{GD}(\omega) + O(\epsilon^2) \\ \tilde{C}_{GD}(\omega) &= C(\omega) + (\epsilon/4)||\nabla C(\omega)||^2\end{aligned}$$

Original loss

Implicit regularizer



$$\mathcal{L}(w, X) = \sum_i l_i(w, x_i)$$

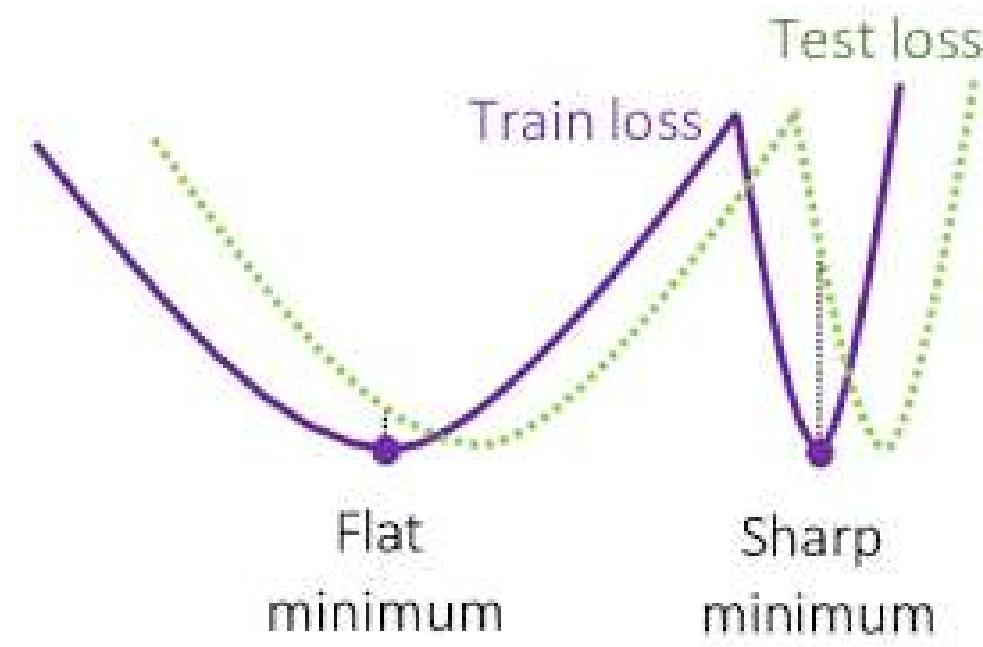
$$w_{t+1} = w_t - \eta \nabla_w \ell_i \qquad \text{Stochastic GD}$$

$$w_{t+1} = w_t - \eta \sum_{i \in I} \nabla_w \ell_i \qquad \text{Mini-batch GD}$$

$$w_{t+1} = w_t - \eta \sum_i \nabla_w \ell_i \qquad \text{Full-batch GD}$$

# SGD and generalization

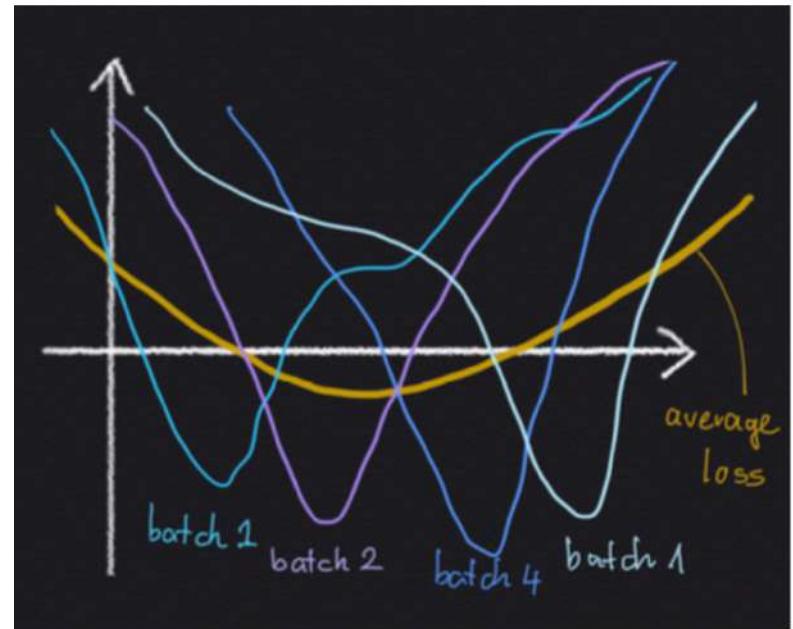
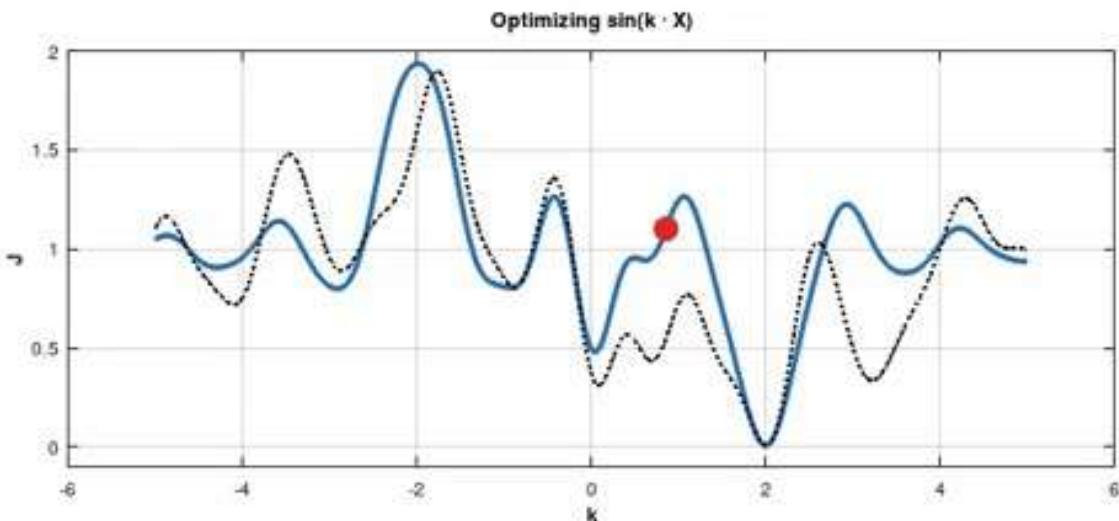
- Skip bad local minima, converge to flat minima (basins of attraction). Good for **generalization, robustness**.



Another advantage?

Can we be more precise?

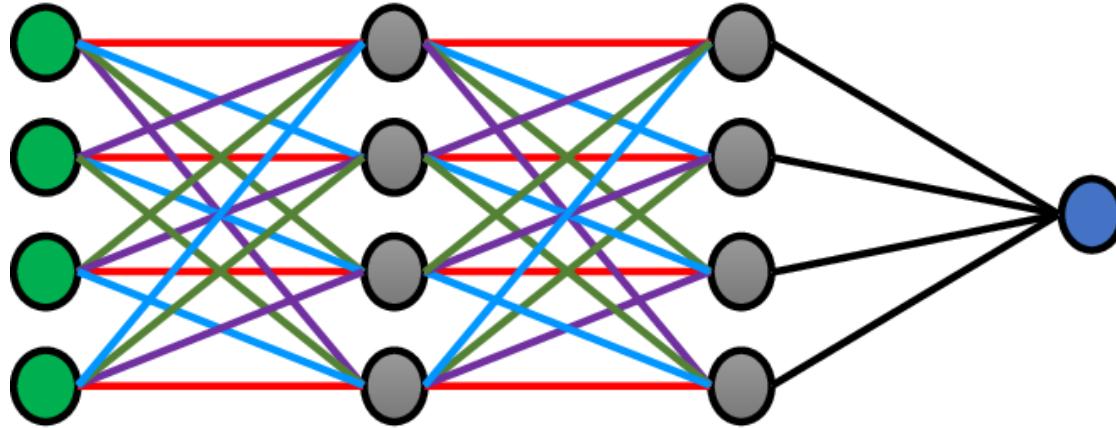
# Intuition



At each step the SGD is looking to a “different, sampled” function

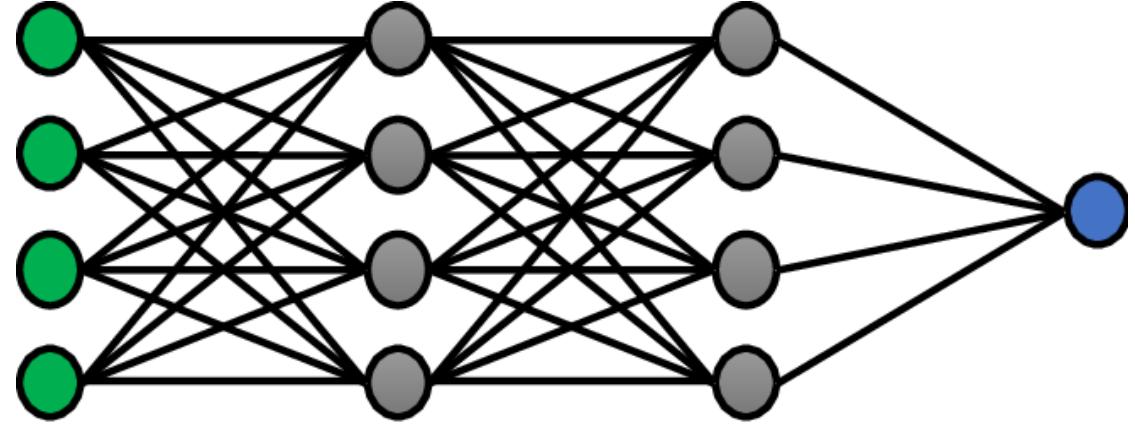
On average the main effect is dominated by big “basins of attraction”

## Weights bias for fully connected and convolutions linear networks



(b) Convolutional network of depth  $L$

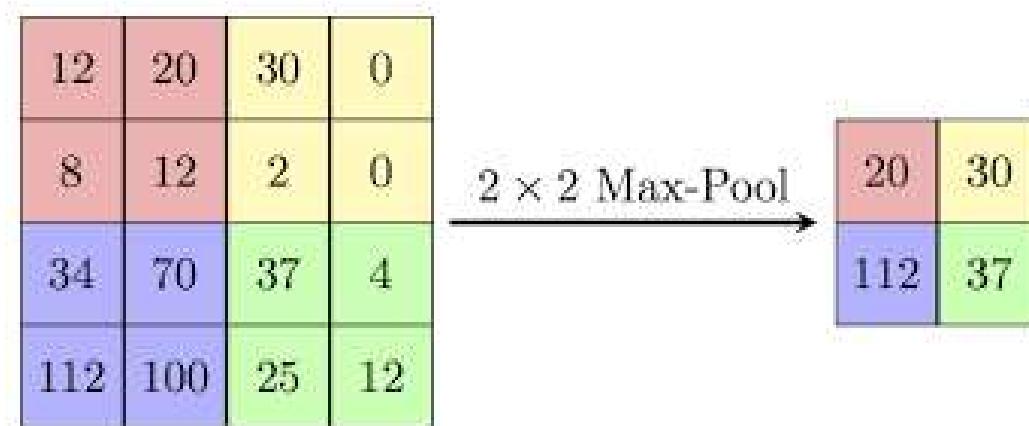
$$\bar{\beta}^\infty \propto \text{first order stationary point of } \underset{\forall n, y_n \langle \mathbf{x}_n, \beta \rangle \geq 1}{\operatorname{argmin}} \|\hat{\beta}\|_{2/L}$$



(a) Fully connected network of depth  $L$

$$\bar{\beta}^\infty \propto \underset{\forall n, y_n \langle \mathbf{x}_n, \beta \rangle \geq 1}{\operatorname{argmin}} \|\beta\|_2 \text{ (independent of } L)$$

# Pooling operation and invariance



# Convolution operation, pooling and invariance: bias on the target function

$$\langle g_i w, x \rangle = (x * w)_i$$

$$f(x) = \sum_i \langle g_i w, x \rangle$$

$$f(gx) = f(x)$$

$$f(gx) = \sum_i \langle g_i w, gx \rangle = \sum_i \langle g^T g_i w, x \rangle = \sum_i \langle g'_i w, x \rangle$$

$$g^T g_i = g'_i$$

$$f(gx) = \begin{bmatrix} \sigma \langle g_1 w, gx \rangle \\ \vdots \\ \sigma \langle g_N w, gx \rangle \end{bmatrix} = \begin{bmatrix} \sigma \langle w, g_1^T gx \rangle \\ \vdots \\ \sigma \langle w, g_N^T gx \rangle \end{bmatrix} = P_g \begin{bmatrix} \sigma \langle w, gx \rangle \\ \vdots \\ \sigma \langle w, gx \rangle \end{bmatrix} = P_g f(x)$$

$$f(gx) = \Pi_g f(x)$$

Cnns are translation equivariant

$$\phi(T_\alpha x)(\beta) = \sigma\left(\int k(\beta - \gamma)T_\alpha x(\gamma)d\gamma\right) = T_\alpha\phi(x)(\beta)$$

Let's prove it

# Equivariance and convolutionality are one to one:

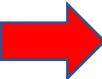
## On the Generalization of Equivariance and Convolution in Neural Networks to the Action of Compact Groups

---

Risi Kondor<sup>1</sup> Shubhendu Trivedi<sup>2</sup>

### Abstract

Convolutional neural networks have been extremely successful in the image recognition domain because they ensure equivariance to translations. There have been many recent attempts to generalize this framework to other domains, including graphs and data lying on manifolds. In this paper we give a rigorous, theoretical treatment of convolution and equivariance in neural networks with respect to not just translations, but the action of any compact group. Our main result is to prove that (given some natural constraints) convolutional structure is not just a sufficient, but also a necessary condition for equivariance to the action of a compact group. Our exposition makes use of concepts from representation theory and noncommutative harmonic analysis and derives new generalized convolution formulae.



Recently, there has been considerable interest in extending neural networks to more exotic types of data, such as graphs or functions on manifolds (Niepert et al., 2016; Defferrard et al., 2016; Duvenaud et al., 2015; Li et al., 2016; Cohen et al., 2018; Monti et al., 2017; Masci et al., 2015). In these domains, equivariance and multiscale structure are just as important as for images, but finding the right notion of convolution is not obvious.

On the other hand, mathematics does offer a sweeping generalization of convolution tied in deeply with some fundamental ideas of abstract algebra: if  $G$  is a compact group and  $f$  and  $g$  are two functions  $G \rightarrow \mathbb{C}$ , then the convolution of  $f$  with  $g$  is defined

$$(f * g)(u) = \int_G f(uv^{-1}) g(v) d\mu(v). \quad (1)$$

Note the striking similarity of this formula to the ordinary notion of convolution, except that in the argument of  $f$ ,  $u - v$  has been replaced by the group operation  $uv^{-1}$  and

How do we learn invariance or equivariance  
w.r.t. generic label-preserving transformations?

# Effect of Data Augmentation: invariant loss for one layer network

$$\mathcal{L}(w) = \frac{1}{N} \sum_{i=1}^N \int \ell(\sigma\langle w, g_\theta x_i \rangle; y_i) d\theta$$

$$\begin{aligned}\mathcal{L}(g_{\bar{\theta}} w) &= \frac{1}{N} \sum_{i=1}^N \int \ell(\sigma\langle w, g^* g_\theta x_i \rangle; y_i) d\theta \\ &= \sum_{i=1}^N \int \ell(\sigma\langle w, g_\theta x_i \rangle; y_i) d\theta \\ &= \mathcal{L}(w)\end{aligned}$$

The Loss is an  
invariant function!

$$g_\theta g_{\theta'} = g_{\theta''}$$

## Intuition with translation invariant/robust functions

- **Invariance:** Let  $\underline{f} : R \rightarrow R$  s.t.  $\underline{f}(w + q) = \underline{f}(w)$ ,  $\forall q \in R$ .  $\underline{f}$  is constant and  $\hat{\underline{f}}(0)$  is the only non zero Fourier transform.

$$\underline{f}(w) = \int_{-\infty}^{\infty} dt f(w - t), \quad f : R \rightarrow R.$$

- **Robustness:** Suppose we relax this invariance properties asking for robustness/approximate invariance. One way is to integrate in the interval  $[-a, a]$ .

$$\underline{f}(w) = \int_{-a}^a dt f(w - t) = \int_{-\infty}^{+\infty} dt \text{Ind}_{[-a,a]}(t) f(w - t).$$

Taking the Fourier and using  $\widehat{f(\cdot - t)} = e^{ikt} \hat{f}(k)$

$$\hat{\underline{f}}(k) = \left( \int_{-\infty}^{\infty} dt e^{ikt} \text{Ind}_{[-a,a]}(t) \right) \hat{f}(k) = 2a \text{sinc}(2ka) \hat{f}(k).$$

# Effect on the gradient descent

From the group structure we have:

$$\mathcal{L}(U_\alpha w; X, y) = \mathcal{L}(w; X, y), \quad \alpha \in R.$$

Then

$$\frac{d}{d\alpha} \mathcal{L}(U_\alpha w) = 0 \rightarrow \langle \nabla_w \mathcal{L}(w), \partial_\alpha U_\alpha w \rangle = 0,$$

i.e. **some directions in the gradient are not allowed by the invariances of the Loss.**

This is important since:

$$\frac{dw}{dt} = \nabla \mathcal{L}(w).$$

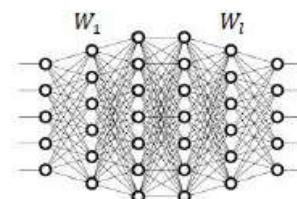
Learned weights reflects the symmetries of  
the data augmentation

Goal:  $y_i = \varphi_W(x_i)$ ,  $\varphi_{\mathbf{W}}(\mathbf{g}\mathbf{x}_i) = \varphi_{\mathbf{W}}(\mathbf{x}_i)$

## Data augmentation:

$$X = \left\{ \begin{array}{c} 00000000000000000000 \\ 111111111111111111 \\ 222222222222222222 \\ 333333333333333333 \\ 444444444444444444 \\ 555555555555555555 \\ 666666666666666666 \\ 777777777777777777 \\ 888888888888888888 \\ 999999999999999999 \end{array} \right\}, Y = \left\{ \begin{array}{c} 0 \\ 1 \\ 2 \\ 3 \\ \dots \\ \dots \\ 9 \end{array} \right\}$$

## Rotation augmented



$$\mathcal{L}(w) = \frac{1}{N} \sum_{i=1}^N \int \ell(\sigma\langle w, g_\theta x_i \rangle; y_i) d\theta$$

# Equivariance and learned weights

*Theorem 1 (see also, [43]):* Consider the Loss in eq. (5). The gradient of the Loss is an equivariant function in the  $\mathcal{G}$  transformations, i.e.:

$$\nabla_W \mathcal{L}_a(gW; X, y) = g \nabla_W \mathcal{L}_a(W; X, y), \quad \forall g \in \mathcal{G}.$$

*Proof 1:* Taking the gradient of eq. (5) we have:

$$\begin{aligned} & \nabla_W \mathcal{L}_a(W; X, y) = \\ &= \frac{1}{Q} \sum_{i=1}^R \sum_{j=1}^{|\mathcal{G}|} \ell'(v^T \sigma \langle W, g_j x_i \rangle, y_i) v^T \sigma' \langle W, g_j x_i \rangle g_j x_i \end{aligned}$$

Calculating  $\nabla_W \mathcal{L}_a(gW; X, y)$  we have:

$$\begin{aligned} & \nabla_W \mathcal{L}_a(gW; X, y) = \\ &= \frac{1}{Q} \sum_{i=1}^R \sum_{j=1}^{|\mathcal{G}|} \ell'(v^T \sigma \langle gW, g_j x_i \rangle, y_i) v^T \sigma' \langle gW, g_j x_i \rangle g_j x_i \\ &= \frac{1}{Q} \sum_{i=1}^R \sum_{j=1}^{|\mathcal{G}|} \ell'(v^T \sigma \langle W, g^T g_j x_i \rangle, y_i) v^T \sigma' \langle W, g^T g_j x_i \rangle g_j x_i \\ &= \frac{1}{Q} \sum_{i=1}^R \sum_{j=1}^{|\mathcal{G}|} \ell'(v^T \sigma \langle W, \hat{g}_j x_i \rangle, y_i) v^T \sigma' \langle W, \hat{g}_j x_i \rangle g \hat{g}_j x_i = \\ &= g \nabla_W \mathcal{L}_a(W; X, y) \end{aligned}$$

where in the fourth line we redefined  $g^T g_j = \hat{g}_j$  to get  $g_j = g \hat{g}_j$ . The key property here is the closure of the group transformations. ■

# Equivariance and learned weights

*Corollary 1:* Consider the Loss in eq. 5 i.e. the case of data augmentation with full orbits, i.e.  $gX = X$  for all  $g \in \mathcal{G}$ , and suppose we initialize the network with symmetric weights i.e.  $gW_0 = W_0$  for all  $g \in \mathcal{G}$ . Then the learned weights are symmetric  $W_T = gW_T$ .

*Proof 2:* We first notice that due to the property  $gX = X$ ,  $\nabla_W \mathcal{L}_a(W; X, y) = \nabla_W \mathcal{L}_a(W; gX, y)$ . Then by the dot product structure of the Loss we can see the transformation as applied to the weights i.e.  $\langle W, gx \rangle = \langle g^T W, x \rangle$ . Thus, using the equivariance property in theorem 1, we have that for all  $g \in \mathcal{G}$ :

$$\begin{aligned}\nabla_W \mathcal{L}_a(W; X, y) &= \nabla_W \mathcal{L}_a(W; gX, y) \\ &= \nabla_W \mathcal{L}_a(g^T W; X, y) \\ &= g^T \nabla_W \mathcal{L}_a(W; gX, y).\end{aligned}$$

Therefore each weights update in the gradient descent is group invariant and as a consequence the learned weight at time  $T$  is invariant, being:

$$W_T = W_0 - \gamma \sum_{k=1}^T \nabla_W \mathcal{L}_a(W_k; X, y).$$

■

# Permutation equivariant layer

Our goal is to design neural network layers that are equivariant to permutations of elements in the input  $\mathbf{x}$ . The function  $\mathbf{f} : \mathcal{X}^M \rightarrow \mathcal{Y}^M$  is **equivariant** to the permutation of its inputs iff

$$\mathbf{f}(\pi\mathbf{x}) = \pi\mathbf{f}(\mathbf{x}) \quad \forall \pi \in \mathcal{S}_M$$

where the symmetric group  $\mathcal{S}_M$  is the set of all permutation of indices  $1, \dots, M$ .

Consider the standard neural network layer

$$\mathbf{f}_\Theta(\mathbf{x}) \doteq \sigma(\Theta\mathbf{x}) \quad \Theta \in \mathbb{R}^{M \times M} \tag{20}$$

where  $\Theta$  is the weight vector and  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is a nonlinearity such as sigmoid function. The following lemma states the necessary and sufficient conditions for permutation-equivariance in this type of function.

**Lemma 3** *The function  $\mathbf{f}_\Theta : \mathbb{R}^M \rightarrow \mathbb{R}^M$  as defined in (20) is permutation equivariant if and only if all the off-diagonal elements of  $\Theta$  are tied together and all the diagonal elements are equal as well. That is,*

$$\Theta = \lambda\mathbf{I} + \gamma(\mathbf{1}\mathbf{1}^\top) \quad \lambda, \gamma \in \mathbb{R} \quad \mathbf{1} = [1, \dots, 1]^\top \in \mathbb{R}^M$$

where  $\mathbf{I} \in \mathbb{R}^{M \times M}$  is the identity matrix.

1. **Goal:** To show that the matrix  $\Theta = \lambda I + \gamma(11^T)$  commutes with all permutation matrices  $\pi \in S_M$ , ensuring permutation equivariance of the function  $f_\Theta(x) = \sigma(\Theta x)$ .

2. **Linearity of Commute:**

- The linear combination  $\Theta = \lambda I + \gamma(11^T)$  commutes with any permutation matrix  $\pi$  because both the identity matrix  $I$  and the constant matrix  $11^T$  commute with any permutation matrix.

3. **Diagonal Elements are Identical:**

- For any transposition  $\pi_{k,l}$ , it is shown that the diagonal elements  $\Theta_{i,i}$  must be identical for all  $i$ , implying  $\Theta_{i,i} = \lambda$ .

4. **Off-Diagonal Elements are Identical:**

- By considering pairs of off-diagonal elements, it is shown that the off-diagonal elements  $\Theta_{i,j}$  must be identical, concluding that  $\Theta_{i,j} = \gamma$  for all off-diagonal  $i \neq j$ .

5. **Conclusion:** The necessary and sufficient conditions for permutation equivariance are met, as  $\Theta = \lambda I + \gamma(11^T)$  commutes with all permutation matrices  $\pi \in S_M$ .

From definition of permutation equivariance  $\mathbf{f}_\Theta(\pi \mathbf{x}) = \pi \mathbf{f}_\Theta(\mathbf{x})$  and definition of  $\mathbf{f}$  in (20), the condition becomes  $\sigma(\Theta \pi \mathbf{x}) = \pi \sigma(\Theta \mathbf{x})$ , which (assuming sigmoid is a bijection) is equivalent to  $\Theta \pi = \pi \Theta$ . Therefore we need to show that the necessary and sufficient conditions for the matrix  $\Theta \in \mathbb{R}^{M \times M}$  to commute with all permutation matrices  $\pi \in \mathcal{S}_M$  is given by this proposition. We prove this in both directions:

- To see why  $\Theta = \lambda \mathbf{I} + \gamma (\mathbf{1}\mathbf{1}^\top)$  commutes with any permutation matrix, first note that commutativity is linear – that is

$$\Theta_1 \pi = \pi \Theta_1 \wedge \Theta_2 \pi = \pi \Theta_2 \Rightarrow (a\Theta_1 + b\Theta_2)\pi = \pi(a\Theta_1 + b\Theta_2).$$

Since both Identity matrix  $\mathbf{I}$ , and constant matrix  $\mathbf{1}\mathbf{1}^\top$ , commute with any permutation matrix, so does their linear combination  $\Theta = \lambda \mathbf{I} + \gamma (\mathbf{1}\mathbf{1}^\top)$ .

- We need to show that in a matrix  $\Theta$  that commutes with “all” permutation matrices

- *All diagonal elements are identical:* Let  $\pi_{k,l}$  for  $1 \leq k, l \leq M, k \neq l$ , be a transposition (i.e. a permutation that only swaps two elements). The inverse permutation matrix of  $\pi_{k,l}$  is the permutation matrix of  $\pi_{l,k} = \pi_{k,l}^\top$ . We see that commutativity of  $\Theta$  with the transposition  $\pi_{k,l}$  implies that  $\Theta_{k,k} = \Theta_{l,l}$ :

$$\pi_{k,l} \Theta = \Theta \pi_{k,l} \Rightarrow \pi_{k,l} \Theta \pi_{l,k} = \Theta \Rightarrow (\pi_{k,l} \Theta \pi_{l,k})_{l,l} = \Theta_{l,l} \Rightarrow \Theta_{k,k} = \Theta_{l,l}$$

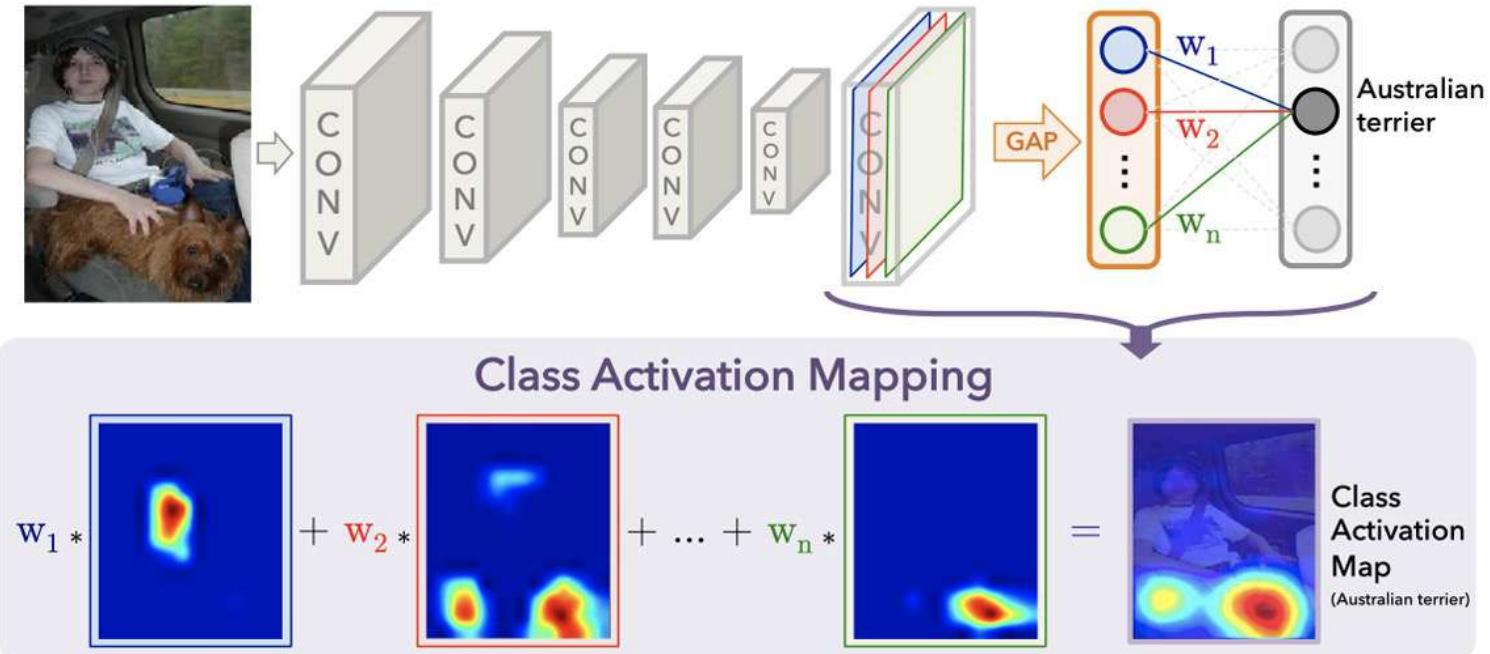
Therefore,  $\pi$  and  $\Theta$  commute for any permutation  $\pi$ , they also commute for any transposition  $\pi_{k,l}$  and therefore  $\Theta_{i,i} = \lambda \forall i$ .

- *All off-diagonal elements are identical:* We show that since  $\Theta$  commutes with any product of transpositions, any choice two off-diagonal elements should be identical. Let  $(i, j)$  and  $(i', j')$  be the index of two off-diagonal elements (i.e.  $i \neq j$  and  $i' \neq j'$ ). Moreover for now assume  $i \neq i'$  and  $j \neq j'$ . Application of the transposition  $\pi_{i,i'} \Theta$ , swaps the rows  $i, i'$  in  $\Theta$ . Similarly,  $\Theta \pi_{j,j'}$  switches the  $j^{\text{th}}$  column with  $j'^{\text{th}}$  column. From commutativity property of  $\Theta$  and  $\pi \in \mathcal{S}_n$  we have

$$\begin{aligned} \pi_{j',j} \pi_{i,i'} \Theta &= \Theta \pi_{j',j} \pi_{i,i'} \Rightarrow \pi_{j',j} \pi_{i,i'} \Theta (\pi_{j',j} \pi_{i,i'})^{-1} = \Theta \Rightarrow \\ \pi_{j',j} \pi_{i,i'} \Theta \pi_{i',i} \pi_{j,j'} &= \Theta \Rightarrow (\pi_{j',j} \pi_{i,i'} \Theta \pi_{i',i} \pi_{j,j'})_{i,j} = \Theta_{i,j} \Rightarrow \Theta_{i',j'} = \Theta_{i,j} \end{aligned}$$

where in the last step we used our assumptions that  $i \neq i'$ ,  $j \neq j'$ ,  $i \neq j$  and  $i' \neq j'$ . In the cases where either  $i = i'$  or  $j = j'$ , we can use the above to show that  $\Theta_{i,j} = \Theta_{i'',j''}$  and  $\Theta_{i',j'} = \Theta_{i'',j''}$ , for some  $i'' \neq i, i'$  and  $j'' \neq j, j'$ , and conclude  $\Theta_{i,j} = \Theta_{i',j'}$ .

# Where is a ANN looking at? Class activation mapping



How does a network make a decision? Fourier  
important features.

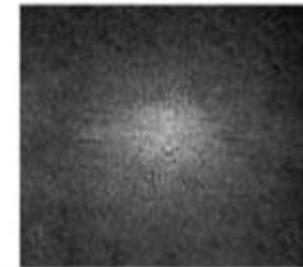
$x$

$M_\Phi \odot \mathcal{F}x$

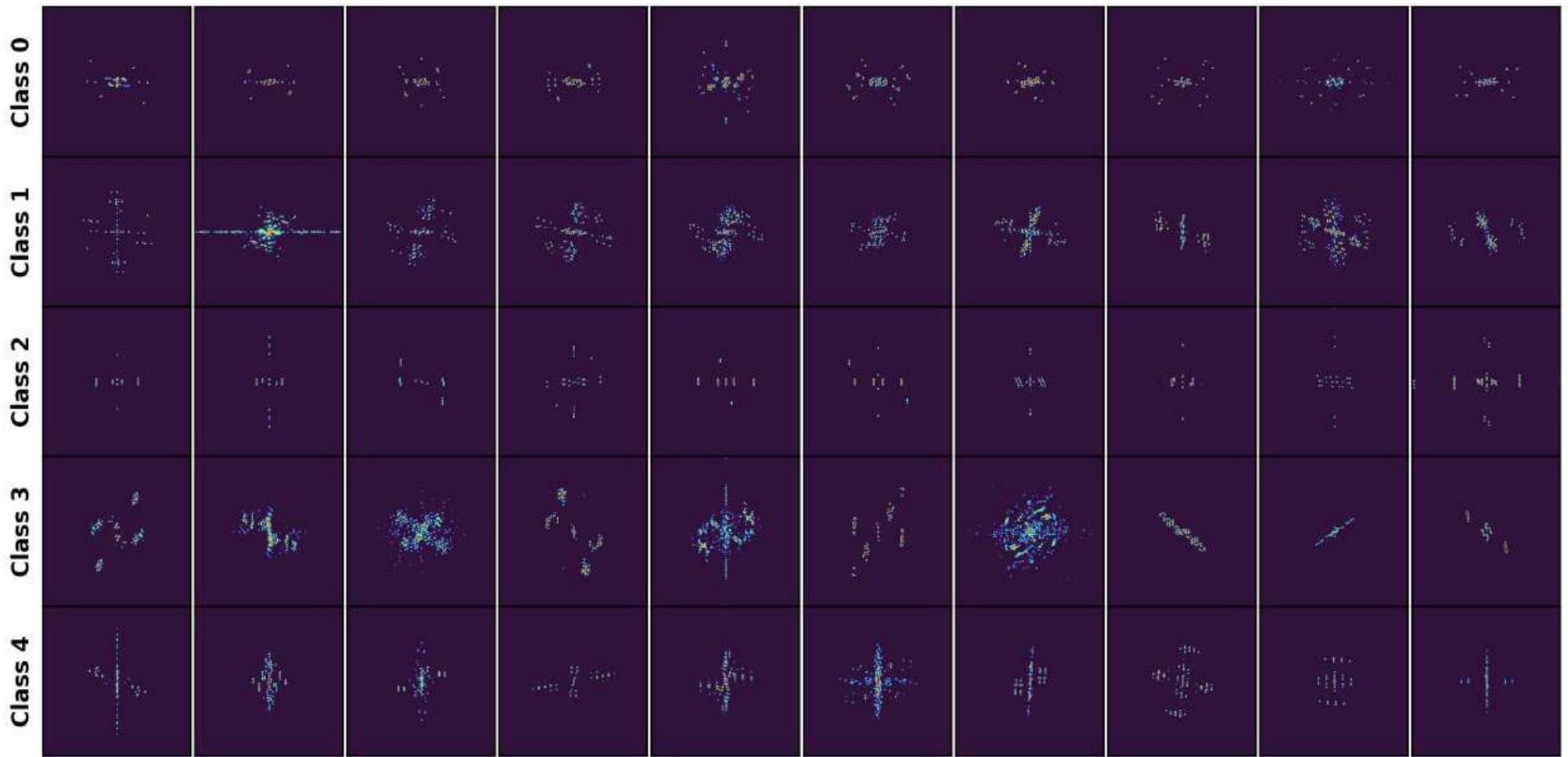
$\mathcal{F}^{-1}(M_\Phi \odot \mathcal{F}x)$



•



$$M_\Phi(\lambda, p) = \operatorname{argmin}_{M_\Phi} \sum_{x \in \mathcal{X}_V} e^{[\mathcal{L}(\Phi(\bar{x}), y) - \mathcal{L}(\Phi(x), y)]^2} + \lambda \|M_\Phi\|_p, \quad \lambda \in \mathbb{R}_+$$



# Or texture?

IMAGENET-TRAINED CNNS ARE BIASED TOWARDS  
TEXTURE; INCREASING SHAPE BIAS IMPROVES  
ACCURACY AND ROBUSTNESS

**Robert Geirhos**  
University of Tübingen & IMPRS-IS  
[robert.geirhos@bethgelab.org](mailto:robert.geirhos@bethgelab.org)

**Patricia Rubisch**  
University of Tübingen & U. of Edinburgh  
[p.rubisch@sms.ed.ac.uk](mailto:p.rubisch@sms.ed.ac.uk)

**Claudio Michaelis**  
University of Tübingen & IMPRS-IS  
[claudio.michaelis@bethgelab.org](mailto:claudio.michaelis@bethgelab.org)

**Matthias Bethge\***  
University of Tübingen  
[matthias.bethge@bethgelab.org](mailto:matthias.bethge@bethgelab.org)

**Felix A. Wichmann\***  
University of Tübingen  
[felix.wichmann@uni-tuebingen.de](mailto:felix.wichmann@uni-tuebingen.de)

**Wieland Brendel\***  
University of Tübingen  
[wieland.brendel@bethgelab.org](mailto:wieland.brendel@bethgelab.org)



(a) Texture image  
81.4% **Indian elephant**  
10.3% indri  
8.2% black swan

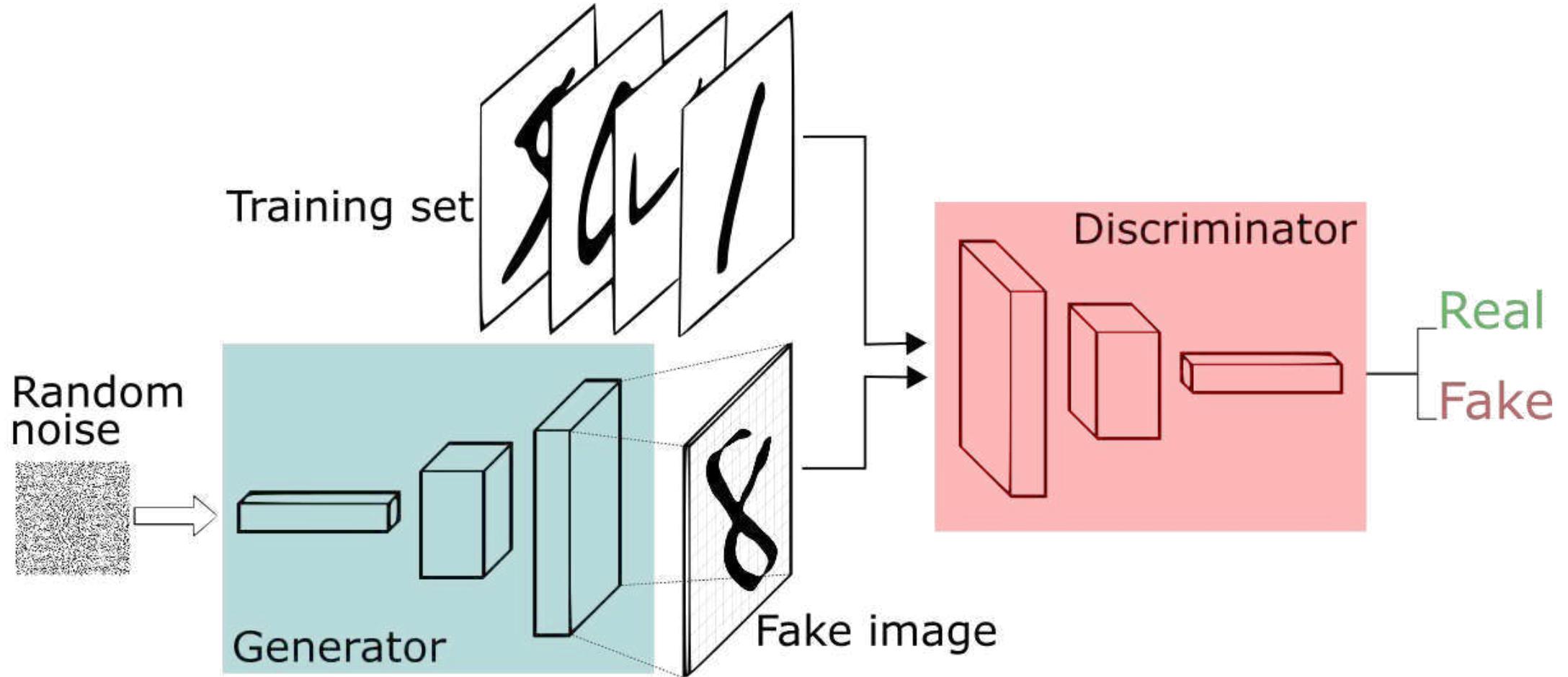


(b) Content image  
71.1% **tabby cat**  
17.3% grey fox  
3.3% Siamese cat



(c) Texture-shape cue conflict  
63.9% **Indian elephant**  
26.4% indri  
9.6% black swan

# Generative adversarial networks



# Gan Loss

gradient ascent    predict well on real images  
=> want probability close to 1

predict well on fake images  
=> want probability close to 0

$$\nabla_{\mathbf{W}_D} \frac{1}{n} \sum_{i=1}^n \left[ \overbrace{\log D(\mathbf{x}^{(i)})}^{\text{predict well on real images}} + \log \left( 1 - D(G(\mathbf{z}^{(i)})) \right) \right]$$

Discriminator objective in the neg. log-likelihood (binary cross entropy) perspective:

Real images,  $y = 1$

$$\mathcal{L}(\mathbf{w}) = \boxed{-y^{(i)} \log (\hat{y}^{(i)})} - (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})$$

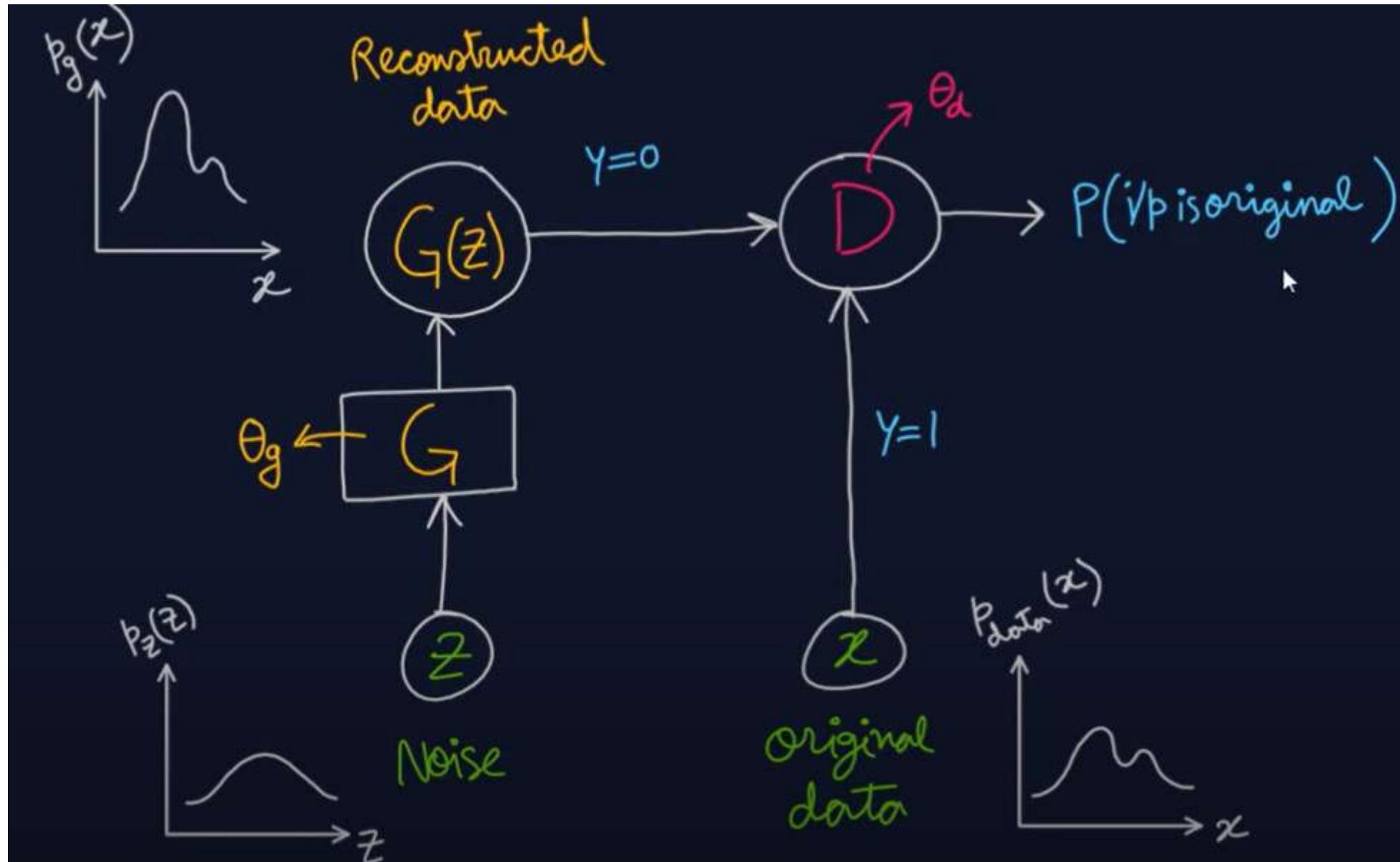
Want,  $\hat{y} = 1$

Fake images,  $y = 0$

$$\mathcal{L}(\mathbf{w}) = -y^{(i)} \log (\hat{y}^{(i)}) \boxed{- (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})}$$

Want,  $\hat{y} = 0$

## General setting



The problem: the discriminator try to mimic data distribution; the generator tries to fool the discriminator

Value Function:

$$\min_G \max_D V(G, D) = E_{x \sim p_{\text{data}}} [\ln(D(x))]$$

+

$$E_{z \sim p_z} [\ln(1 - D(G(z)))]$$

## Analogy with Binary Cross-Entropy

Binary Cross-Entropy Function

$$\mathcal{L} = -\sum \hat{y} \ln \hat{y} + (1-\hat{y}) \ln (1-\hat{y})$$

when  $y=1, \hat{y}=D(x) \Rightarrow \mathcal{L} = \ln [D(x)]$

when  $y=0, \hat{y}=D(G(z)) \Rightarrow \mathcal{L} = \ln [1-D(G(z))]$

Adding,  $\mathcal{L} = \ln [D(x)] + \ln [1-D(G(z))]$

## Expectation

$$E(\mathcal{L}) = E(\ln[D(x)]) + E(\ln[1 - D(G(z))])$$

$$\sum p_{\text{data}}(x) \ln[D(x)] + \sum p_z(z) \ln[1 - D(G(z))]$$

$$\int p_{\text{data}}(x) \ln[D(x)] dx + \int p_z(z) \ln[1 - D(G(z))] dz$$

$$V(G, D) = E_{x \sim p_{\text{data}}} [\ln(D(x))] + E_{z \sim p_z} [\ln(1 - D(G(z)))]$$

Training loop :

\* fix the learning of G \*

Inner loop for D :

- take m data samples & m fake data samples
- update  $\theta_d$  by grad. ascent

$$\frac{\partial}{\partial \theta_d} \frac{1}{m} \left[ \ln [D(x)] + \ln [1 - D(G(z))] \right]$$

\* fix the learning of D \*

take m fake data samples

update  $\theta_g$  by grad. descent

$$\frac{\partial}{\partial \theta_g} \frac{1}{m} \left[ \ln [1 - D(G(z))] \right]$$

What is the Optimization doing?

for fixed  $G$ ,

$$V(G, D) = \int_{\mathcal{X}} p_{\text{data}}(x) \ln[D(x)] + p_g(x) \ln[1-D(x)] dx$$

will be maximum for  $D(x) =$

$$\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}$$

\*

[https://www.youtube.com/watch?v=Gib\\_kiXgnvA](https://www.youtube.com/watch?v=Gib_kiXgnvA)

We have fixed  $D(x)$

$$\min_G V = E_{x \sim p_{\text{data}}} \ln \left( \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \right) + E_{x \sim p_g} \ln \left( 1 - \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \right)$$

$$= E_{x \sim p_{\text{data}}} \ln \left( \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \right) + E_{x \sim p_g} \ln \left( \frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)} \right)$$

$$= \mathbb{E}_{x \sim p_{\text{data}}} \ln \left( \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \right) + \mathbb{E}_{x \sim p_g} \ln \left( \frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)} \right)$$

JS divergence

$$\text{JS}(p_1 || p_2) = \frac{1}{2} \mathbb{E}_{x \sim p_1} \ln \left( \frac{p_1}{\frac{p_1 + p_2}{2}} \right) + \frac{1}{2} \mathbb{E}_{x \sim p_2} \left( \frac{p_2}{\frac{p_1 + p_2}{2}} \right)$$

$$\min_G V = E_{x \sim p_{\text{data}}} \ln \left( \frac{p_{\text{data}}(x)}{\frac{p_{\text{data}}(x) + p_g(x)}{2}} \right) + E_{x \sim p_g} \ln \left( \frac{p_g(x)}{\frac{p_{\text{data}}(x) + p_g(x)}{2}} \right) - 2 \ln 2$$

$$\min_G V = 2 \text{JS}(p_{\text{data}} || p_g) - 2 \ln 2$$

The optimization tries to match the training data  
and the generator data probability distributions

# Importance of ANNs in Neuroscience

2021 Special Issue on AI and Brain Science: Perspective

## Natural and Artificial Intelligence: A brief introduction to the interplay between AI and neuroscience research



Tom Macpherson <sup>a</sup>, Anne Churchland <sup>b</sup>, Terry Sejnowski <sup>c,d</sup>, James DiCarlo <sup>e</sup>,  
Yukiya Kamitani <sup>f,g</sup>, Hidehiko Takahashi <sup>h</sup>, Takatoshi Hikida <sup>a,\*</sup>

<sup>a</sup> Laboratory for Advanced Brain Functions, Institute for Protein Research, Osaka University, Osaka, Japan

<sup>b</sup> Cold Spring Harbor Laboratory, Neuroscience, Cold Spring Harbor, NY, USA

<sup>c</sup> Computational Neurobiology Laboratory, Salk Institute for Biological Studies, CA, USA

<sup>d</sup> Division of Biological Sciences, University of California San Diego, CA, USA

<sup>e</sup> Brain and Cognitive Sciences, Massachusetts Institute of Technology, MA, USA

<sup>f</sup> Department of Neuroinformatics, ATR Computational Neuroscience Laboratories, Kyoto, Japan

<sup>g</sup> Graduate School of Informatics, Kyoto University, Kyoto, Japan

<sup>h</sup> Department of Psychiatry and Behavioral Sciences, Tokyo Medical and Dental University Graduate School, Tokyo, Japan

---

### ARTICLE INFO

---

#### Article history:

Available online 28 September 2021

---

#### Keywords:

Artificial intelligence

Neuroscience

Neural imaging

Visual processing

Working memory

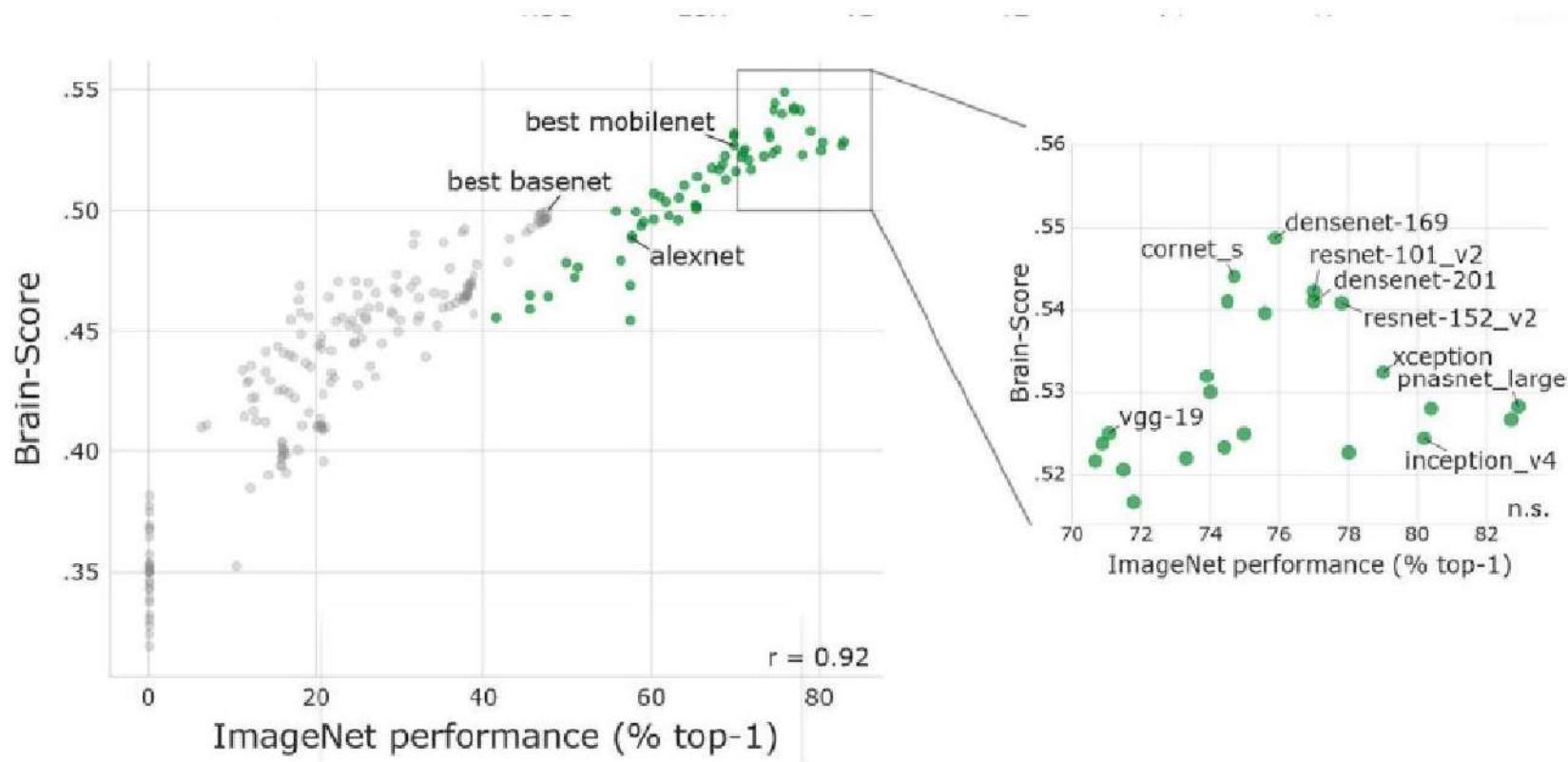
Computational psychiatry

---

### ABSTRACT

Neuroscience and artificial intelligence (AI) share a long history of collaboration. Advances in neuroscience, alongside huge leaps in computer processing power over the last few decades, have given rise to a new generation of *in silico* neural networks inspired by the architecture of the brain. These AI systems are now capable of many of the advanced perceptual and cognitive abilities of biological systems, including object recognition and decision making. Moreover, AI is now increasingly being employed as a tool for neuroscience research and is transforming our understanding of brain functions. In particular, deep learning has been used to model how convolutional layers and recurrent connections in the brain's cerebral cortex control important functions, including visual processing, memory, and motor control. Excitingly, the use of neuroscience-inspired AI also holds great promise for understanding how changes in brain networks result in psychopathologies, and could even be utilized in treatment regimes. Here we discuss recent advancements in four areas in which the relationship between neuroscience and AI has led to major advancements in the field; (1) AI models of working memory, (2) AI visual processing, (3) AI analysis of big neuroscience datasets, and (4) computational psychiatry.

# Which ANN?: Brain score



"Brain-Score: Which Artificial Neural Network for Object Recognition is most Brain-Like?" M. Schrimpf et al. 2020, <https://www.brain-score.org/>

Score based on: neural activity matching and behavioural tasks (recognition etc.).

## Invariance w.r.t. to what?: group transformations

Group: set  $\mathcal{G}$  with composition/multiplication operation satisfying:

- ① closure:  $gg' \in \mathcal{G}, \forall g, g' \in \mathcal{G}$
- ② associativity:  $(gg')g'' = g(g'g'') \in \mathcal{G}, \forall g, g', g'' \in \mathcal{G}$
- ③ identity: there exists  $Id \in \mathcal{G}$  such that  $Idg = gId = g \forall g \in \mathcal{G}$
- ④ invertibility:  $\forall g \in \mathcal{G}$  there exists  $g^{-1} \in \mathcal{G}$  :  $gg^{-1} = Id$ .

Rotation



Scale

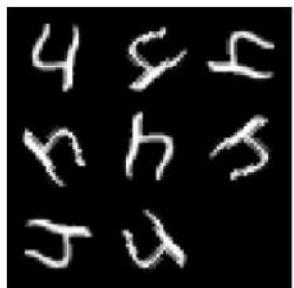


└

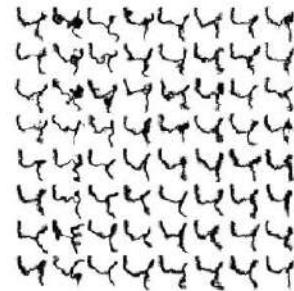
## Examples of Transformations/Groups

- Compact: rotation
- Locally compact: translation, scaling (multiplication), affine  
e.g., Translation group: linear operator  $T_\tau : \mathcal{X} \rightarrow \mathcal{X}$

$$T_\tau x(p) = x(p - \tau), \quad \forall p, \tau \in \mathbb{R}, x \in \mathcal{X}$$



0	1	2	3	4	5
6	7	8	9	-	τ
~	∞	∞	0	-	∞
4	5	1	7	∞	2
9	0	9	8	4	∞
-	-	3	2	1	0



- Non-locally compact: diffeomorphisms, local/global deformations  
e.g. given a smooth map  $d : \mathbb{R} \rightarrow \mathbb{R}$ : linear operator  $D_d : \mathcal{X} \rightarrow \mathcal{X}$

$$D_d x(p) = x(d(p)), \quad \forall p \in \mathbb{R}, x \in \mathcal{X}$$

- Non-group transformations
  - e.g. 3D rotations in depth, illumination.

—

## Defining invariance w.r.t. group transformations of stimuli: equivalence classes

Let  $\mathcal{G}$  a group. Let  $X$  set is a collection of orbits of stimuli:

$$O_t = \{t' \mid t' = gt, g \in \mathcal{G}, t \in \mathcal{X}\}, \quad X = (O_{t_1}, \dots, O_{t_M})$$



The group  $\mathcal{G}$  induces a **partition on the data**:

$$x' \sim x \Leftrightarrow x' = gx, \quad \exists g \in \mathcal{G}$$

## Invariant and Selective Representations

Let  $\Phi : \mathcal{X} \rightarrow \mathcal{F}$ .

Invariance:

$$x' \sim x \Rightarrow \Phi(x') = \Phi(x), \quad \forall x, x' \in \mathcal{X}$$

$$x' = gx \Rightarrow \Phi(x') = \Phi(x), \quad \forall x, x' \in \mathcal{X}$$

(trivial invariant representations are possible, e.g. constant function)

Selectivity:

$$x' \sim x \Leftarrow \Phi(x') = \Phi(x), \quad \forall x, x' \in \mathcal{X}$$

$$x = gx' \Leftarrow \Phi(x) = \Phi(x'), \quad \forall x, x' \in \mathcal{X}$$

## Invariant representations: Haar measure

Example (Translation group/Lebesgue measure)

$$d(T_{s'} s) = d(s - s') = ds$$

Translation invariance of the integral:

$$\int f(s)ds = \int f(s - s')ds$$

Key observation: a (locally) compact group has associated an **invariant Haar measure**  $dg$  i.e.

$$d(gg') = dg, \quad \forall g' \in \mathcal{G}$$

## Building invariants by group averaging (1)

Haar invariance allows to build invariants in the form of group averages:

$$I_f = \int dg f(g) = \int dg f(ge) \equiv I_f(e)$$

In fact

$$\begin{aligned} I_f(\bar{g}) &= \int dg f(g\bar{g}) \\ &= \int d(\hat{g}\bar{g}^{-1}) f(\hat{g}) \quad \hat{g} = g\bar{g}, \quad g = \hat{g}\bar{g}^{-1} \\ &= \int d\hat{g} f(\hat{g}), \quad d(\hat{g}\bar{g}^{-1}) = d\hat{g} \\ &= \int d\hat{g} f(\hat{g}) = I_f(e) \end{aligned}$$

where we used: 1) reparameterization of group elements 2) Haar invariance 3) group closure under composition.

## Choosing a specific representation

We choose  $f$  as:

$$f_{x,t}(g) = \eta \langle gx, t \rangle$$

where  $\eta : \mathbb{R} \rightarrow \mathbb{R}$  and  $t \in \mathcal{X}$ .

This choice is good for:

- Biological plausibility (we will see later)
- Allows to build an invariant representation of  $x$

## Invariant representations by group averaging (2)

### Theorem

Invariance  $\mu_t : \mathcal{X} \rightarrow \mathbb{R}$  defined as

$$\mu_t(x) = \int dg \eta \langle gx, t \rangle, \quad x, t \in \mathcal{X}$$

is invariant w.r.t. transformations  $x \rightarrow \bar{g}x, \forall \bar{g} \in \mathcal{G}$

Indeed

$$\begin{aligned}\mu_t(\bar{g}x) &= \int dg \eta \langle g\bar{g}x, t \rangle \\ &= \int dg \eta \langle \hat{g}x, t \rangle \quad \hat{g} = g\bar{g}, \quad g = \hat{g}\bar{g}^{-1} \\ &= \int d\hat{g} \eta \langle \hat{g}x, t \rangle, \quad d(\hat{g}\bar{g}^{-1}) = d\hat{g} \\ &= \mu_t(x)\end{aligned}$$

## Enough projections ensure selectivity

If we have enough  $t$ -projections we can discriminate signals s.t.  $x \not\sim x'$

Define a metric given  $K$  templates/projections:

$$D_K(x, x') = \frac{1}{K} \sum_{k=1}^K |\mu_{t_k}(x) - \mu_{t_k}(x')|.$$

A simple concentration inequality shows that:

$$|D_K(x, x') - D_\infty(x, x')| \leq \epsilon, \text{ with probability } 1 - \delta^2$$

if  $K > (\frac{2}{c}\epsilon^{-2}) \ln(\frac{Q}{\delta})$  with  $Q$  the number of equivalence classes.

This can be interpreted as a random projections/Johnson-Lindenstrauss type theorem but on 1) On quotient space w.r.t.  $\mathcal{G}$ . 2) With non-linear projections.

$$\mu_t(x) = \int dg \ \eta \langle gx, t \rangle$$

Do we need the full orbit of  $x$ ?

## Transferability and sample complexity

$\mu_t$  uses full orbit of  $x$  to build an invariant representation. However:

$$\langle gx, t \rangle = \langle x, g^{-1}t \rangle, \forall x, t \in \mathcal{X}$$

Then we have mathematical equivalence of

- $\mu_t$  calculated having  $O_x$

- $\mu_t$  calculated having  $O_t$

i.e.

$$\mu_t(x) = \int dg \eta \langle gx, t \rangle = \int dg \eta \langle x, g^{-1}t \rangle = \int dg \eta \langle x, gt \rangle$$

We need only one example of  $x$  and a set of transformed templates (memorized or learned) to have an invariant representation of  $x$ .

## Main result

### Theorem (Invariance and Selectivity for Finite Templates)

For  $Q$  orbits of  $\mathcal{G}$  in  $\mathcal{X}$ , with  $K > (\frac{2}{c}\epsilon^{-2}) \ln(\frac{Q}{\delta})$  templates,  $\epsilon, \delta > 0$ , the representation

$$\Phi(x) = (\mu_1(x), \dots, \mu_K(x))$$

is

- **Invariant** i.e.  $\Phi(gx) = \Phi(x)$ ,  $\forall g \in G$
- **(Almost) selective** i.e.  $|D_K(x, x') - D_\infty(x, x')| \leq \epsilon$  with probability  $1 - \delta^2$ .

In other words:

$$\Phi(x) = \Phi(x') \iff x \sim x'$$

## Algorithmic remarks

$$\mu_t(x) = \int dg\eta \langle x, gt \rangle$$

- We consider a finite group:  $\mathcal{G} = \{g_1, \dots, g_N\}$ .
- Activation function: step functions, sigmoids, max, moments, with different thresholds e.g.

$$\eta_b \langle gx, t \rangle = \sigma(\langle x, gt \rangle - b)$$

- Templates are generic (not reflecting data distribution) and finite  
 $\mathcal{T} = \{t^k\}_{k=1}^K$ .

$$\mu_b^k(x) = \sum_i \sigma(\langle g_i x, t_k \rangle - b)$$

## Algorithmic steps

Template sampling  $t \in \mathcal{T} = \{t^k\}_{k=1}^K$

Projections on transforming templates  $\{\langle x, g_j t^k \rangle\}, j = 1, \dots, |G|$

Pooling (sum) using  $B$  non-linear functions (e.g., sigmoids)

$$\mu_b^k(x) = \sum_{j=1}^{|G|} \eta_b \langle x, g_j t^k \rangle, t^k \in \mathcal{T}, g_j \in G, b = 1 \dots B$$

Signature from components

$$\Phi(x) = (\{\mu_1^1(x)\}, \{\mu_2^1(x)\}, \dots, \{\mu_N^K(x)\}) \in \mathbb{R}^{N \times |\mathcal{T}|}$$

Number of templates (to separate  $|\hat{Y}|$  classes with probability  $1 - \delta^2$ )

$$K \geq \frac{2}{c\epsilon^2} \log \frac{|\hat{Y}|}{\delta}$$

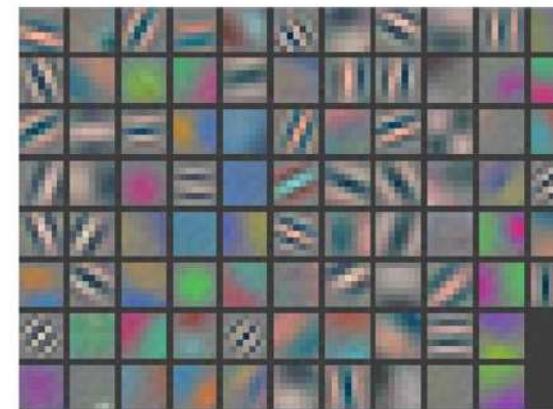
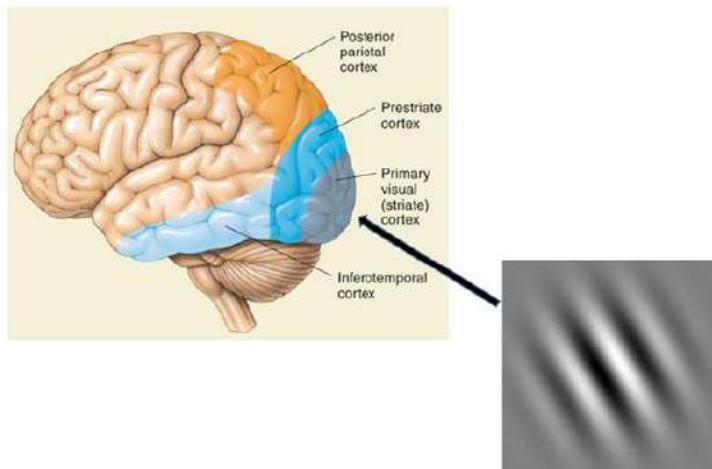
## What if the group is locally compact? approximate invariance and Gabor templates

Suppose we want to implement translation and scale invariance

$$\mu_t(x) = \sum_{\lambda,y} \sigma \langle x, D_\lambda T_y t \rangle$$

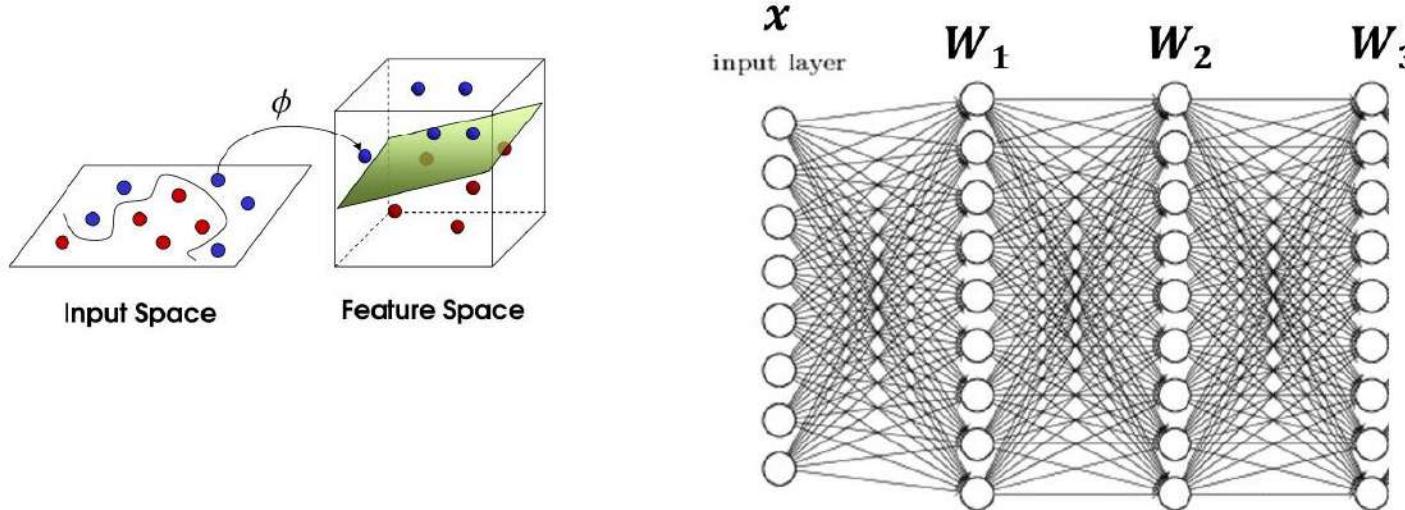
Which is the best choice for  $t = t^*$  to guarantee approximate invariance?

$$t^* = \arg \min_{t, \forall \lambda, y} |\mu_t(D_\lambda T_y x) - \mu_t(x)|$$



## How is this translated in the context of deep networks?

Deep networks: ‘smart’ way to parameterize functions.



$$x \in \mathbb{R}^d, \quad \Phi : \mathbb{R}^d \rightarrow \mathbb{R}^L, \quad \Phi \equiv \Phi_L \circ \dots \circ \Phi_1$$

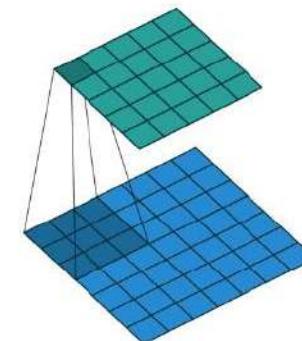
$$\Phi_1(x) = \sigma \langle W_1, x \rangle$$

$$\Phi_\ell(x) = \sigma \langle W_\ell, \Phi_{\ell-1}(x) \rangle$$

$$\sigma : \mathbb{R} \rightarrow \mathbb{R}, \quad \text{nowadays } \sigma(\cdot) = \max(0, \cdot) \equiv |\cdot|_+ \quad (w \equiv t).$$

## Constraints on network operations/topology

$$W = \begin{bmatrix} w_1 & w_d & w_1 & \dots & w_2 \\ w_2 & w_1 & w_2 & \dots & w_3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ w_d & w_{d-1} & w_{d-3} & \dots & w_1 \end{bmatrix} =$$

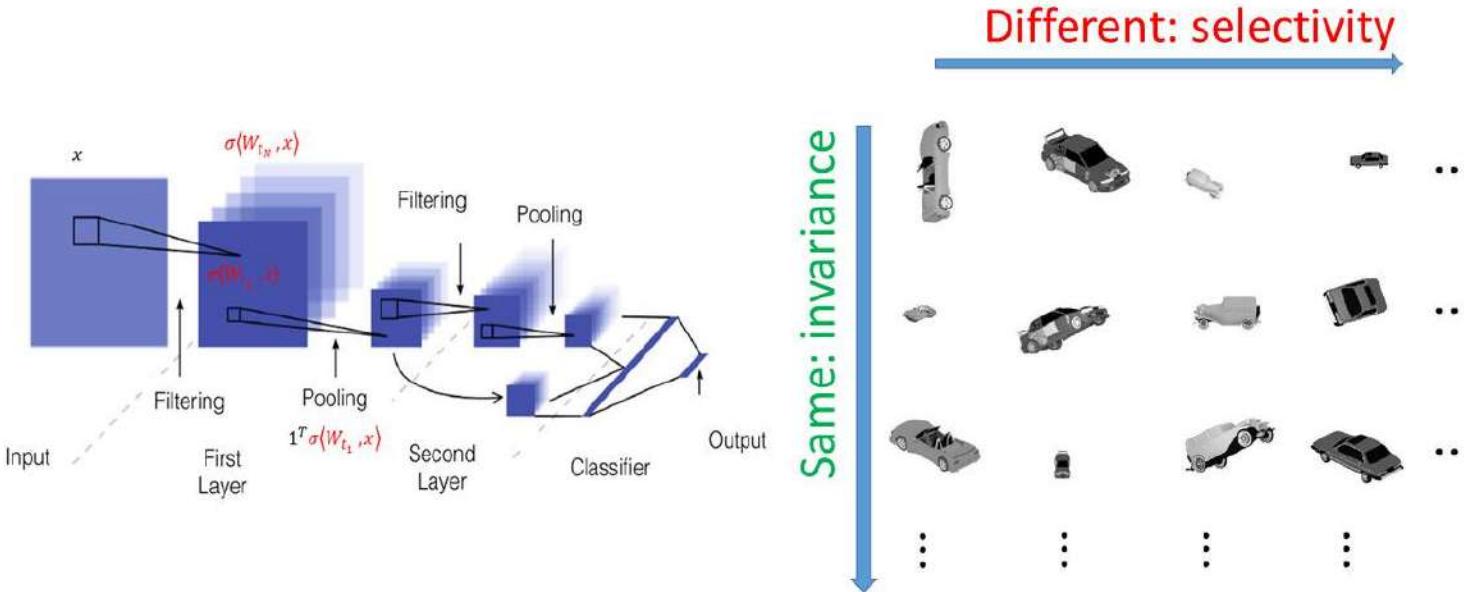


an orbit w.r.t. the cyclic group  $\mathcal{G}$ . This is equivalent to a convolution  
 $x * w = W^T x$ .

### Remark

*This construction can be extended to convolutions w.r.t. other groups.*

# Invariant selective Data Representation



## Theorem

Suppose data are generated by  $\mathcal{G}$ . A CNN with convolutions w.r.t.  $\mathcal{G}$  is implementing a data representation  $\Phi$  that is invariant and selective i.e.  $\mathbf{x} \sim \mathbf{x}' \Leftrightarrow \Phi(\mathbf{x}) = \Phi(\mathbf{x}')$  with (one layer)  $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$  and:

$$\Phi_w(x) = \sum_i |(x *_{\mathcal{G}} w)_i|_+ = \sum_i |(W^T x)_i|_+ = \sum_i |\langle x, g_i w \rangle|_+$$

## Orbits and probability distribution: some definitions

Consider the map  $x \rightarrow gx$  where  $g$  is the random variable. This defines a probability distribution  $x \rightarrow P_x(g)$  the probability distribution supported on each orbit

$$O_x = \{x' \mid x' = gx, g \in \mathcal{G}\}.$$

How do we compare different distributions  $P_x, P_{x'}$ ? We use arbitrary (continuous) test functions  $f \in C_c(\mathcal{X})$  and the use the following definition

$$I(f, P_x) \equiv \int_{O_x} f(y) dP_x(y) = \int_{\mathcal{X}} f(y) dP_x(y) = \int_{\mathcal{G}} f(gx) dg, \quad \forall f \in C_c(\mathcal{X})$$

with  $dg$  the Haar measure on the group. Then  $P_x = P_{x'}$  iff  $I(f, P_x) = I(f, P_{x'})$  for all  $f$ .

## Demonstration steps

We sketch the demonstration in 3 steps:

$$① \quad x \sim x' \iff O_x = O_{x'}$$

$$② \quad O_x = O_{x'} \iff P_x = P_{x'}$$

$$③ \quad P_x = P_{x'} \iff \Phi(x) = \Phi(x')$$

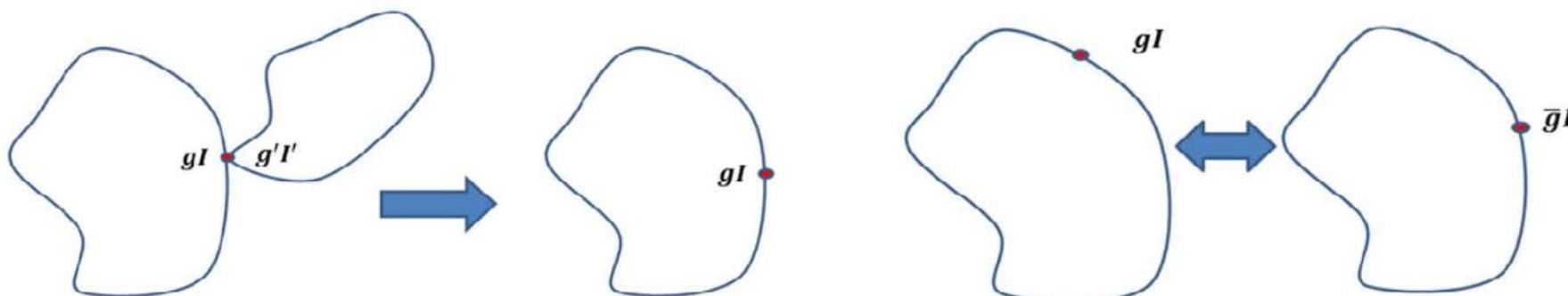
where  $\Phi$  is the CNN representation.

## Equivalence classes and orbits

$$x \sim x' \iff O_x = O_{x'}$$

①  $x \sim x' \Rightarrow x = gx' \Rightarrow O_x = O_{gx'} = O_{x'}$

②  $O_x = O_{x'} \Rightarrow x = gx' \Rightarrow x \sim x'$



## Equivalence Orbits-probability distribution

### Theorem

$$O_x = O_{x'} \Leftrightarrow P_x = P_{x'}$$

$\Rightarrow$ : If  $O_x = O_{x'} \Rightarrow x = hx', \exists h \in \mathcal{G}$

Then using the definition

$$\begin{aligned} \int_{O_x} f(y) dP_x(y) &= \int_{\mathcal{G}} f(gx) dg = \int_{\mathcal{G}} f(ghx') dg \\ &= \int_{\mathcal{G}} f(gx') dg = \int_{O_{x'}} f(y) dP_{x'}(y) \end{aligned}$$

for all  $f$  and thus  $P_x = P_{x'}$ .

$\Leftarrow$ : If  $P_x = P_{x'} \Rightarrow \text{supp}(P_x) = \text{supp}(P_{x'}) \Rightarrow O_x = O_{x'}$ .

## Mean probability embedding (3)

Theorem (Probability mean embedding, Smola et al. '07)

Let  $x \rightarrow P_x(g)$  induced by  $\mathcal{G}$  and  $\Psi(x)$  be a “rich enough” representation <sup>a</sup>. Then

$$P_x = P_{x'} \iff \Phi(x) = \int dg \Psi(gx) = \int dg \Psi(gx') = \Phi(x')$$

i.e. the distribution is one to one with the representation map.

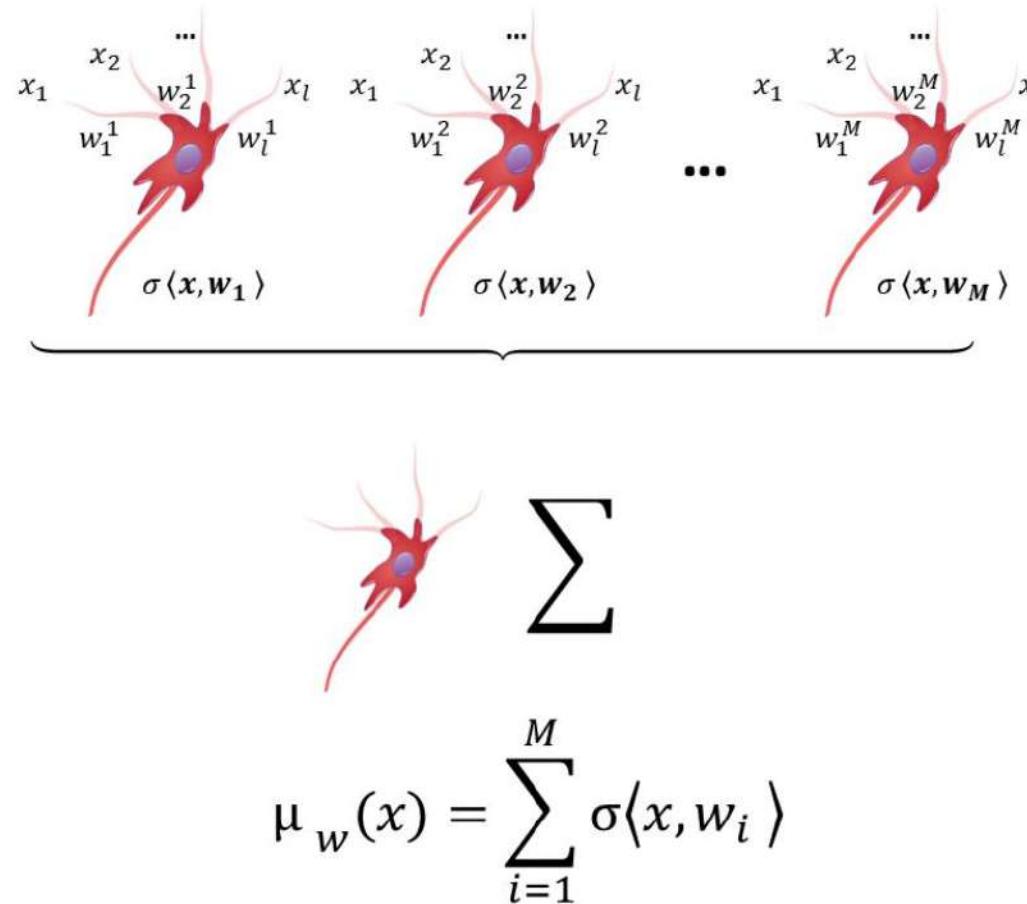
In the case of CNNs  $(\Psi(g_i x))_{w,i} = |\langle x, g_i w \rangle|_+ = |(x * w)_i|_+$  and

$$(\Phi(x))_w = \sum_i |\langle W_t, x \rangle_i|_+ = \sum_i |\langle g_i w, x \rangle|_+$$

---

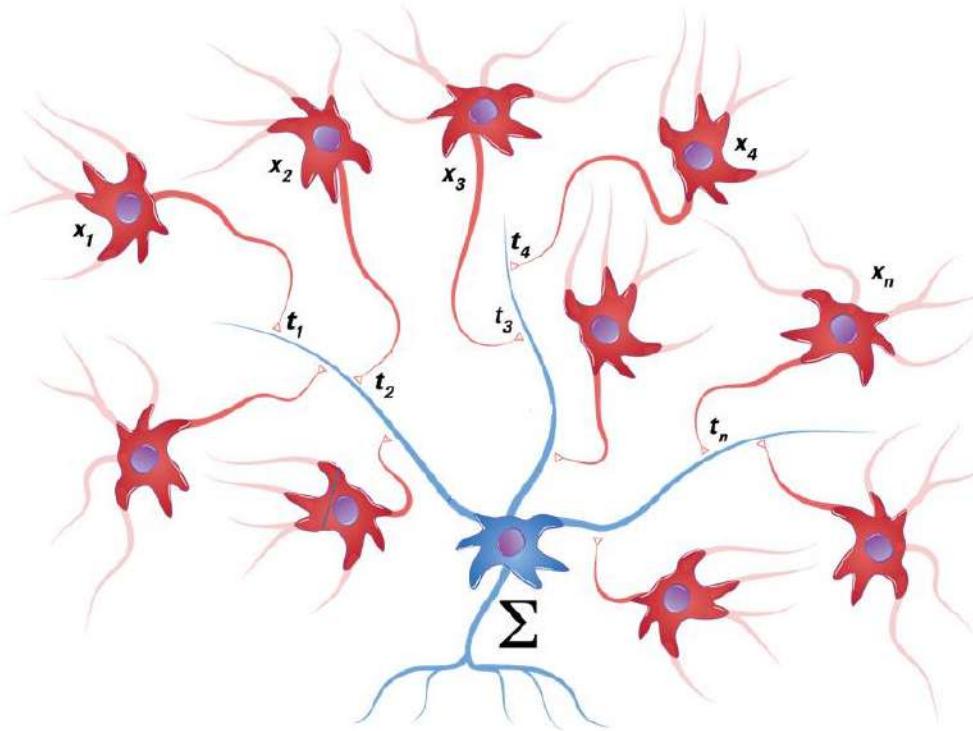
<sup>a</sup>The corresponding kernel universal, i.e. dense in the space of continuous functions

## Aggregating single neurons responses: Hubel Wiesel model



This module is often called in visual cortex *Simple-Complex* and  $\mu(x)$  is a neural representation of the input  $x$ .

## Translating our results



A cascade of simple and complex cells is implementing an invariant and selective representation.

## Robustness to more general transformations and locality



Let  $s$  a  $C^\infty$  transformation depending on  $\Theta = (\theta_1, \dots, \theta_P)$  parameters.  
Expanding, e.g., around e.g. 0:

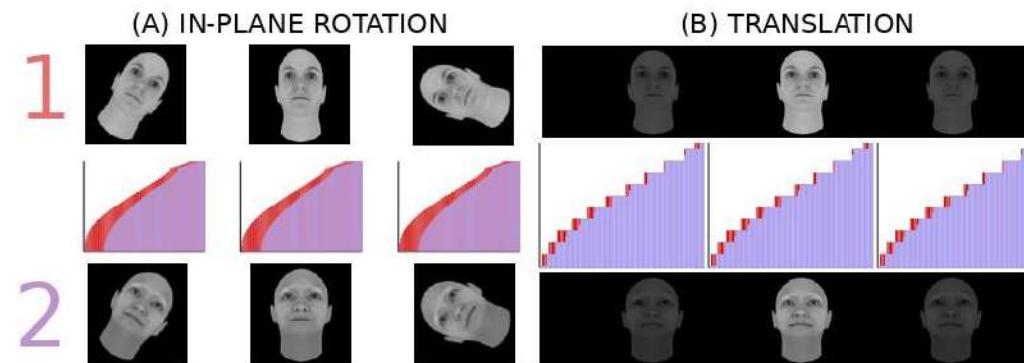
$$s(x, \Theta) = s(x, \mathbf{0}) + \sum_{i=1}^P \frac{\partial s(x, \Theta)}{\partial \theta_i} \theta_i + o(\|\Theta\|^2) = x + \sum_{i=1}^P \theta_i L_{\theta_i}(x) + o(\|\Theta\|^2)$$

where  $L_{\theta_i}$  are the infinitesimal generators and Therefore locally

$$g(\Theta) = \exp(\theta_1 L_{\theta_1} + \theta_2 L_{\theta_2} + \dots + \theta_P L_{\theta_P}).$$

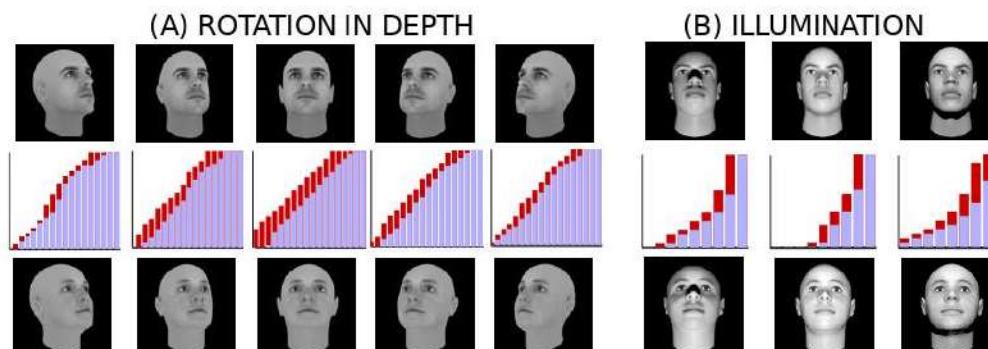
The local learned weights will be **orbits w.r.t. the local group approximating the non-group global transformation.**

## Beyond groups



$$\overline{\Phi}^t(x) = \sum_i |\langle g_i t, x \rangle|_+$$

As a processing procedure the principles extend beyond compact group transformations (Leibo et al. '15)



## Selectivity: probabilistic argument

(For compact groups) a probability distribution  $P_x$  on  $O_x$  is defined by measure  $dg$  via random variable  $g \mapsto gx$

### Theorem (Orbit-probability distribution equivalence)

$$(x \sim x' \iff O_x = O'_x \iff P_x \sim P_{x'})$$

“1D” probability distributions  $P_x^t$  are defined by projections  $gx \mapsto \langle gx, t \rangle$

### Theorem (Cramer-Wold)

*Probability distributions are uniquely determined by all 1D projections*

$$P_x \equiv P_{x'} \iff P_x^t \equiv P_{x'}^t, \quad \forall t \in \mathbb{S}^d$$

1D projection can be characterized by its CDF  $P_x^t \mapsto \int dg H(b - \langle x, gt \rangle)$

## Learning - Transferability

Unitary groups:  $\langle gx, t \rangle = \langle x, g^{-1}t \rangle, \forall x, t \in \mathcal{X}$

$$\mu(x) = \int \eta(\langle gx, t \rangle) dg = \int \eta(\langle x, g^{-1}t \rangle) = \int \eta(\langle x, gt \rangle) dg$$

### Memory-based learning:

- View based/augmentation: Need  $\{gx\}, \forall g \in \mathcal{G}$  ( $\forall x$ )
- Memory-based: Store  $\{gt\}, \forall g \in \mathcal{G}$  (once)
- Trading time with space requirements
- *Biological learning, learning from sequences, CNNs*

## Selective and invariant signature (recap)

Signature:

$$\mu_t(x) = \int dg\eta \langle x, gt \rangle$$

Or its computational-friendly form

$$\mu_b^k(x) = \sum_i \sigma(\langle x, g_i t_k \rangle - b)$$

The signature is invariant and selective.

## From previous class: stimuli hypothesis

Let  $\mathcal{G}$  a group. Let  $X$  set is a collection of orbits of stimuli:

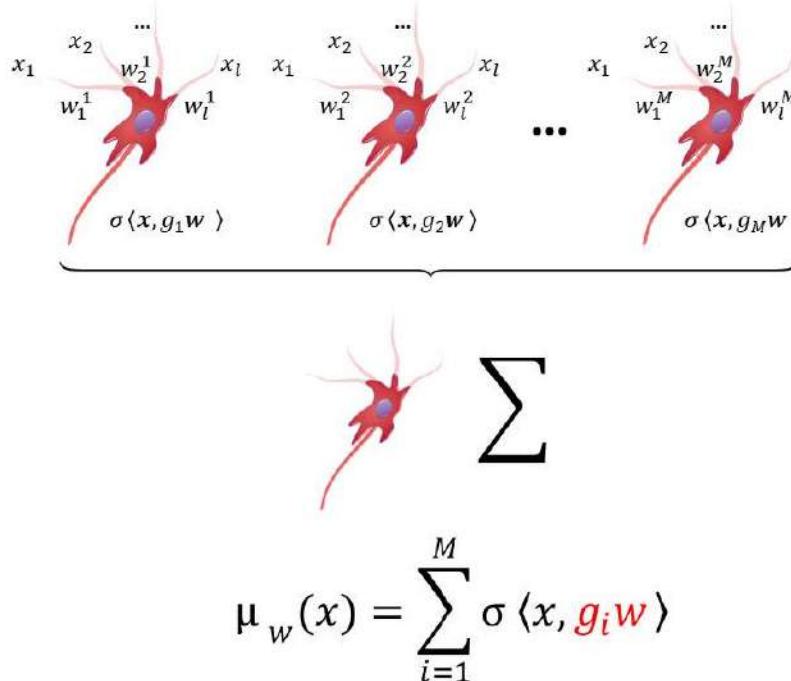
$$O_t = \{t' \mid t' = gt, g \in \mathcal{G}, t \in \mathcal{X}\}, \quad X = (O_{t_1}, \dots, O_{t_M})$$



The group  $\mathcal{G}$  induces a **partition on the data**:

$$x' \sim x \Leftrightarrow x' = gx, \quad \exists g \in \mathcal{G}$$

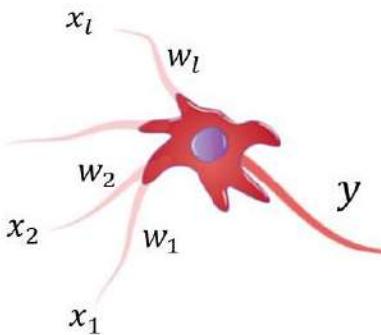
## Learning question



How can we learn  $g_i w$ ? How can we learn the correct pooling over orbits?

## Learning hypothesis: Hebbian learning of simple neurons synaptic weights

Neurons updates their weights during stimuli presentation. We consider a weights-normalized version of the hebbian rule  $\Delta\mathbf{w} = \eta y_t \mathbf{x}$

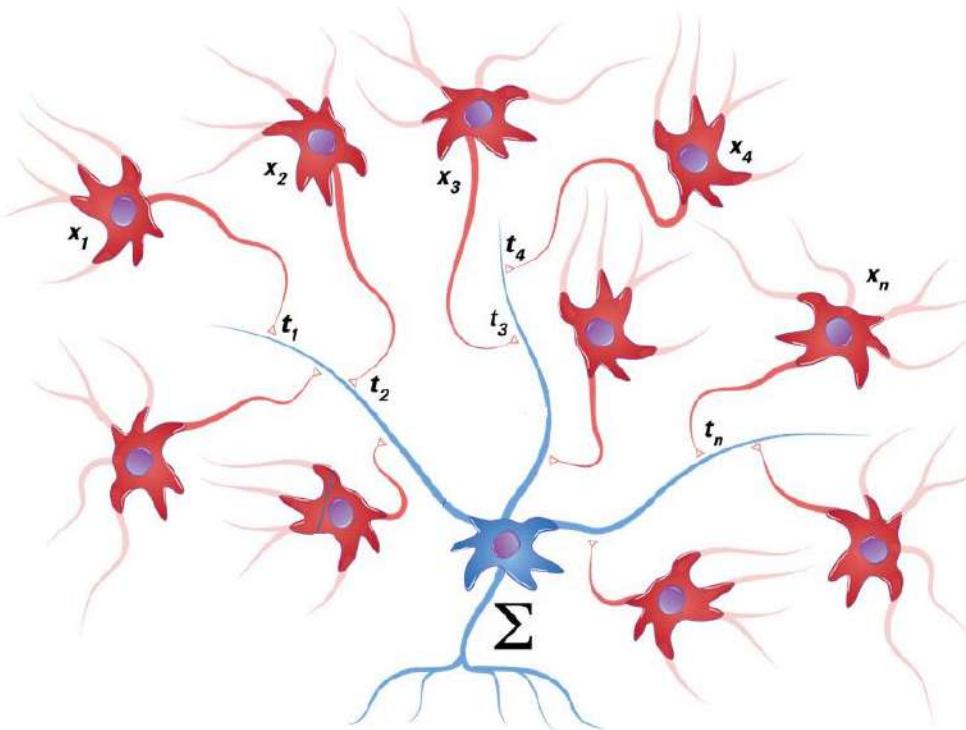


$$y_t = \langle \mathbf{x}_t, \mathbf{w}_t \rangle, \quad \Delta\mathbf{w} = \eta y_t (\mathbf{x}_t - y_t \mathbf{x}_t)$$

$$XX^T \mathbf{w} = \lambda_{max} \mathbf{w}$$

Simple cells weights converge to the first eigenvector of the covariance of the stimuli.

## Question



Can a simple-complex module,  $\mu_w(x)$ , under the hypotheses above, learn invariant/selective properties w.r.t.  $\mathcal{G}$ ?

## Stimuli + learning: simple neurons weights

### Theorem

*The weights of simple cells converge to linear combinations of elements of an orbit of  $\mathcal{G}$*

Simple proof:

- $X = \{g_1 t_1, \dots, g_{|\mathcal{G}|} t_1, \dots, g_1 t_Q, \dots, g_{|\mathcal{G}|} t_Q\}$

$$C = X X^T = \sum_i g_i T T^T g_i^T, \quad [C, g] = Cg - gC = 0$$

and

$$Cw = \lambda w \Rightarrow Cgw = gCw = \lambda gw, \quad \forall g \in \mathcal{G}, w \in E_\lambda$$

- In particular

$$E_{\max} = \text{span}(\mathbf{O}_w), \quad \forall w \in E_{\max}.$$

## Extension: equivariant plasticity rules

The result is true for a much broader class of plasticity rules those deriving from a Loss function of the form:

$$\mathcal{L}(X, W) = \sum_{i,j} f(\sigma \langle w_i, x_j \rangle), \quad f \in \mathcal{C}^1.$$

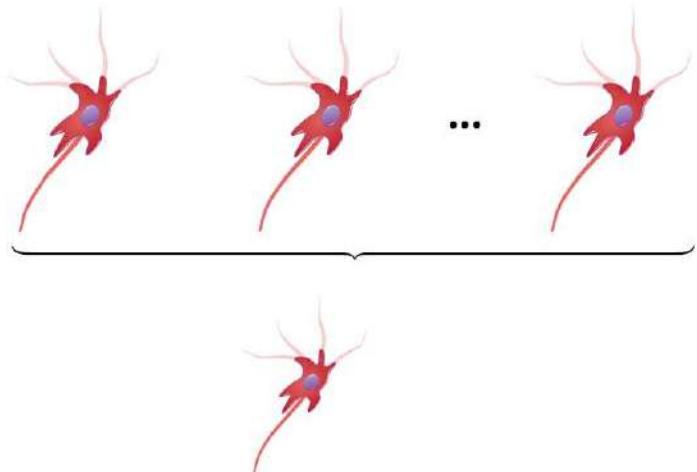
We have that the associated plasticity rule is equivariant w.r.t.  $\mathcal{G}$  transformations:

$$\nabla_{w_i} \mathcal{L}(\{w_1, \dots, \textcolor{red}{g}w_i, \dots, w_N\}, X) = \textcolor{red}{g} \nabla_{w_i} \mathcal{L}(\{w_1, \dots, w_i, \dots, w_N\}, X).$$

If  $\mathbf{w}$  is solution so is  $g_i \mathbf{w}$ .

## Which simple cells is a complex cell aggregating?

Suppose simple weights are learned.



$$O_x = \{x, g_2x, \dots, g_Mx\}$$

$$\bar{w} = \arg \max_w \langle w, x \rangle$$

$$g_2\bar{w} = \arg \max_w \langle w, g_2x \rangle$$

$$\langle x, \bar{w} \rangle = \langle gx, g\bar{w} \rangle$$

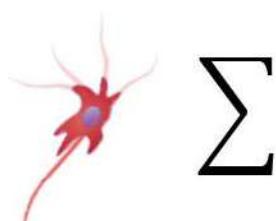
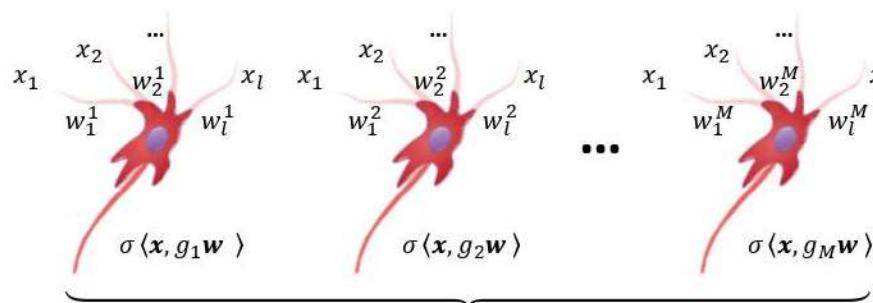
The principle is Hebbian: "fire together link together".

**Theorem (Complex cells pooling and invariance)**

A complex cell learns to aggregates over simple cells whose weights form an orbit with respect to the group  $\mathcal{G}$  i.e.  $w^i = g_i w$ .

## Summarizing:

- ① The symmetries in the stimuli are "inherited" by the weights of the simple cells:  $E_{\max} = \text{span}(\mathbf{O}_w)$ ,  $\forall w \in E_{\max}$ .
- ② The set of simple cells aggregated by a complex cell have weights that are orbits i.e.  $w^i = g_i w$ .



$$\sum$$

$$\mu_w(x) = \sum_{i=1}^M \sigma \langle x, g_i w \rangle$$

## Simple-complex modulus equivariant-invariant properties

- Simple cells are permutation-equivariant to  $g \in \mathcal{G}$  transformations.

$$\sigma(W^T g x) = P_g \sigma(W^T x), \quad W = (g_1 w, \dots, g_{|\mathcal{G}|} w)$$

- Complex cells are invariant to  $g \in \mathcal{G}$  transformations.

$$\mu_{\mathbf{w}}(\mathbf{x}) = \sum_i \sigma \langle x, g_i w \rangle = \mu_{\mathbf{w}}(\mathbf{gx}) \quad \forall g \in \mathcal{G}$$

Same properties hold for a wild class of aggregating functions (e.g. max) and pointwise non-linearities.

## What about selectivity/discriminability?

- Invariance:

$$x' \sim_{\mathcal{G}} x \Rightarrow \mu(x') = \mu(x), \quad \forall x, x' \in X$$

Elements within an orbit have the same  $\mu$ .

- Selectivity:

$$x' \sim_{\mathcal{G}} x \Leftarrow \mu(x') = \mu(x), \quad \forall x, x' \in X$$

Elements belonging to different orbits have different  $\mu$ .

We want to collapse elements within orbits and tear apart different orbits,  
i.e. build a proper quotient space.

## Discriminability: importance of the non-linearity

The information loss can be recovered considering a family of non-linearities e.g.  $\{\sigma_z(\cdot) \equiv H(\cdot - z), z \in \mathbb{R}\}$ .

$$c_z(x) \equiv (c(x))_z = \sum_{i=1}^{|\mathcal{G}|} H(\langle x, g_i w \rangle - z), \quad z \in \mathbb{R}.$$

We have:

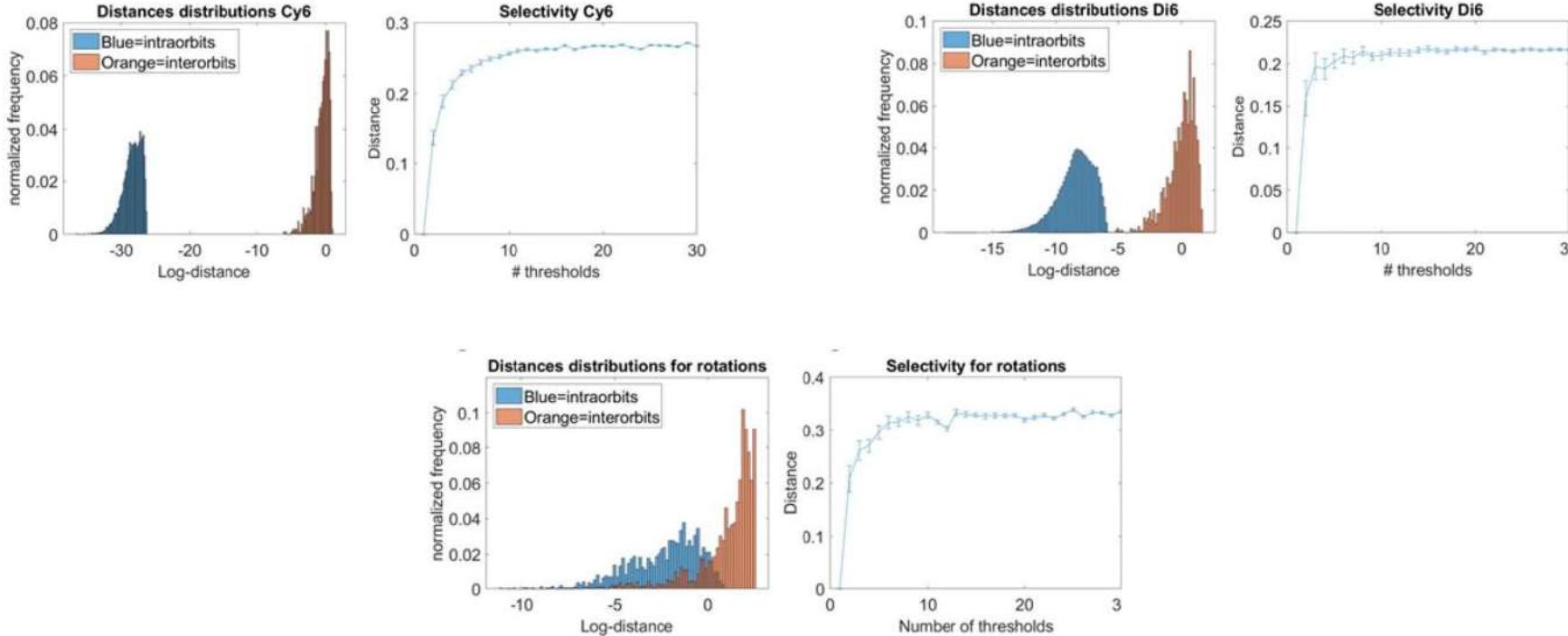
### Theorem (Complex cells selectivity)

Let  $x, x' \in \mathbb{R}^d$  be two input and  $c(x), c(x')$  their complex cells responses. Then the distance defined as

$$\text{dist}(x, x') := \|c(x) - c(x')\|_{\ell_2}.$$

is zero iff  $x \sim x'$ .

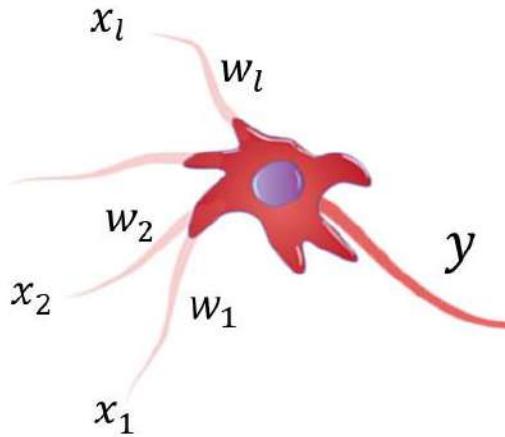
## Some simulations



- $X$  generated by: 1) Cyclic group, 2) Dihedral group, 3) Finite group of rotations on natural image patches.
- Plots: Intra and inter orbit distances for  $N$  orbits. Discriminability vs number of thresholds.

## Oja's learning on orbits

Oja's learning:



$$\begin{aligned}w_i(t+1) &= w_i(t) + \lambda y(x)x_i \\||w|| &= 1\end{aligned}$$

$w \rightarrow \text{first PCA}(X)$

The covariance matrix of  $D = \{(f_1, Rf_1), \dots, (f_N, Rf_N)\}$  is

$$C = D^T D = FF^T + RFF^T R^T$$

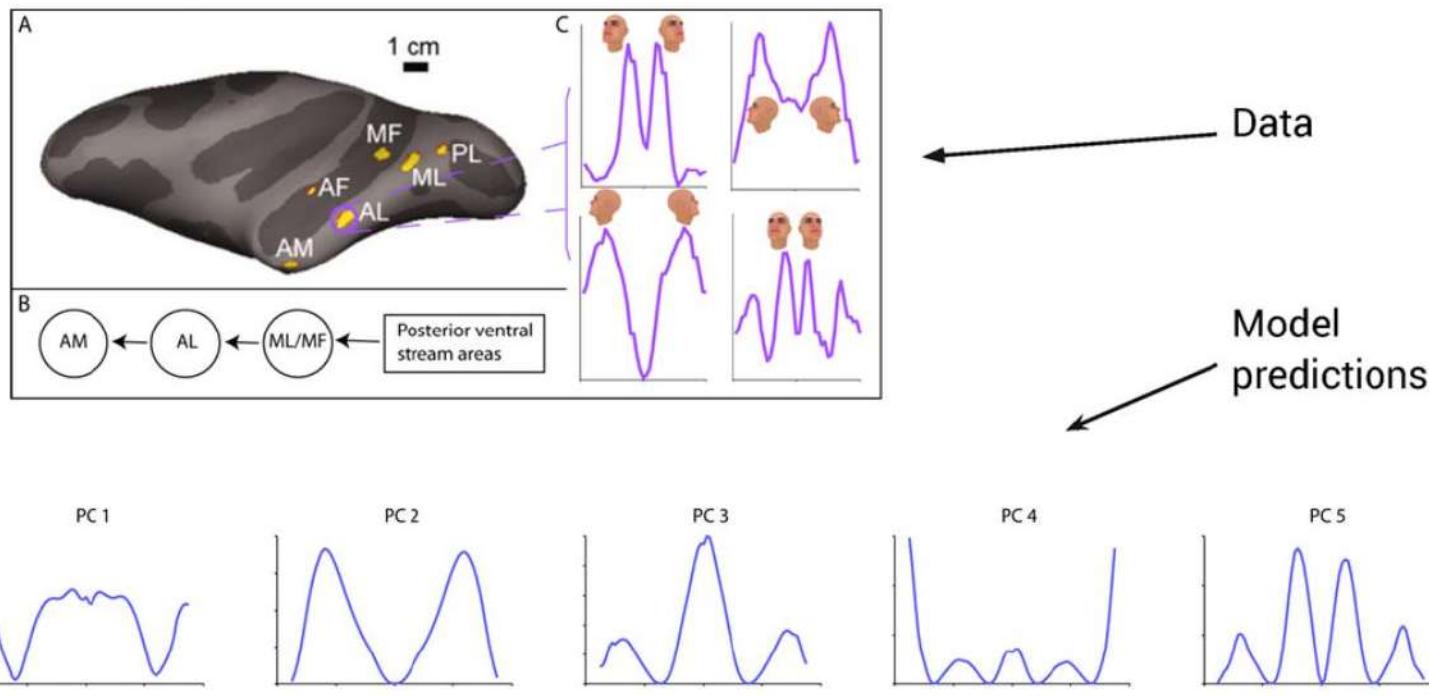
with  $F = \{f, \dots, f_N\}$  and we have

$$[C, R] = 0$$

i.e. *the eigenfunctions are such that  $Rw = w$  since  $Rw$  and  $w$  are eigenfunctions with the same eigenvalue.*

# Mirror symmetric neuronal receptive fields

Calculating the eigenfunctions of the covariance matrix:



**View-invariance implies mirror symmetric orientation tuning curves**

# Geometric deep learning

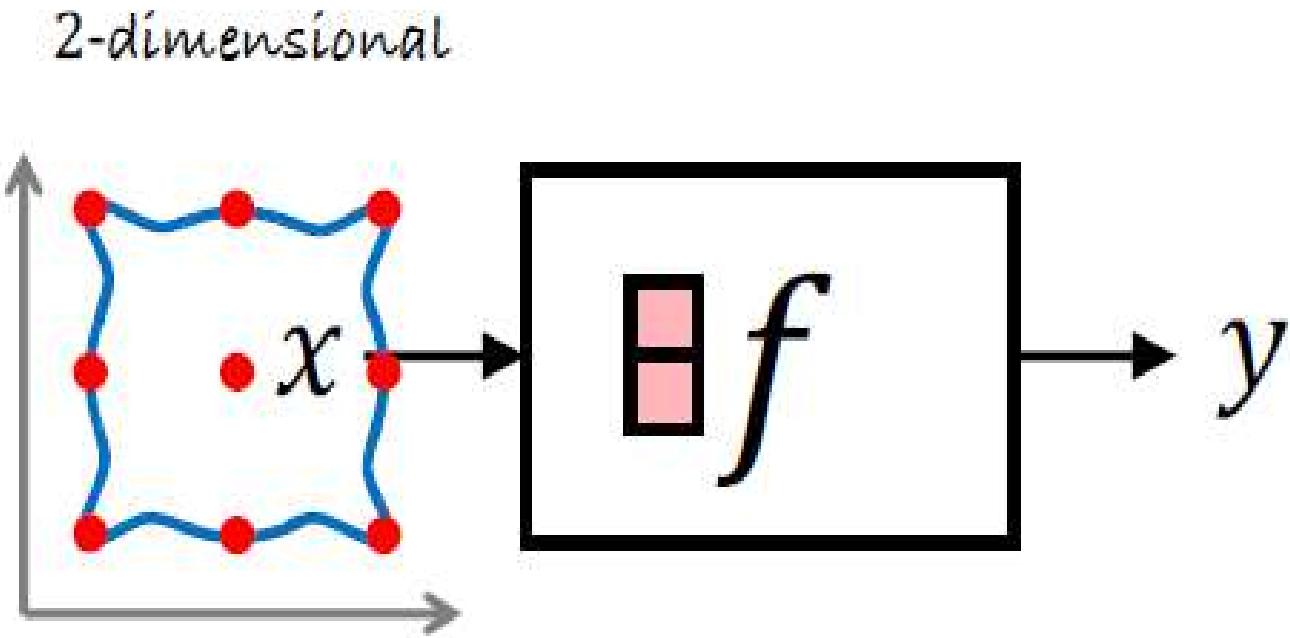
- Geometric deep learning involves developing algorithms and models for analyzing **data structured in non-euclidean domains** (graphs, networks, or manifolds)
- It combines principles from graph theory, differential geometry, and topology to enable learning and processing of complex shapes and structures that are difficult to represent in traditional deep learning frameworks.
- One important direction is the understanding of the geometry of the data and the **design of meaningful priors** to decrease the sample complexity of the learning

# Class 1

- Curse of dimensionality and geometric deep learning
- Fighting the curse with priors, examples
- Implicit regularization
- Symmetry

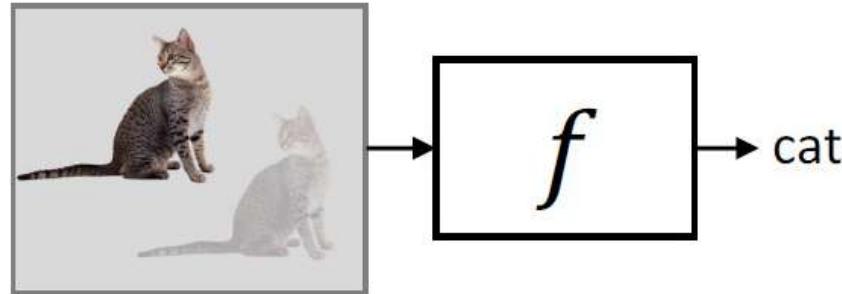
See: <https://geometricdeeplearning.com/>

# Curse of dimensionality

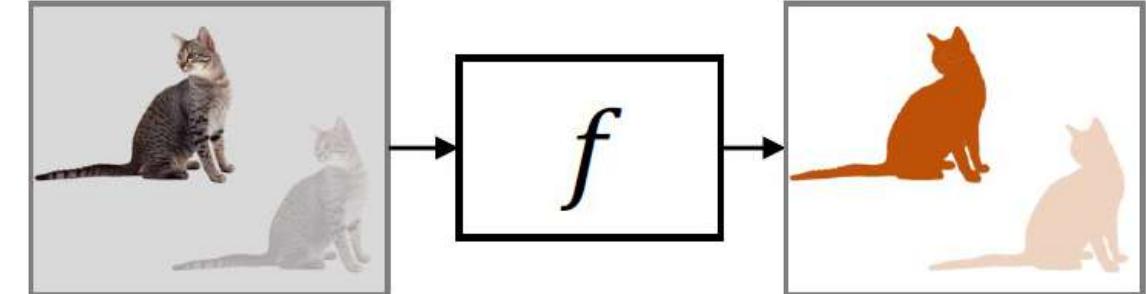


# Fight against the curse: priors on the function space. **Equivariance**

$\mathfrak{G}$ -invariance  $f(\rho(g)x) = f(x)$

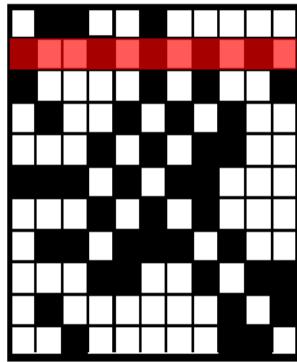


$\mathfrak{G}$ -equivariance  $f(\rho(g)x) = \rho(g)f(x)$



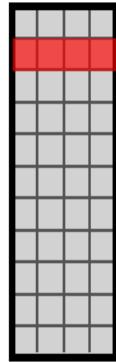
# Example with graphs, graph isomorphism

Adjacency  
matrix  $n \times n$

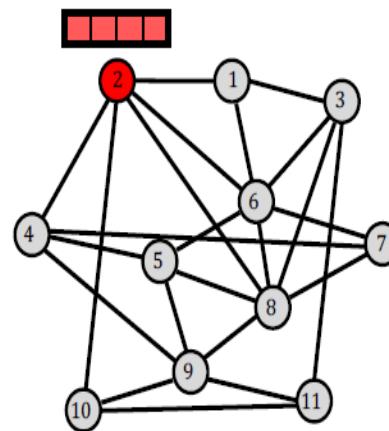


$PAP^T$

Feature  
matrix  $n \times d$

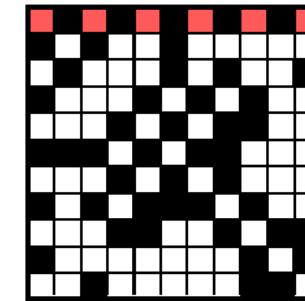


$PX$



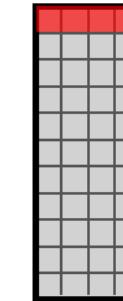
arbitrary ordering of nodes

Adjacency  
matrix  $n \times n$



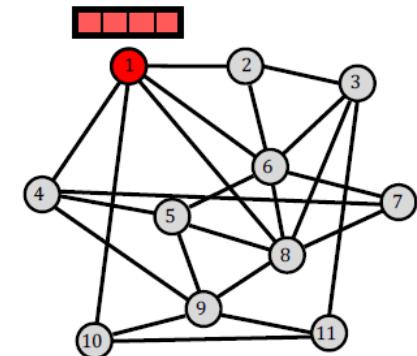
$A$

Feature  
matrix  $n \times d$



$X$

arbitrary ordering of nodes

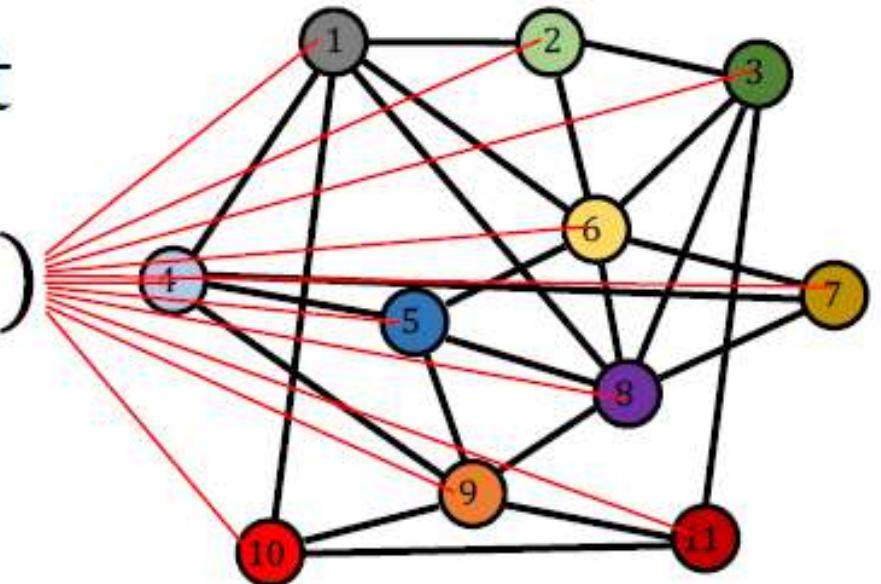


It is a very good idea to look for permutation invariant functions!

# *Invariant Graph Functions*

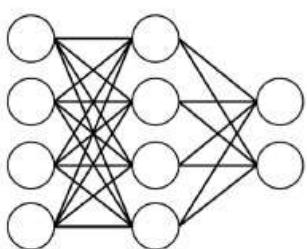
permutation-invariant

$$f(\mathbf{P}\mathbf{X}, \mathbf{P}\mathbf{A}\mathbf{P}^\top) = f(\mathbf{X}, \mathbf{A})$$

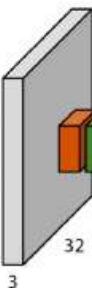


*How do you construct a permutation Invariant layer?*

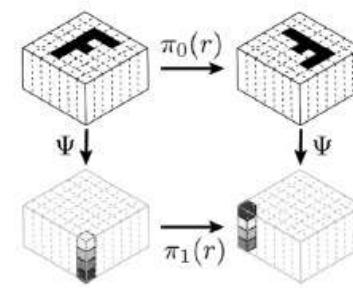
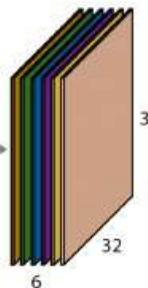
# Other examples



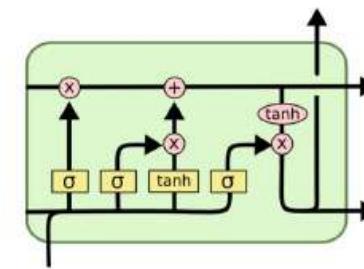
**Perceptrons**  
Function regularity



**CNNs**  
Translation



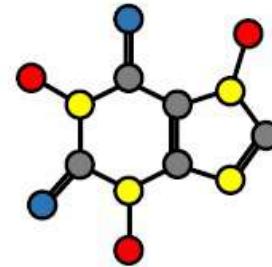
**Group-CNNs**  
Translation+Rotation,  
Global groups



**LSTMs**  
Time warping



**DeepSets / Transformers**  
Permutation



**GNNs**  
Permutation



**Intrinsic CNNs**  
Isometry / Gauge choice

# *Symmetries of the Label Function*

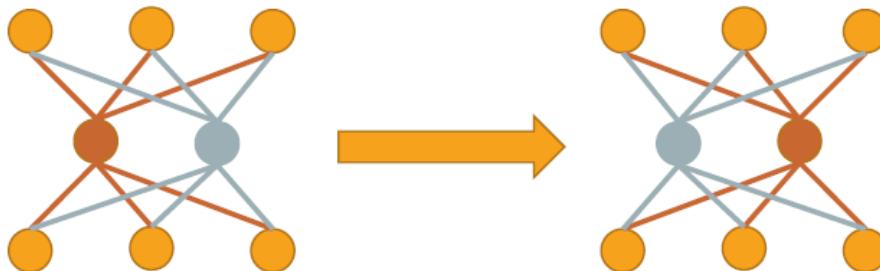
- Let  $\mathcal{X}$  denote the input space and  $\mathcal{Y}$  the label space
- Let  $L : \mathcal{X} \rightarrow \mathcal{Y}$  be the ground-truth label function
- A transformation  $g : \mathcal{X} \rightarrow \mathcal{X}$  is a symmetry if  $L \circ g = L$

$$L(\text{dog}) = L(\text{dog}) = \text{"dog"}$$


## *Symmetries of the Parameterization*

- Let  $\mathcal{X}$  denote the input space,  $\mathcal{Y}$  the label space, and  $\mathcal{W}$  the weight space
- Let  $f : \mathcal{X} \times \mathcal{W} \rightarrow \mathcal{Y}$  denote a model (e.g. neural network)
- A transformation  $g : \mathcal{W} \rightarrow \mathcal{W}$  is a symmetry of the parameterization if

$$f(x, g w) = f(x, w) \quad \text{for all } x \in \mathcal{X} \text{ and } w \in \mathcal{W}$$



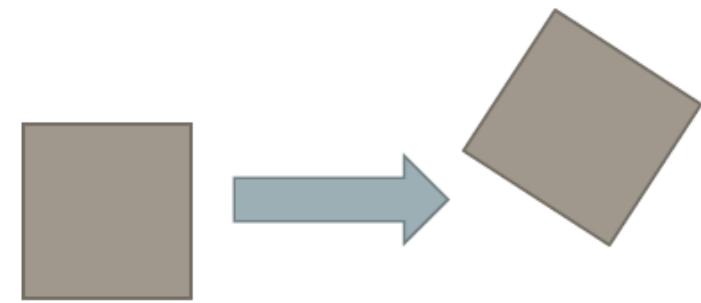
*Swapping the incoming and outgoing connections of two neurons in the same layer does not change the input-output map  $f(\cdot, w)$*

*What is the structure that allows to bridge symmetries of parameterization and data as equivalent?*

# Symmetries of the data: transformations, groups

Example: Euclidean planar motions acting on  $\mathbb{R}^2$ :

$$((\theta, t_x, t_y), (x, y)) \mapsto \begin{bmatrix} \cos \theta & \sin \theta & t_x \\ \sin \theta & \cos \theta & t_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$



# Transformations groups

- ▶ compact: rotation
- ▶ locally compact: translation, scaling (multiplication), affine  
e.g., Translation group: linear operator  $T_\tau : \mathcal{X} \rightarrow \mathcal{X}$

$$T_\tau x(p) = x(p - \tau), \quad \forall p, \tau \in \mathbb{R}, x \in \mathcal{X}$$



0	1	2	3	4	5
6	7	8	9	0	1
2	3	4	5	6	7
8	9	0	1	2	3
4	5	6	7	8	9
9	0	1	2	3	4
5	6	7	8	9	0



- ▶ non-locally compact: diffeomorphisms, local/global deformations  
e.g. given a smooth map  $d : \mathbb{R} \rightarrow \mathbb{R}$ : linear operator  $D_d : \mathcal{X} \rightarrow \mathcal{X}$

$$D_d x(p) = x(d(p)), \quad \forall p \in \mathbb{R}, x \in \mathcal{X}$$

- ▶ non-group transformations
  - e.g. 3D rotations, illumination

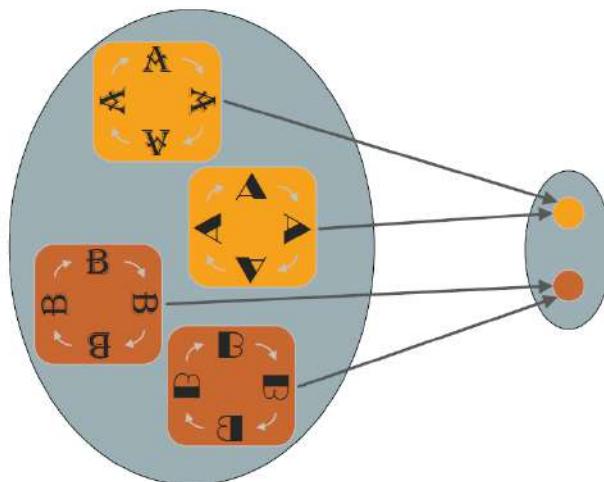
# Group properties

A *group* is a set  $\mathfrak{G}$  with a binary operation denoted  $gh$  satisfying the following properties:

- *Associativity*:  $(gh)f = g(hf)$  for all  $g, h, f \in \mathfrak{G}$
- *Identity*: there exists a unique  $e \in \mathfrak{G}$  satisfying
$$ge = eg = g$$
- *Inverse*: for each  $g \in \mathfrak{G}$  there is a unique inverse  $g^{-1} \in \mathfrak{G}$ , such that  $gg^{-1} = g^{-1}g = e$
- *Closure*: for every  $g, h \in \mathfrak{G}$ , we have  $gh \in \mathfrak{G}$

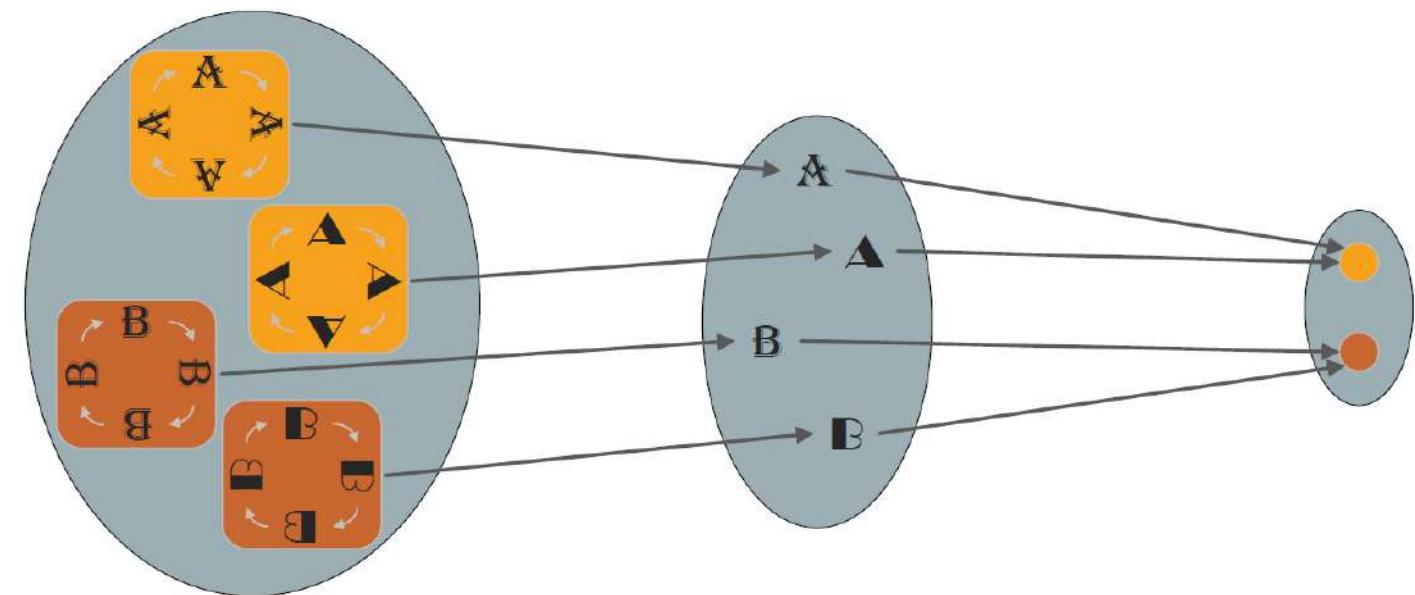
# Invariant representations

*Orbits & equivalence relations*



$$O_x = \{gx \mid x \in X, g \in G\}$$

*Invariant Representations*



# Invariant representations

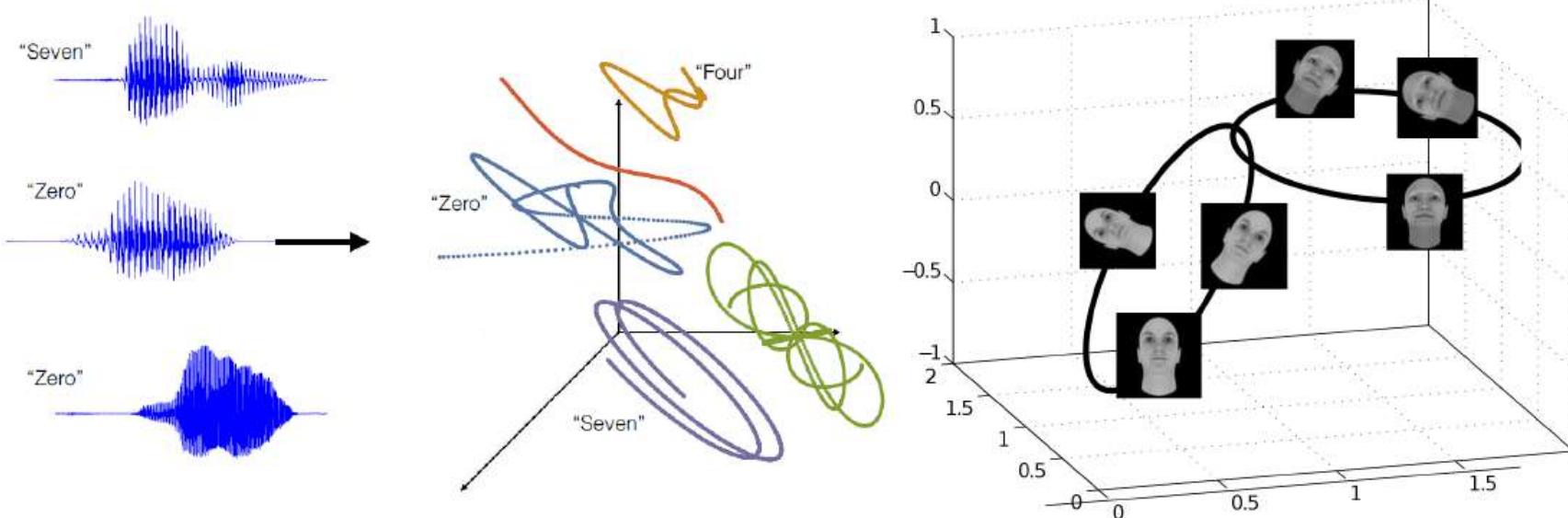
**Orbit**  $O_s = \{gs \in \mathcal{S} \mid g \in G\} \in \mathcal{S}$

- Transformations as actions of a group element  $g \in G$  on  $s$

$$s' = gs(x) = s(g^{-1}x), x \in \mathbb{R}$$

- Equivalence relation, partition of  $\mathcal{S}$

$$s \sim s' \iff \exists g \in G : s' = gs, \quad \forall s, s' \in \mathcal{S}$$



# Invariant and selective representations

Invariance:

$$x' \sim x \Rightarrow \Phi(x') = \Phi(x), \quad \forall x, x' \in \mathcal{X}$$

Selectivity:

$$x' \sim x \Leftarrow \Phi(x') = \Phi(x), \quad \forall x, x' \in \mathcal{X}$$

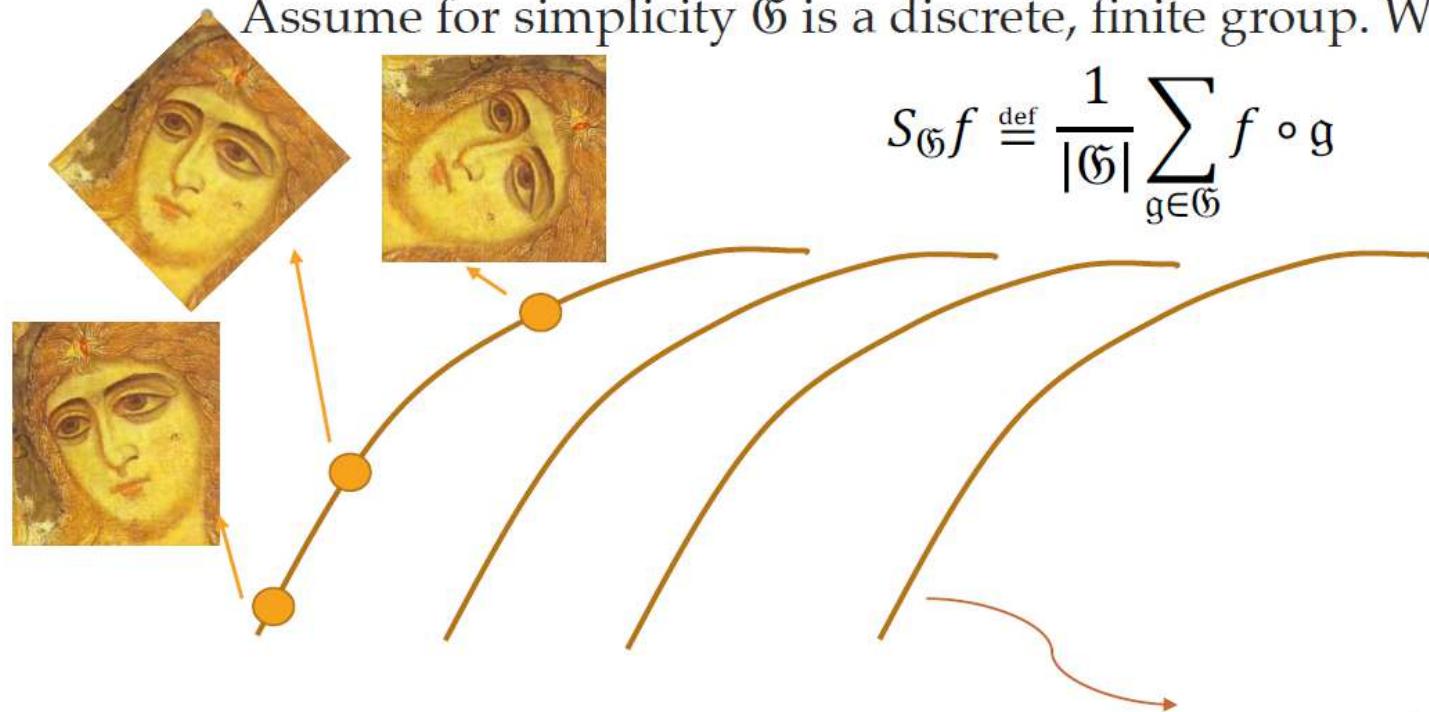
How to define the equivalence classes  $x' \sim x$ ?

Equivalence under the action of a transformation  $x' = gx, \quad g \in G$

# Invariant Function Classes

- We may first consider an abstract (non-algorithmic) option, using  $\mathfrak{G}$ -smoothing operators:

Assume for simplicity  $\mathfrak{G}$  is a discrete, finite group. We define



$$S_{\mathfrak{G}} f \stackrel{\text{def}}{=} \frac{1}{|\mathfrak{G}|} \sum_{g \in \mathfrak{G}} f \circ g$$

$$S_{\mathfrak{G}} f(x) = \frac{1}{|\mathfrak{G}|} \sum_{g \in \mathfrak{G}} f(g \cdot x)$$

$S_{\mathfrak{G}}$  thus averages a function over group orbits.

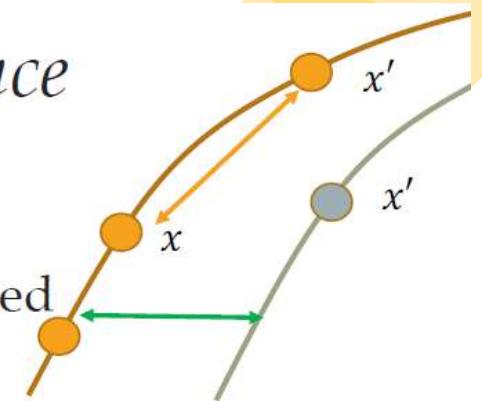
Observe that  $S_{\mathfrak{G}} f^* = f^*$ .

Group Orbits  $\mathfrak{G}.x = \{g.x; g \in \mathfrak{G}\}$

# On the Sample Complexity of Learning under Invariance

- We can consider, as before, the Lipschitz class  $\mathcal{F}$ , and its  $\mathfrak{G}$ -smoothed version.

$$f \in \mathcal{F}: |f(x) - f(x')| \leq \beta \|x - x'\|; \quad f \in S_{\mathfrak{G}} \mathcal{F}: |f(x) - f(x')| \leq \beta \inf_{g \in \mathfrak{G}} \|x - g \cdot x'\|$$



- **Theorem [BVB'21]:** Using a  $\mathfrak{G}$ -invariant kernel ridge regression, the generalization error of learning a Lipschitz,  $\mathfrak{G}$ -invariant function  $f^*$  satisfies

$$\mathbb{E} \mathcal{R}(\tilde{f}) \lesssim (\|\mathfrak{G}\| n)^{-\frac{1}{d}}$$

- Group size  $\|\mathfrak{G}\|$  can be exponential in dimension (local translations, or  $\mathcal{S}_d$ ).

## Invariant kernels and Convolutional representations

**Convolution layer representation** (with average pooling),

$$\Phi(x) = \left( \sum_g s(\langle x, gt_1 \rangle), \dots, \sum_g s(\langle x, gt_T \rangle) \right), \quad x \in \mathcal{X}.$$

The **associated kernel** can be written as

$$K(x, x') = \langle \Phi(x), \Phi(x') \rangle = \sum_t \left( \sum_g s(\langle x, gt_1 \rangle) \cdot \sum_g s(\langle x', gt_1 \rangle) \right)$$

Using linearity

$$K(x, x') = \sum_g \sum_{g'} \underbrace{\left( \sum_t s(\langle x, gt_1 \rangle) \cdot s(\langle x', g't_1 \rangle) \right)}_{K'(x, x')}.$$

## Convolution kernels

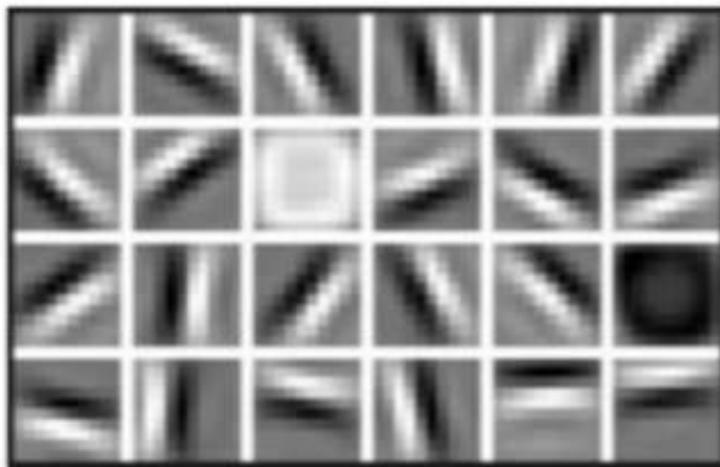
$$K(x, x') = \sum_g \sum_{g'} K'(gx, g'x')$$

- ▶ Kernel of the above form are a special case of **convolution kernels**.
- ▶ If  $g_1, \dots, g_G$  form a **group** then the kernel is invariant  
 $K(x, x') = K(x, gx')$ , i.e.  $\Phi(x) = \Phi(gx)$ .

# Compositionality prior

- What is the computational advantage?

**Low Level Features**



Lines & Edges

**Mid Level Features**



Eyes & Nose & Ears

**High Level Features**

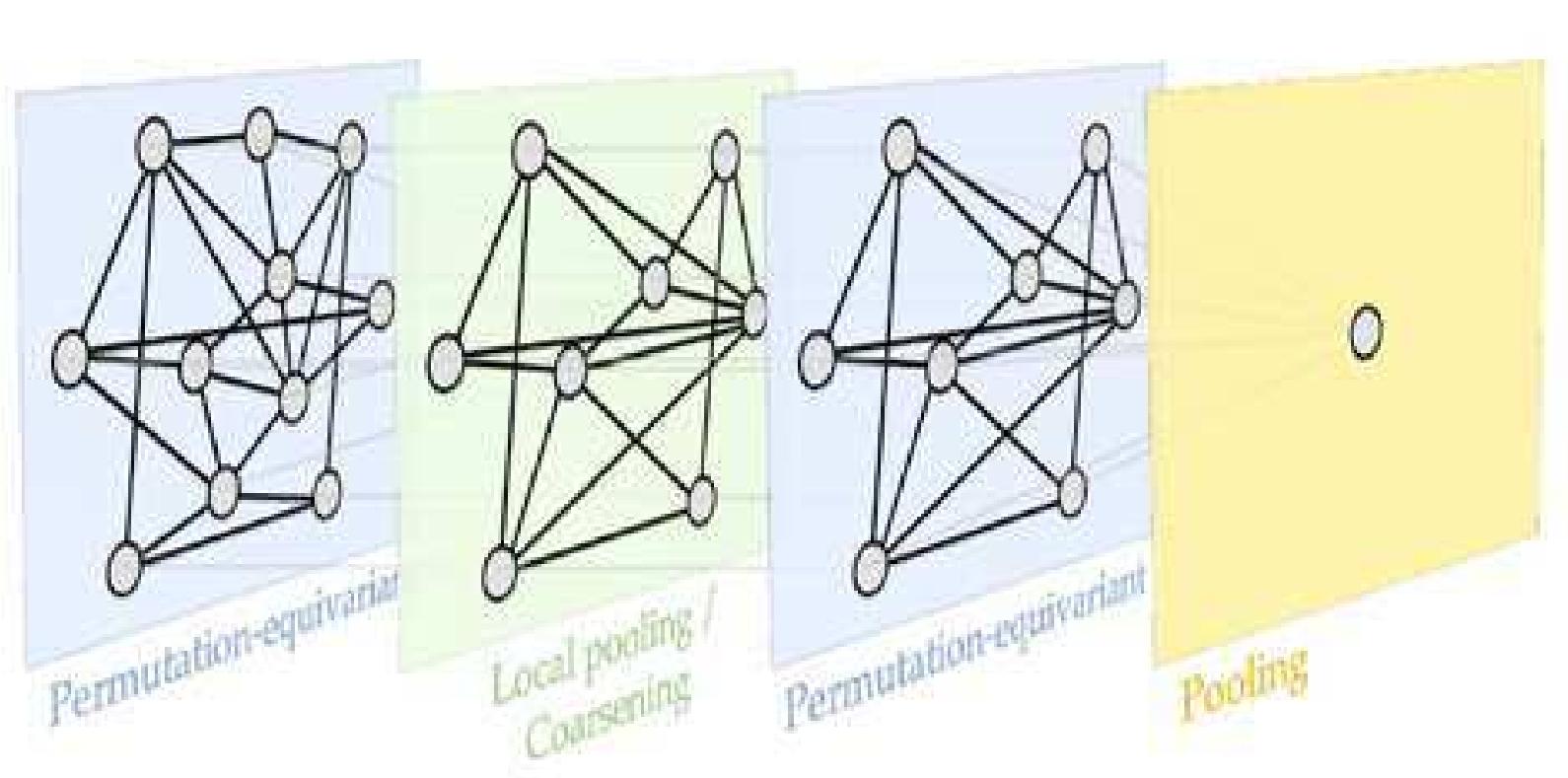


Facial Structure

# How?

- This suggests a constructive approach to build rich invariants with multiscale structure, with building blocks:
  - *Linear  $\mathfrak{G}$ -equivariant layer*  $B : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{X}(\Omega', \mathcal{C}')$ , satisfying  $B(g \cdot x) = g \cdot B(x)$  for all  $g \in \mathfrak{G}$  and  $x \in \mathcal{X}(\Omega, \mathcal{C})$ .
  - *Nonlinearity*  $\sigma : \mathcal{C} \rightarrow \mathcal{C}'$  applied element-wise as  $(\sigma(x))(u) = \sigma(x(u))$ .
  - *Local pooling (coarsening)*  $P : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{X}(\Omega', \mathcal{C})$ , such that  $\Omega' \subseteq \Omega$ .
  - *$\mathfrak{G}$ -invariant layer (global pooling)*  $A : \mathcal{X}(\Omega, \mathcal{C}) \rightarrow \mathcal{Y}$ , satisfying  $A(g \cdot x) = A(x)$  for all  $g \in \mathfrak{G}$  and  $x \in \mathcal{X}(\Omega, \mathcal{C})$ .
- These blocks can be defined under very mild conditions on the geometric domain  $\Omega$ .

# An example with graphs

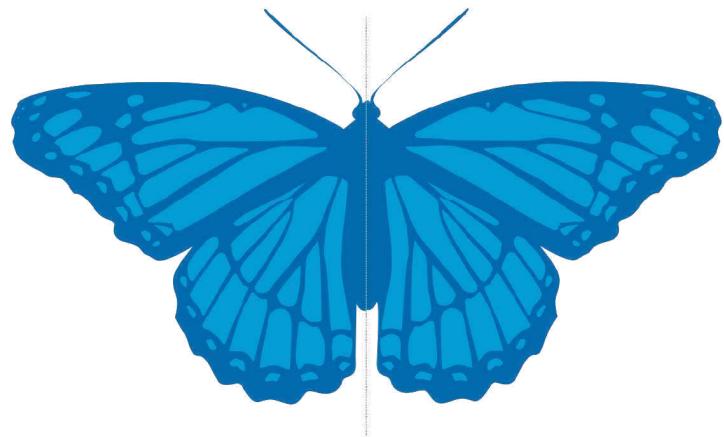


# Takeaway

---

- Geometric Priors enable a novel class of hypothesis spaces that together may tame the curse of dimensionality.
- Symmetries must be combined with scale separation.
- The language of group invariance and equivariance provides a principled architecture compatible with these priors, while preserving large approximation power.
- Theory side: what precise assumptions on non-linear coarsening operators?

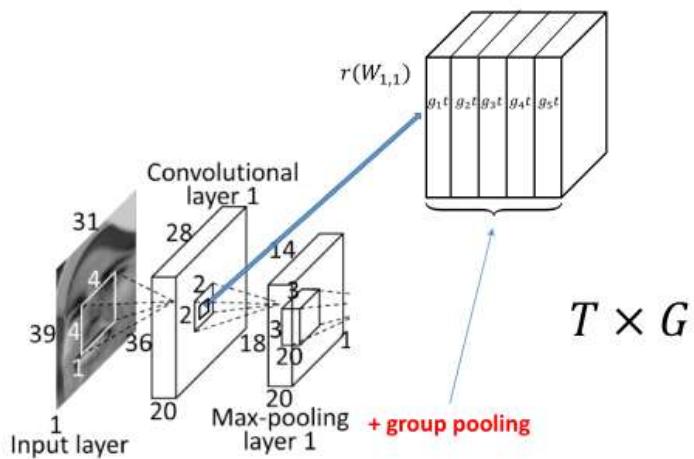
# Class 1/2: Learning symmetries



What if the symmetry is  
not evident?

## Rationale

Standard CNNs use translational convolution. A reason why they need such big training sets might be that the translational invariance is not the optimal choice to produce the best invariant representation. We want to learn the symmetry from the data.



Message: **the statistics of the data might be invariant to a broader set of symmetries than translations and to learn them will reduce the sample complexity.**

Let  $\mathcal{G}$  a group. Let  $X$  set is a collection of orbits of stimuli:

$$O_t = \{t' \mid t' = gt, g \in \mathcal{G}, t \in \mathcal{X}\}, \quad X = (O_{t_1}, \dots, O_{t_M})$$



How do we learn the group of symmetry?

## Main observation to built the regularization term

If  $W = \{g_1t, \dots, g_{|\mathcal{G}|}t\}$ ,  $t \in \mathbb{R}^d$ ,  $g \in \mathcal{G}$ ,  $\mathbb{R}^{d \times d}$

$$(W^T W)_{i,j} = \langle w_i, w_j \rangle = \langle g_i t, g_j t \rangle = \langle t, g_i^{-1} g_j t \rangle$$

Let  $\mathcal{G} = \{e, g\}$ ,  $O_t = \{t, gt\}$ :

$$W^T W = \begin{bmatrix} \langle t, e^{-1}et \rangle & \langle t, e^{-1}gt \rangle \\ \langle t, g^{-1}et \rangle & \langle t, g^{-1}gt \rangle \end{bmatrix} = \begin{bmatrix} \langle t, et \rangle & \langle t, gt \rangle \\ \langle t, gt \rangle & \langle t, et \rangle \end{bmatrix}$$

The Gramian's **columns are permutations of a single vector**.

- ① The problem can be chosen to be low dimensional.
- ② We map the problem from  $\mathcal{G}$  to the permutation group.

## Explanation: group multiplication tables and Cayley theorem

Let  $\mathcal{G} = \{e, a, b, c, d, e, f\}$ .

*	e	a	b	c	d	f	permutation
e	e	a	b	c	d	f	e
a	a	e	d	f	b	c	(12)(35)(46)
b	b	f	e	d	c	a	(13)(26)(45)
c	c	d	f	e	a	b	(14)(25)(36)
d	d	c	a	b	f	e	(156)(243)
f	f	b	c	a	e	d	(165)(234)

### Theorem (Cayley)

*Every group  $\mathcal{G}$  is isomorphic to a subgroup of the symmetric group acting on  $\mathcal{G}$ .*

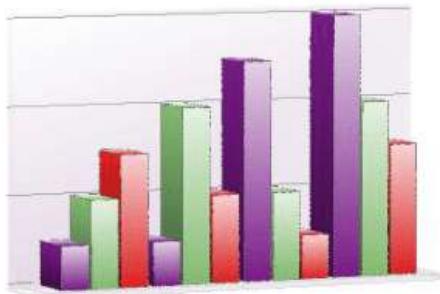
## Separating permutation orbits

We need to test if the Gramian columns  $G_{:,i}$  are permuted version of the same vector. A choice is to consider the distribution function of the vector values:

$$\mu_\lambda(G_{:,i}) = \sum_{k=1}^{|\mathcal{G}|} \delta(G_{k,i} - \lambda)$$

and impose equality of the distribution:

$$\|\mu(G_{:,i}) - \mu(G_{:,j})\|_2^2 = \sum_{k,l=1}^{|\mathcal{G}|} \int d\lambda (\mu_\lambda(G_{k,i}) - \mu_\lambda(G_{l,j}))^2 = 0$$



## Separating permutation orbits

For all pairs of Gramian columns we apply:

$$\|\mu(G_{:,i}) - \mu(G_{:,j})\|_2^2 \quad \forall i, j \Leftrightarrow \sum_{ikjl=1}^{|\mathcal{G}|} (1 - |\mathcal{G}| \delta_{ij}) \delta(w_i^T w_k - w_j^T w_l) = 0$$

or if we approximate  $\delta$  by its smooth form (letting  $\sigma \rightarrow 0$ )

$$\sum_{ikjl=1}^{|\mathcal{G}|} (1 - |\mathcal{G}| \delta_{ij}) e^{-\frac{(w_i^T w_k - w_j^T w_l)^2}{\sigma^2}} = 0.$$

More compact:

$$r(W) = \tau^T g_\sigma(C\text{vec}(W^T W)) = 0$$

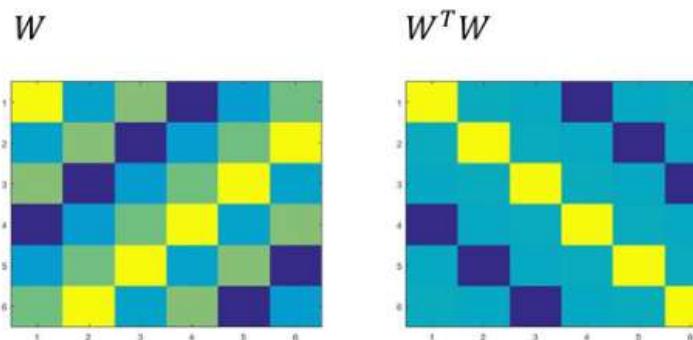
## Enforcing the Gramian constraint

### Theorem (Symmetry regularization)

Let matrix  $W \in \mathbb{R}^{d \times |\mathcal{G}|}$  and  $r : \mathbb{R}^{|\mathcal{G}| \times |\mathcal{G}|} \rightarrow \mathbb{R}_+$  with

$$r(W) = \tau^T g_\sigma(C \text{vec}(W^T W))$$

$C, \tau$  is a constant matrix and vector and function  $g_\sigma$  acts as the elementwise gaussian function. If  $r(W) = 0$ , then  $W^T W$  is a permuted matrix and viceversa.



# Alternative strategy. Google softmax

```
import sys
!git clone https://github.com/google-research/fast-soft-sort/
softsort_path = './fast-soft-sort'
sys.path.append(softsort_path)

from fast_soft_sort.pytorch_ops import soft_sort
```

## Testing the algorithm

- Test set: randomly permuted orbits of vectors w.r.t. Cyclic group ( $|\mathcal{G}| = 6$ ), Dihedral group ( $|\mathcal{G}| = 12$ ), Pyritohedral group ( $|\mathcal{G}| = 12$ ) in dimension 6.



- We calculate the intra and infra orbit elements distance with a ground truth orbit vs learned orbit. If the learned weights are an orbit of  $\mathcal{G}$  then:
  - ➊  $\|\Phi(x) - \Phi(x')\| = 0 \Rightarrow O_x = O_{x'}$
  - ➋  $\|\Phi(x) - \Phi(x')\| \neq 0 \Rightarrow O_x \neq O_{x'}$

To test invariance and selectivity we plot the distances in a confusion matrix and distance distributions.

# Recap

The following approximates the distribution of the values  $\langle g_i I, t^k \rangle$  for one template  $t^k$ :

$$\mu_n^k(I) = \sum_i (\langle g_i I, t^k \rangle)^n$$

$\forall I, I'$  images we have:

**Theorem 2** *The signature  $\mu(I) = (\mu_1^1(I), \dots, \mu_N^K(I))$*

- is **invariant** i.e.  $\mu(g_i I) = \mu(I)$ ,  $\forall g_i \in G$
- is **selective (among classes)** i.e.  $\mu(I) = \mu(I')$  iff  $I \sim I'$

If  $G$  is unitary, we have  $\langle g_i I, t^k \rangle = \langle I, g_i^{-1} t^k \rangle \Rightarrow$  we have only one orbit of an arbitrary template to implement invariance of an image seen only once.

---

## Framework for illumination transformations

\* **Contrast Functions:** continuous, monotonic, and positive functions acting on an image.

We assume the contrast function  $e : \mathbb{R}^2 \times \mathbb{R} \rightarrow \mathbb{R}^+$  has the following properties:

1.  $e(x, 0) = 1$
2.  $e(x, \zeta)e(x, \zeta') = e(x, \zeta \circ \zeta')$ ,  $\zeta, \zeta' \in \mathbb{R}_+$

where  $\circ$  is a group composition. The transformation of an image  $I$  is:

$$I_\zeta(x) \equiv I(x, \zeta) = e(x, \zeta)I(x, 0).$$

and similarly for template  $t$ . We have the following signature for an image under change in illumination:

$$\mu_n^k(I) = \int_0^\infty (\langle I_0, t_\zeta^k \rangle)^n d\zeta$$

# Protocol

\* **Invariance:** We calculated the euclidean distance between images and their transformations. For invariance we tested:

$$\|\mu(I_0) - \mu(I_\zeta)\|^2 \sim 0$$

where we fixed 1 template, with  $\mu(I) \in \mathbb{R}^{n \times k}$ , and averaged the euclidean distance over images and their transformations.

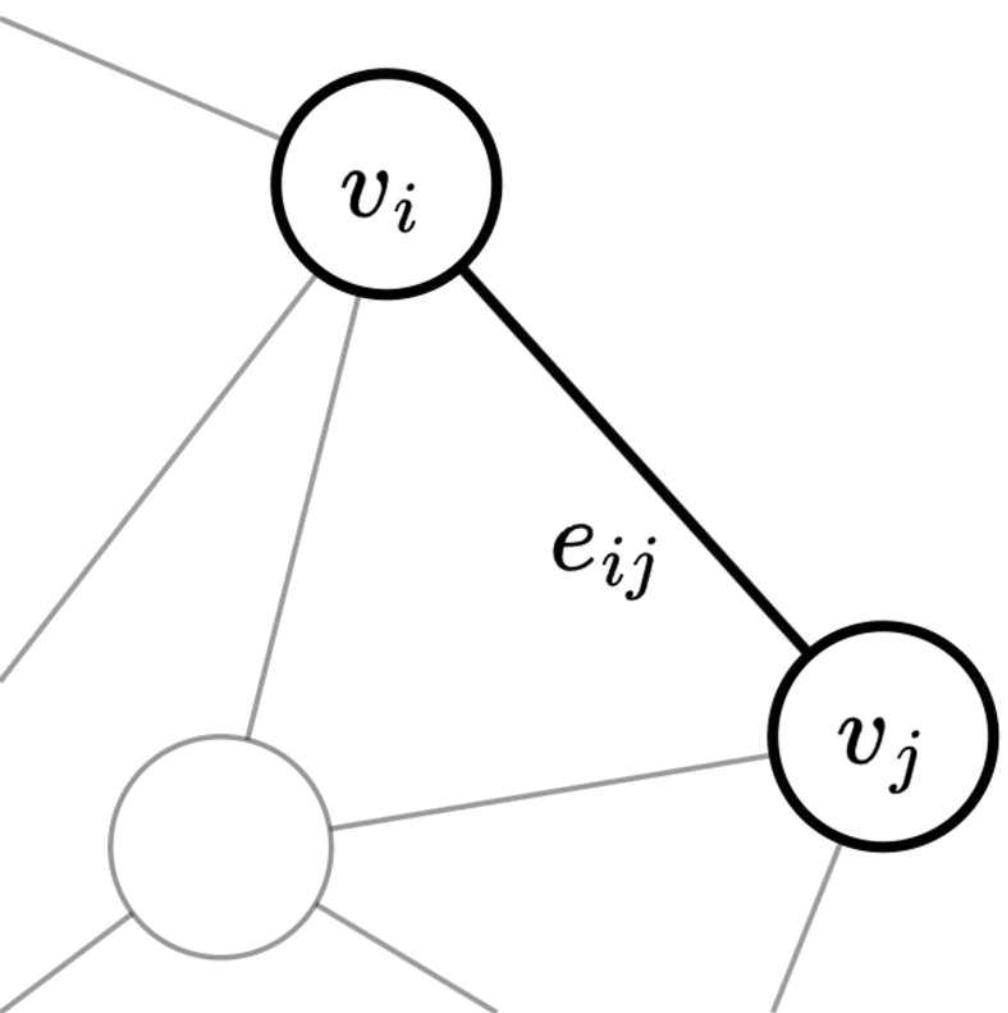
\* **Control:** We defined a “fake orbit” as a set of randomly selected images from the datasets.

\* **Selectivity:** We calculated the euclidean difference between different images. For selectivity we tested:

$$\|\mu(I_0^i) - \mu(I_0^j)\|^2 \not\sim 0, \quad i \neq j$$

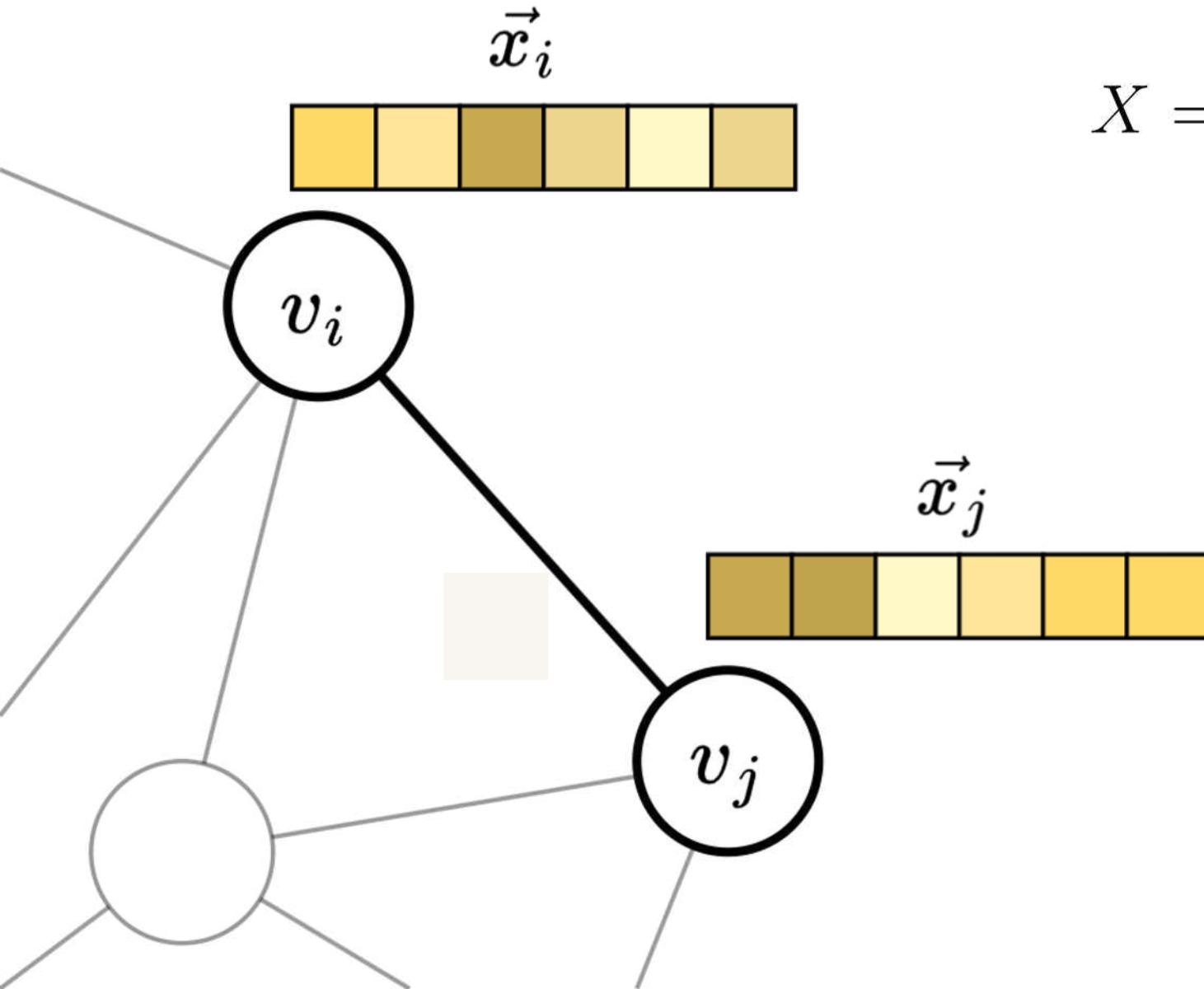
For each of the datasets, we let  $I_0$  be a face under one illumination condition and  $I_\zeta$  be all other illumination conditions of the face. Our set of templates  $\{t_\zeta^k\}$ , contained faces under all illumination conditions with  $t^k \neq I$ .

$$\mathcal{G} = (V, E)$$



**Adjacency matrix of  $\mathcal{G}$**

$$A_{ij} = \begin{cases} 0 & \text{if } (i, j) \notin E \\ 1 & \text{if } (i, j) \in E \end{cases}$$



$$X = \begin{pmatrix} \vec{x}_1 \\ \vdots \\ \vec{x}_i \\ \vdots \\ \vec{x}_n \end{pmatrix}$$

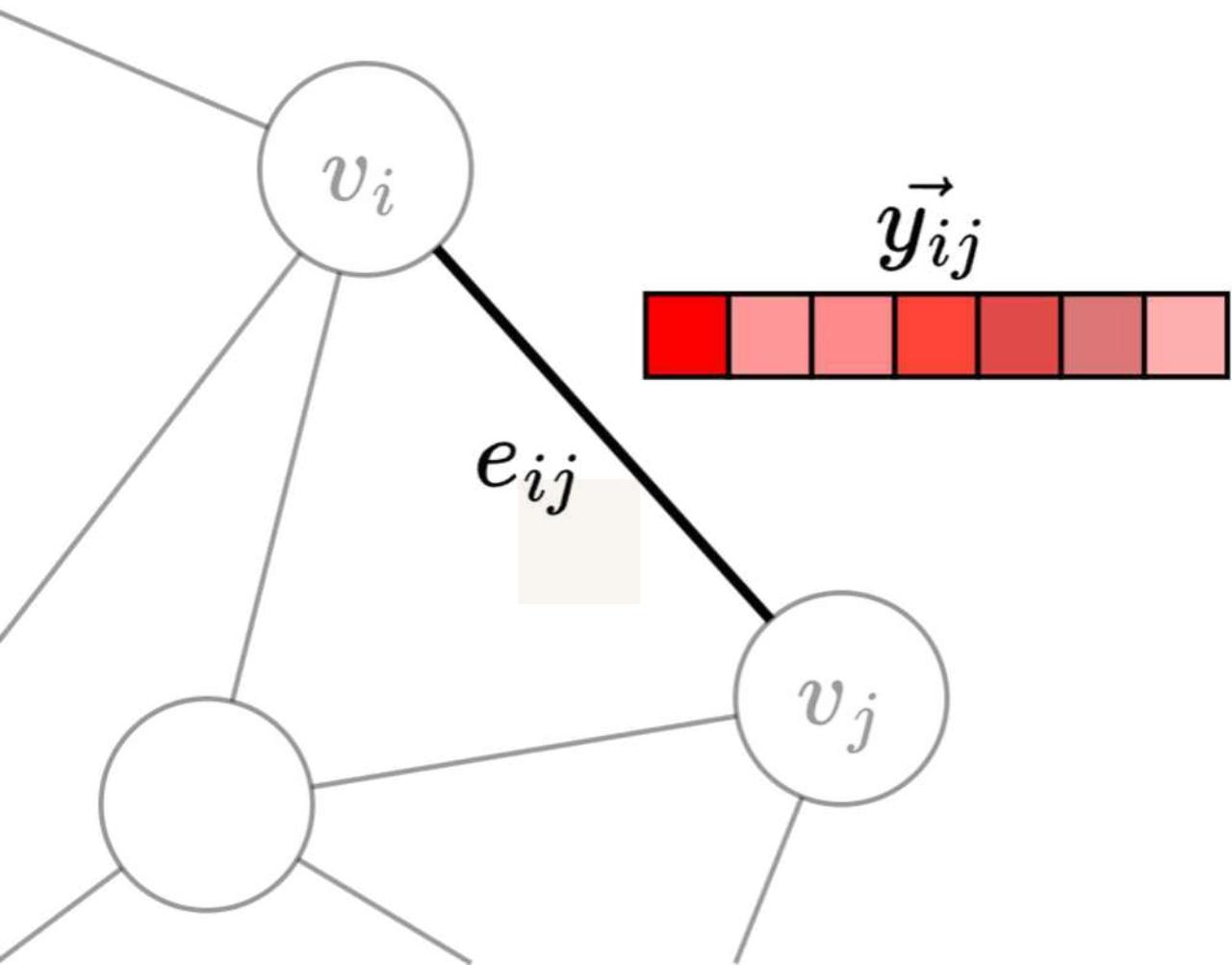
**node features of  $\mathcal{G}$**

for example:

- properties of the atom
- data of the account
- number/type of crossroads

$$Y_{ij} = \vec{y}_{ij}$$

edge features of  $\mathcal{G}$



for example:

- type of bond
- date of friend request, relationship
- average traffic, number of car parks

$$\vec{g}$$



global features of  $\mathcal{G}$

for example:

- known properties
- year
- day of the week

# Tasks on graphs

## Global tasks

tasks on the **whole graph**.

Both classification and regression.

Examples:

- predicting properties of a molecule
- predicting the hobby shared by all the people
- predicting the likelihood of a car accident

# Tasks on graphs

## Node tasks

tasks at the level of the **single node**.  
Both classification and regression.

Examples:

- predicting if an atom is part of an aromatic ring
- predicting political stance of each users on a topic
- classifying intersections based on predicted traffic

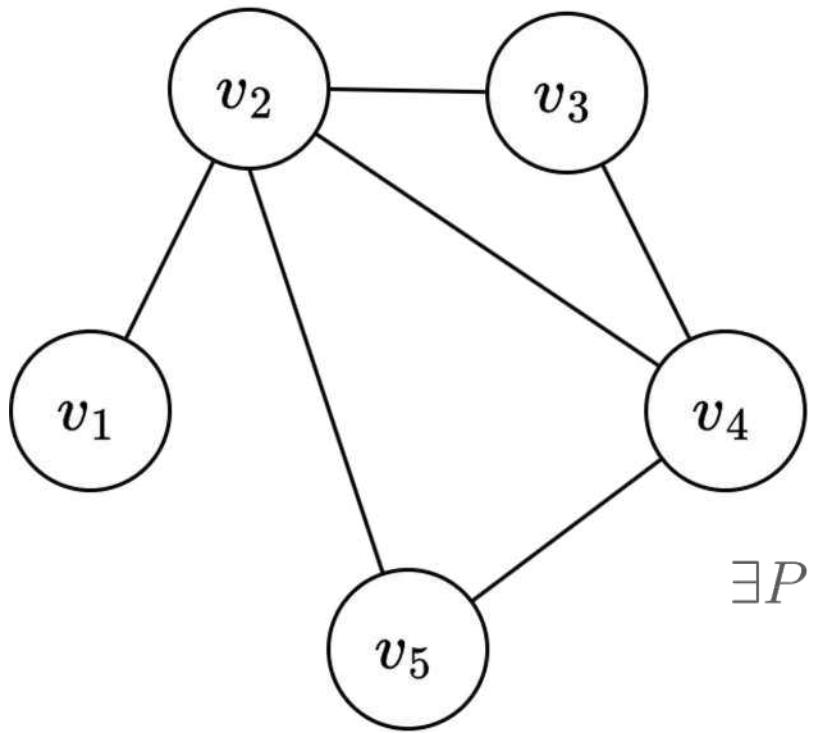
# Tasks on graphs

tasks at the level of the **single edge**.  
Both classification and regression.

Examples:

- predicting which bond is easier to break
- predicting how likely two users will become friends
- predicting how busy a road will be in an hour

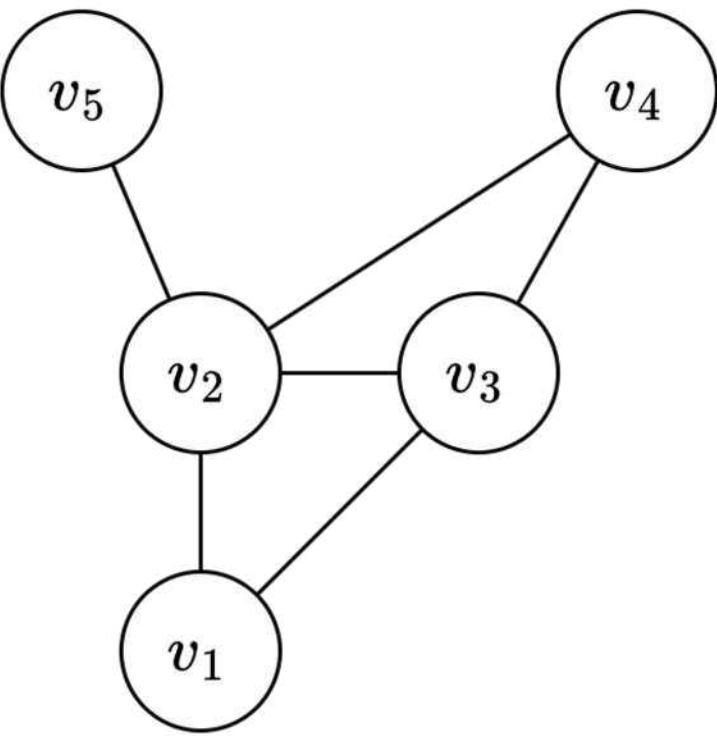
Edge tasks



$\Leftarrow \Rightarrow$

$\exists P$  permutation matrix s.t.

$$PA_1P^{-1} = A_2$$



$$A_1 = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

$$A_2 = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix}$$

$\forall P$  permutation matrix

**invariance** in the context of **global tasks**

$$F_G(A, X, Y, \vec{g}) = F_G(PAP^{-1}, PX, PYP^{-1}, \vec{g}) \in \mathbb{R} \text{ or } \in \mathbb{R}^d$$

**equivariance** in the context of **node** or **edge** tasks

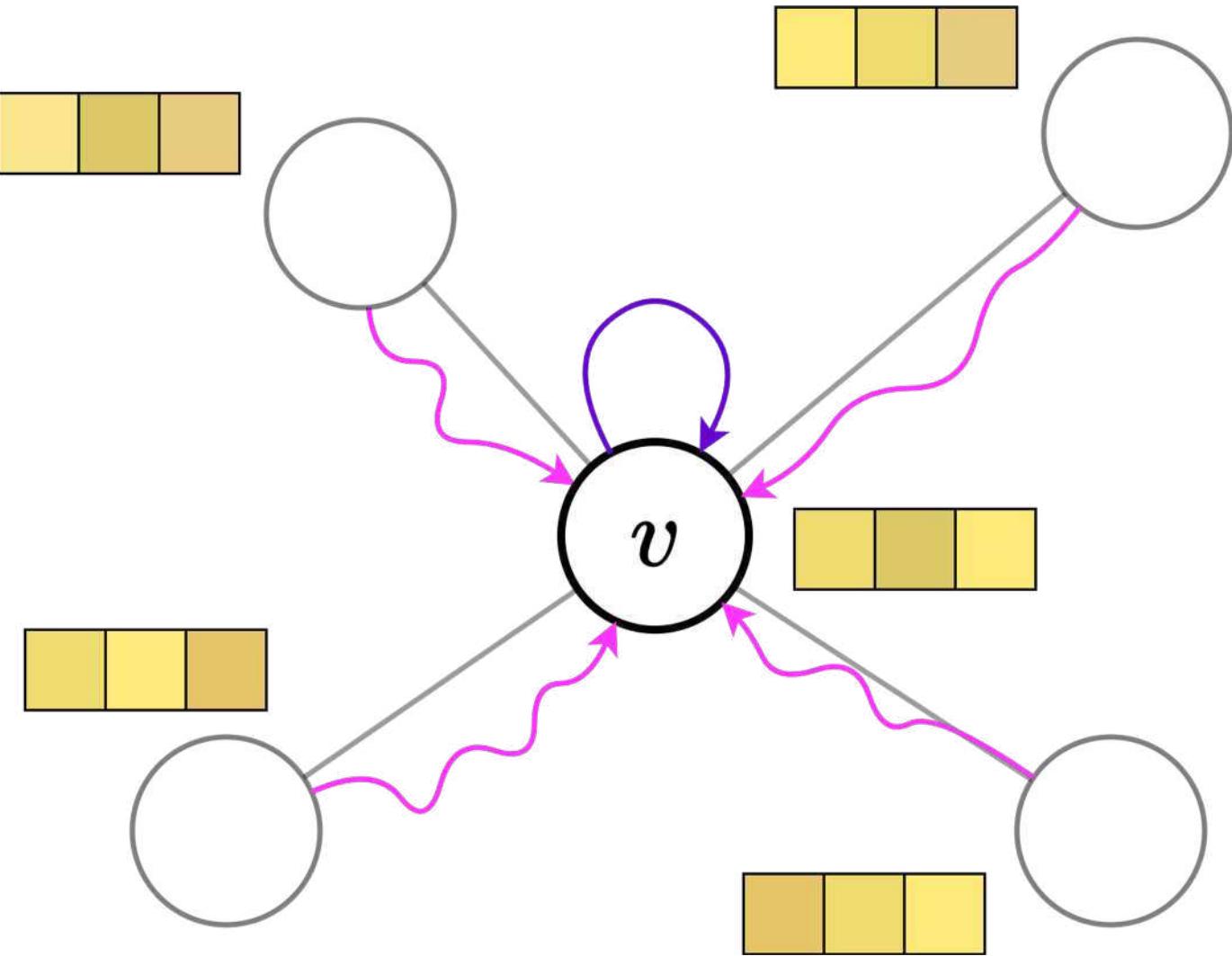
$$\begin{aligned} F_V(A, X, Y, \vec{g}) &= P^{-1} F_V(PAP^{-1}, PX, PYP^{-1}, \vec{g}) \\ &\in \mathbb{R}^n \text{ or } \in \mathbb{R}^{n \cdot d} \end{aligned}$$

$$\begin{aligned} F_E(A, X, Y, \vec{g}) &= P^{-1} F_E(PAP^{-1}, PX, PYP^{-1}, \vec{g}) P \\ &\in \mathcal{M}_{\mathbb{R}}(n, n) \text{ or } \in \mathcal{M}_{\mathbb{R}^d}(n, n) \end{aligned}$$

in the context of node/global tasks

input:  $A, X$

- **Aggregating** over  $\mathcal{N}_v$
- **Updating**  $\vec{x}_v$



# Convolutional approach

$$\vec{x}_v = \phi \left( \vec{x}_v, \bigoplus_{u \in \mathcal{N}_v} c_{vu} \psi(\vec{x}_u) \right)$$

**Update**

$$\phi(\vec{x}, \vec{z}) = \sigma(W\vec{x} + U\vec{z} + \vec{b})$$

learnable

$$\vec{x}_v = \phi \left( \vec{x}_v, \bigoplus_{u \in \mathcal{N}_v} c_{vu} \psi(\vec{x}_u) \right)$$

**coefficient**  
in the case of summation  
can be

$$c_{uv} = \frac{1}{\sqrt{\deg(u) \cdot \deg(v)}}$$

must be a **permutation invariant operator**

- mean
- max
- sum

**transformation of the features**

$$\psi(\vec{x}) = W\vec{x} + \vec{b}$$

\ /  
learnable

# Attentional approach

$$\vec{x}_v = \phi \left( \vec{x}_v, \bigoplus_{u \in \mathcal{N}_v} a(\vec{x}_v, \vec{x}_u) \psi(\vec{x}_u) \right)$$



coefficients are learned  
via NN and softmax  
normalized

$$a(\vec{x}_v, \vec{x}_u) \in \mathbb{R}$$

# Message-passing approach

$$\vec{x}_v = \phi \left( \vec{x}_v, \bigoplus_{u \in \mathcal{N}_v} \psi(\vec{x}_v, \vec{x}_u) \right)$$



$\psi$  here is a learnable  
message function

$$\psi(\vec{x}_v, \vec{x}_u) \in \mathbb{R}^m$$

# 1-neuron ReLu implicit bias

**Notations.** We use bold-faced letters to denote vectors, e.g.,  $\mathbf{x} = (x_1, \dots, x_d)$ . For  $\mathbf{x} \in \mathbb{R}^d$  we denote by  $\|\mathbf{x}\|$  the Euclidean norm.

**Single-neuron networks.** Let  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  be a training dataset, where for every  $i$  we have  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$ . We consider empirical-loss minimization of a single-neuron network, with respect to the square loss. Thus, the objective is given by

$$\mathcal{L}(\mathbf{w}) := \frac{1}{2} \sum_{i=1}^n (\sigma(\langle \mathbf{x}_i, \mathbf{w} \rangle) - y_i)^2 , \quad (1)$$

where  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is a nonlinear activation function. Let  $X \in \mathbb{R}^{n \times d}$  denote the data matrix, i.e., the rows of  $X$  are  $\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top$ , and let  $\mathbf{y} = (y_1, \dots, y_n)$ . We have

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2} \|\sigma(X\mathbf{w}) - \mathbf{y}\|^2 ,$$

where  $\sigma$  is applied component-wise. We assume that the data is realizable, that is,  $\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = 0$ . Moreover, we focus on settings where the network is *overparameterized*, in the sense that  $\mathcal{L}$  has multiple (or even infinitely many) global minima.

We analyze the implicit regularization of gradient flow on the objective given by Eq. 1. This setting captures the behaviour of gradient descent with infinitesimally small step size. Let  $\mathbf{w}(t)$  be the trajectory of gradient flow, where  $\mathbf{w}(0)$  is the initial point. The dynamics of  $\mathbf{w}(t)$  is given by the differential equation

$$\dot{\mathbf{w}}(t) := \frac{d\mathbf{w}(t)}{dt} = -\nabla \mathcal{L}(\mathbf{w}(t)) .$$

If  $\lim_{t \rightarrow \infty} \mathbf{w}(t)$  exists then we denote it by  $\mathbf{w}(\infty)$ . We note that gradient flow is not guaranteed to converge to a global minimum (cf. Yehudai and Shamir [2020]). In the overparameterized setting there can be infinitely many global minima, and we study to which one gradient flow converges, assuming that it converges to a global minimum.

## The implicit regularization of a ReLU neuron is approximately the $\ell_2$ norm

While the implicit regularization of ReLU neuron cannot be expressed as a function  $\mathcal{R}(\mathbf{w})$ , in this section we show that it can be expressed approximately, within a factor of 2, by the  $\ell_2$  norm. This implies that even without early stopping, if the data can be labeled by a ReLU neuron with small  $\ell_2$  norm, then gradient flow will converge to a ReLU neuron whose  $\ell_2$  norm is not much larger. Since a ReLU neuron is just a linear function composed with a fixed nonlinearity, this can be used to derive good statistical generalization guarantees, via standard techniques (cf. Shalev-Shwartz and Ben-David [2014]).

**Theorem 5.1.** *Consider gradient flow on the objective given by Eq. 1, where  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is a monotonically non-decreasing activation function (e.g., ReLU). Assume that  $\mathbf{w}(\infty)$  exists and  $\mathcal{L}(\mathbf{w}(\infty)) = 0$ . Let  $\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w}} \|\mathbf{w} - \mathbf{w}(0)\|$  s.t.  $\mathcal{L}(\mathbf{w}) = 0$ . Then,  $\|\mathbf{w}(\infty) - \mathbf{w}(0)\| \leq 2 \cdot \|\mathbf{w}^* - \mathbf{w}(0)\|$ .*

*Proof.* First, note that

$$\begin{aligned} (\nabla \mathcal{L}(\mathbf{w}(t)), \mathbf{w}(t) - \mathbf{w}^*) &= \left\langle \sum_{i=1}^n \left( \sigma(\mathbf{x}_i^\top \mathbf{w}(t)) - \sigma(\mathbf{x}_i^\top \mathbf{w}^*) \right) \sigma'(\mathbf{x}_i^\top \mathbf{w}(t)) \mathbf{x}_i, \mathbf{w}(t) - \mathbf{w}^* \right\rangle \\ &= \sum_{i=1}^n \left( \sigma(\mathbf{x}_i^\top \mathbf{w}(t)) - \sigma(\mathbf{x}_i^\top \mathbf{w}^*) \right) \sigma'(\mathbf{x}_i^\top \mathbf{w}(t)) (\mathbf{x}_i^\top \mathbf{w}(t) - \mathbf{x}_i^\top \mathbf{w}^*) . \end{aligned}$$

Since  $\sigma$  is monotonically non-decreasing then

$$\left( \sigma(\mathbf{x}_i^\top \mathbf{w}(t)) - \sigma(\mathbf{x}_i^\top \mathbf{w}^*) \right) (\mathbf{x}_i^\top \mathbf{w}(t) - \mathbf{x}_i^\top \mathbf{w}^*) \geq 0 ,$$

and  $\sigma'(\mathbf{x}_i^\top \mathbf{w}(t)) \geq 0$ . Hence,

$$\langle \nabla \mathcal{L}(\mathbf{w}(t)), \mathbf{w}(t) - \mathbf{w}^* \rangle \geq 0.$$

Therefore, we have

$$\frac{d}{dt} \left( \frac{1}{2} \|\mathbf{w}(t) - \mathbf{w}^*\|^2 \right) = \langle \mathbf{w}(t) - \mathbf{w}^*, \dot{\mathbf{w}}(t) \rangle = \langle \mathbf{w}(t) - \mathbf{w}^*, -\nabla \mathcal{L}(\mathbf{w}(t)) \rangle \leq 0.$$

Thus,  $\|\mathbf{w}(\infty) - \mathbf{w}^*\| \leq \|\mathbf{w}(0) - \mathbf{w}^*\|$ .

Hence,

$$\begin{aligned} \|\mathbf{w}(\infty) - \mathbf{w}(0)\| &= \|\mathbf{w}(\infty) - \mathbf{w}^* + \mathbf{w}^* - \mathbf{w}(0)\| \\ &\leq \|\mathbf{w}(\infty) - \mathbf{w}^*\| + \|\mathbf{w}^* - \mathbf{w}(0)\| \\ &= 2 \cdot \|\mathbf{w}^* - \mathbf{w}(0)\|. \end{aligned}$$

□

**Remark 5.1.** Note that Theorem 5.1 implies that if  $\mathbf{w}(0) = 0$  then  $\|\mathbf{w}(\infty)\| \leq 2 \cdot \|\mathbf{w}^*\|$ . By Theorem 4.2, the implicit regularization cannot be expressed as a function of  $\mathbf{w}$ . Hence, while the implicit regularization of a ReLU neuron cannot be expressed exactly as a function of  $\mathbf{w}$ , it can be expressed approximately, within a factor of 2, by the  $\ell_2$  norm.

# Non-linearity choice and implicit bias

For the ReLu expansion of a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  we have:

$$f(x) = \int d\beta c(\beta) ReLu(x - \beta) \quad (6)$$

with

$$f'(x) = \int d\beta c(\beta) Step(x - \beta) \quad (7)$$

Taking into account of saturation we can use instead of a ReLu the following:

$$SReLU(x) = \frac{1}{w} (ReLU(x - \beta) - ReLu(x - (\beta + w))) \quad (8)$$

where the factor  $1/w$  is a normalization factor. Indeed

$$f'(x) = \int d\beta d(\beta) \frac{1}{w} Box_w(x - \beta) = \int d\beta d(\beta) \bar{Box}_w(x - \beta), \quad (9)$$

where  $(1/w)Box = \bar{Box}$  is the Box function at  $\beta + w$  with integral normalized to unit.

How can we pass from one representation to the other? Note that:

$$Box_{\beta,w}(x) = Step_{\beta}(x) - Step_{\beta+w}(x) = (Id - T_w)Step_{\beta}(x) \equiv \Delta_w Step_{\beta}(x) \quad (10)$$

where  $T : \mathbb{R} \rightarrow \mathbb{R}$ ,  $T_w f(x) = f(x + w)$  is a translation operator.

Thus we have

$$d(\beta) = \Delta_w^{-1} c(\beta) \quad (11)$$

where the operator can be explicitly written as a sum of translation operators since:

$$\Delta_w^{-1} = (Id - T_w)^{-1} = \sum_{n \in \mathbb{N}} T_w^n. \quad (12)$$

Indeed:

$$Step_{\beta}(x) = \Delta_w^{-1} Box_{\beta}(x) = \sum_{n \in \mathbb{N}} T_w^n Box_{\beta+w}(x). \quad (13)$$

Explicitly:

$$d(\beta) = \sum_{n \in \mathbb{N}} c(\beta - nw) = \sum_{n \in \mathbb{N}} \frac{1}{w} f''(\beta - nw). \quad (14)$$

Check for constants.

# Non-linearity choice and implicit bias

Supposing an  $N = \text{ceil}(x_{\max}/w)$  maximum (finite range of neural activations):

$$\begin{aligned}\|f''\|_2^2 &= \int d\beta c^2(\beta) = \int d\beta \left| \sum_{n=0}^{N-1} f''(\beta - nw) \right|^2 = \int dk \left| \mathcal{F} \left( \sum_{n=0}^{N-1} \frac{1}{w} T_{nw} \partial_\beta^2 f(\beta) \right) \right|^2 \\ &= \int dk k^4 \left| \sum_{n=0}^{N-1} \frac{e^{inwk}}{w} \right|^2 \hat{f}^2(k).\end{aligned}$$

Using

$$\left| \sum_{n=0}^{N-1} \frac{e^{inwk}}{w} \right|^2 = \frac{\sin^2(\frac{1}{2}Nwk)}{w^2 \sin^2(\frac{1}{2}wk)}$$

we finally have:

$$\boxed{\|f''\|_2^2 = \int dk \frac{k^4}{w^2} \frac{\sin^2(\frac{1}{2}Nwk)}{\sin^2(\frac{1}{2}wk)} \hat{f}^2(k)}$$

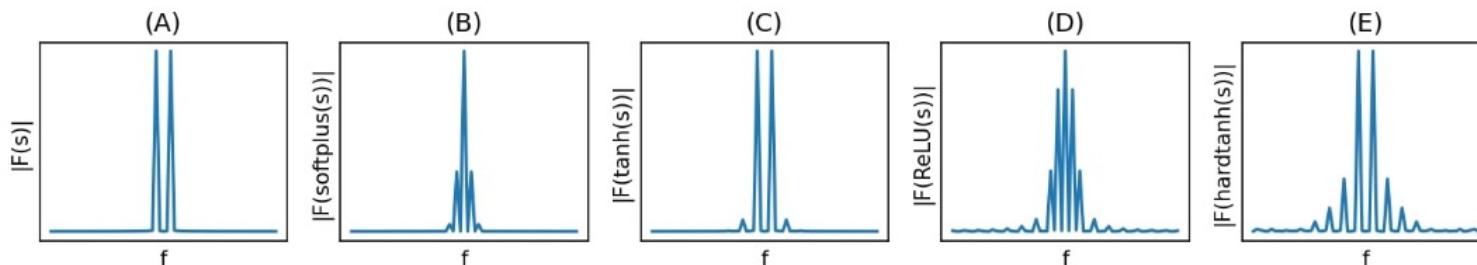
To compare in the ReLu case a similar calculation would lead to:

$$\|f''\|_2^2 = \int dk k^4 \hat{f}^2(k).$$

Suppose  $f(t) = \cos(w_1 t) + \cos(w_2 t)$  is a sound wave and let  $\sigma(t) = t^2$ . Then

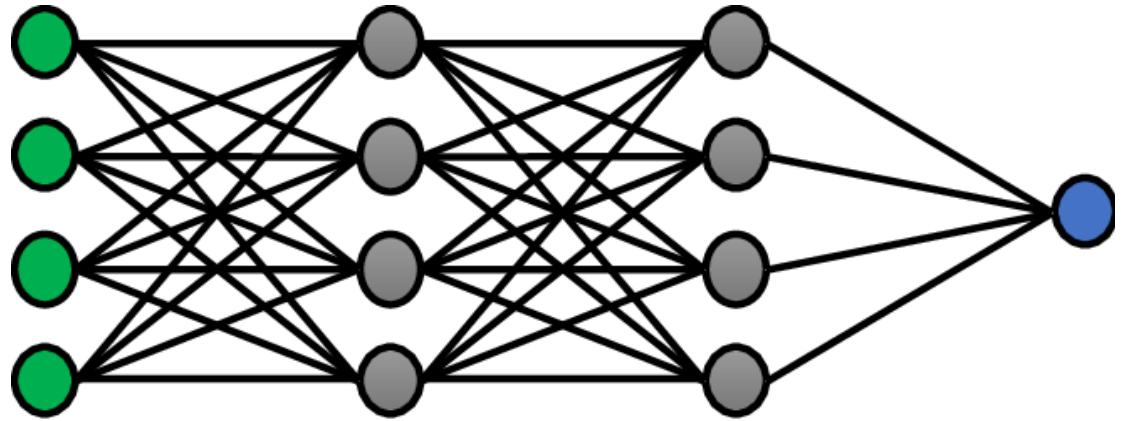
$$(\sigma \circ f)(t) = \frac{1}{2}[2 + \cos(2w_1 t) + \cos(2w_2 t) + 2 \cos((w_1 + w_2)t) + 2 \cos((w_1 - w_2)t)].$$

$$FT\phi(x) = F \sum_n a_n \underbrace{x \odot \cdots \odot x}_{n\text{-times}} = \sum_n a_n \underbrace{\hat{x} * \cdots * \hat{x}}_{n\text{-times}},$$

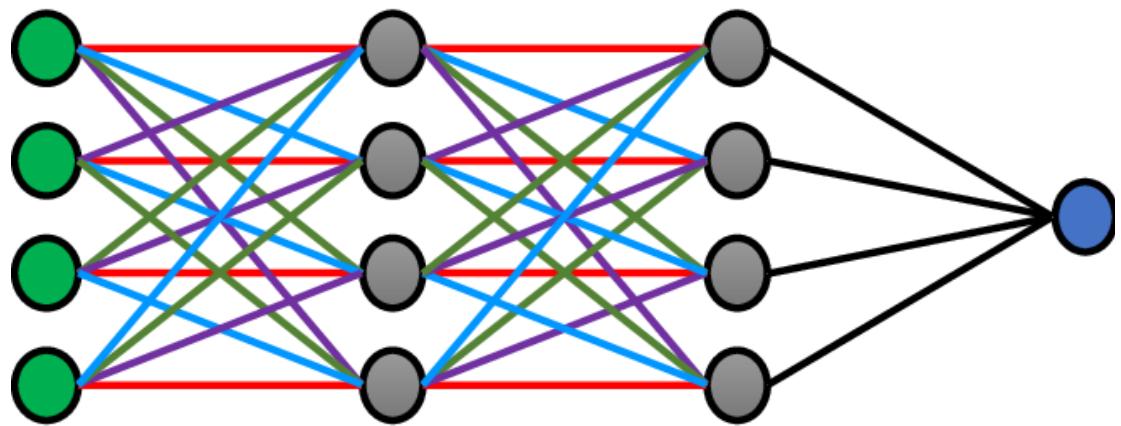


**Figure 1.** Non-linear distortions in the frequency domain due to the application of (B) softplus, (C) tanh, (D) ReLU and (E) hardtanh non-linear activations on  $s(\cdot) = \sin(\cdot)$  (A).

# Implicit bias of convolutional and fully connected networks

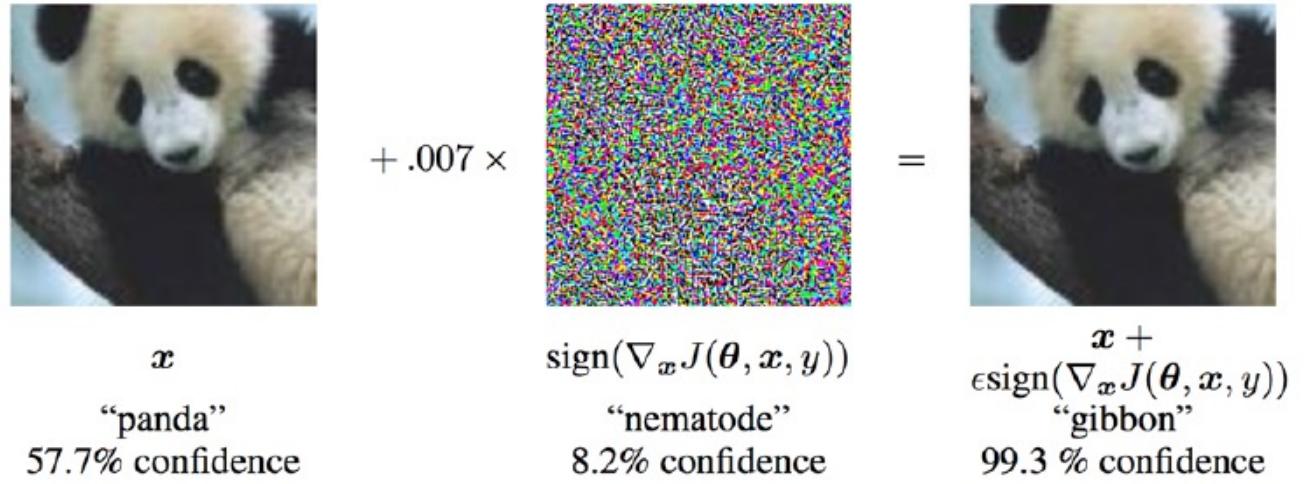


(a) Fully connected network of depth  $L$   
 $\overline{\beta}^\infty \propto \underset{\forall n, y_n \langle \mathbf{x}_n, \beta \rangle \geq 1}{\operatorname{argmin}} \|\beta\|_2$  (independent of  $L$ )



(b) Convolutional network of depth  $L$   
 $\overline{\beta}^\infty \propto \text{first order stationary point of } \underset{\forall n, y_n \langle \mathbf{x}_n, \beta \rangle \geq 1}{\operatorname{argmin}} \|\widehat{\beta}\|_{2/L}$

# Adversarial attacks



**Definition 1** (Adversarial Attack). Let  $x_0 \in \mathbb{R}^d$  be a data point belong to class  $\mathcal{C}_i$ . Define a target class  $\mathcal{C}_t$ . An **adversarial attack** is a mapping  $\mathcal{A} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  such that the perturbed data

$$x = \mathcal{A}(x_0)$$

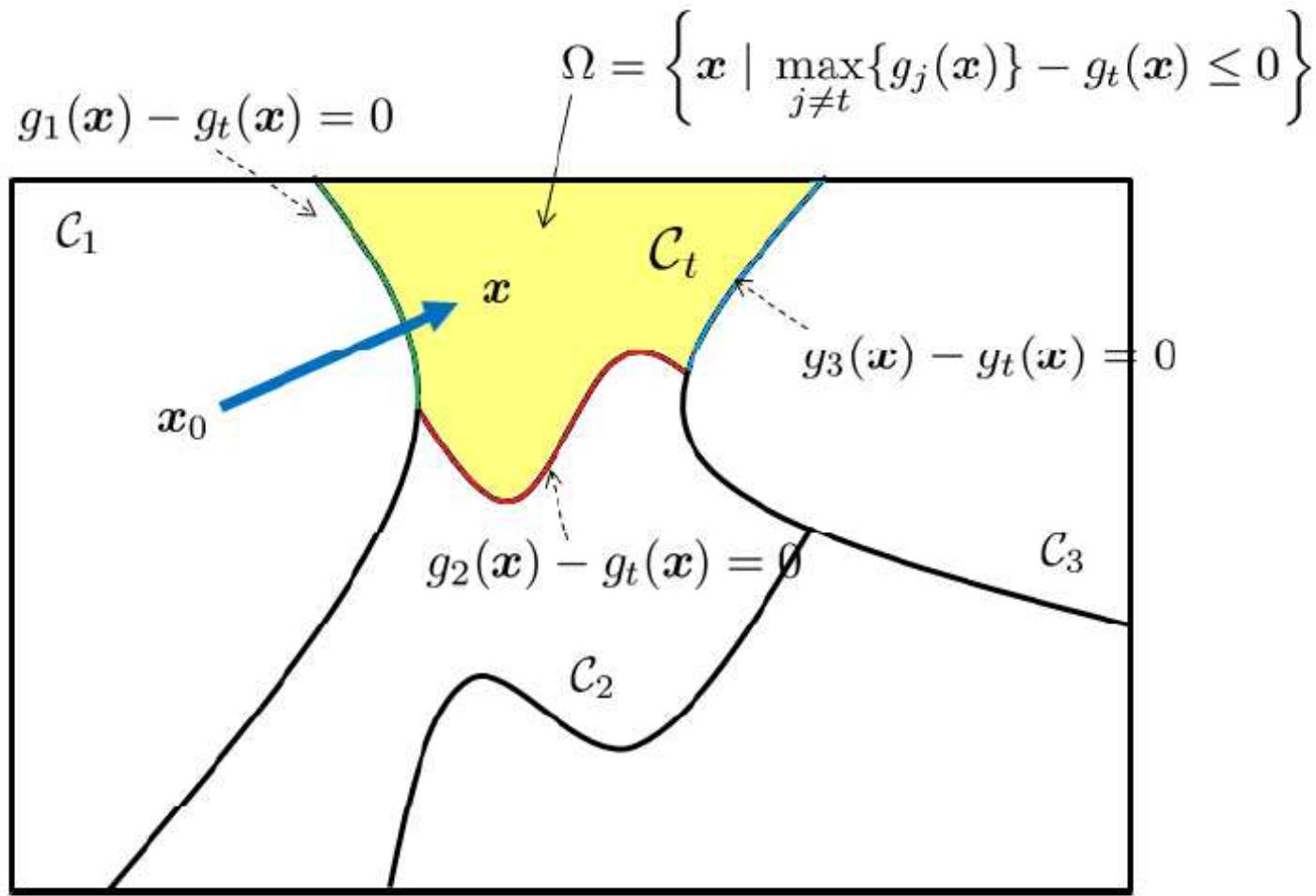
is misclassified as  $\mathcal{C}_t$ .

**Definition 2** (Additive Adversarial Attack). Let  $x_0 \in \mathbb{R}^d$  be a data point belong to class  $\mathcal{C}_i$ . Define a target class  $\mathcal{C}_t$ . An **additive** adversarial attack is an addition of a perturbation  $r \in \mathbb{R}^d$  such that the perturbed data

$$x = x_0 + r$$

is misclassified as  $\mathcal{C}_t$ .

# Geometry



$$\Omega = \left\{ \mathbf{x} \mid \begin{array}{l} g_1(\mathbf{x}) - g_t(\mathbf{x}) \leq 0 \\ g_2(\mathbf{x}) - g_t(\mathbf{x}) \leq 0 \\ \vdots \\ g_k(\mathbf{x}) - g_t(\mathbf{x}) \leq 0 \end{array} \right\}$$

**Definition 3** (Minimum Norm Attack). *The **minimum norm attack** finds a perturbed data  $\mathbf{x}$  by solving the optimization*

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \|\mathbf{x} - \mathbf{x}_0\| \\ & \text{subject to} && \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0, \end{aligned} \tag{3.3}$$

where  $\|\cdot\|$  can be any norm specified by the user.

**Definition 4** (Maximum Allowable Attack). *The **maximum allowable attack** finds a perturbed data  $\mathbf{x}$  by solving the optimization*

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \\ & \text{subject to} && \|\mathbf{x} - \mathbf{x}_0\| \leq \eta, \end{aligned} \tag{3.4}$$

where  $\|\cdot\|$  can be any norm specified by the user, and  $\eta > 0$  denotes the magnitude of the attack.

**Definition 5** (Regularization-based Attack). *The **regularization-based attack** finds a perturbed data  $\mathbf{x}$  by solving the optimization*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\| + \lambda (\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x})) \tag{3.5}$$

where  $\|\cdot\|$  can be any norm specified by the user, and  $\lambda > 0$  is a regularization parameter.

# Attacks

# Linear attacks

Since we want the perturbed data  $\mathbf{x}$  to live in  $\mathcal{C}_t$ , the constraint set becomes

$$\begin{bmatrix} \mathbf{w}_1^T - \mathbf{w}_t^T \\ \vdots \\ \mathbf{w}_{t-1}^T - \mathbf{w}_t^T \\ \mathbf{w}_{t+1}^T - \mathbf{w}_t^T \\ \vdots \\ \mathbf{w}_k^T - \mathbf{w}_t^T \end{bmatrix} \mathbf{x} + \begin{bmatrix} w_{1,0} - w_{t,0} \\ \vdots \\ w_{t-1,0} - w_{t,0} \\ w_{t+1,0} - w_{t,0} \\ \vdots \\ w_{k,0} - w_{t,0} \end{bmatrix} \leq \mathbf{0} \Leftrightarrow \mathbf{A}^T \mathbf{x} \leq \mathbf{b} \quad (3.11)$$

where  $\mathbf{A} = [\mathbf{w}_1 - \mathbf{w}_t, \dots, \mathbf{w}_k - \mathbf{w}_t] \in \mathbb{R}^{d \times (k-1)}$ , and  $\mathbf{b} = [w_{t,0} - w_{1,0}, \dots, w_{t,0} - w_{k,0}]^T$ . This is summarized in the following Lemma.

**Lemma 1** (Constraint Set of Linear Classifier). *Consider a  $k$ -class linear classifier with discriminant function  $g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i,0}$  for  $i = 1, \dots, k$ . Define*

$$\mathbf{A} = [\mathbf{w}_1 - \mathbf{w}_t, \dots, \mathbf{w}_k - \mathbf{w}_t] \in \mathbb{R}^{d \times (k-1)}, \quad \text{and} \quad \mathbf{b} = [w_{t,0} - w_{1,0}, \dots, w_{t,0} - w_{k,0}]^T \in \mathbb{R}^{k-1},$$

*Then, the constraint set is*

$$\Omega = \{\mathbf{x} \mid \mathbf{A}^T \mathbf{x} \leq \mathbf{b}\}. \quad (3.12)$$

To gain insights about these constraints, we consider a  $k$ -class linear classifier. In this case, each discriminant function takes the form

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i,0}. \quad (3.9)$$

The decision boundary between the  $i$ -th class and the  $t$ -th class is therefore

$$g(\mathbf{x}) = (\mathbf{w}_i - \mathbf{w}_t)^T \mathbf{x} + w_{i,0} - w_{t,0} = 0. \quad (3.10)$$

$$\underset{x}{\text{minimize}} \quad \|x - x_0\|^2 \quad \text{subject to} \quad A^T x \leq b$$

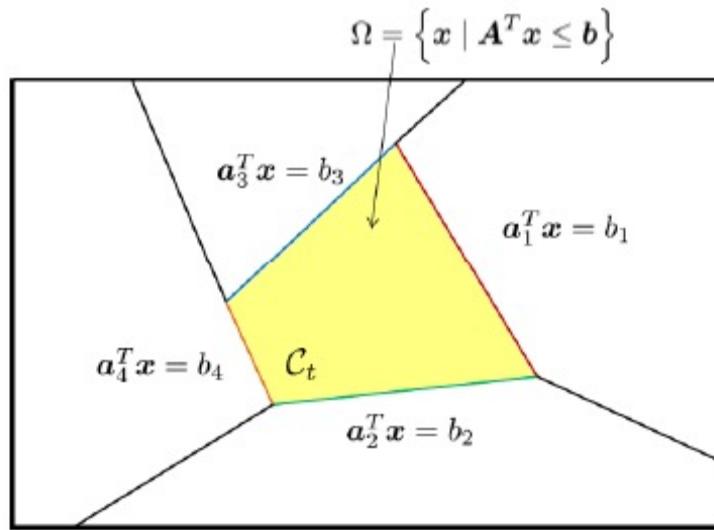


Figure 3.4: Geometry of the constraint set  $\Omega$  for a linear classifier. The constraint set  $\Omega$  is now a polygon with decision boundaries defined by  $a_i^T x = b_i$ , where  $a_i = w_i - w_t$  and  $b_i = w_{i0} - w_{t0}$ .

# Linear attacks

# Minimum norm attack

---

## A. Minimum-Norm Attack

In order to gain insight of the adversarial attack, it is useful to consider a simple linear classifier with only two classes. A linear classifier has a discriminant function

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}, \quad (3.22)$$

and therefore by defining  $\mathbf{w} \stackrel{\text{def}}{=} \mathbf{w}_i - \mathbf{w}_t$  and  $w_0 \stackrel{\text{def}}{=} w_{i0} - w_{t0}$ , we can show that

$$g_i(\mathbf{x}) - g_t(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0. \quad (3.23)$$

Recalling Equation (3.3), the adversarial attack can be formulated as

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \|\mathbf{x} - \mathbf{x}_0\|^2 \\ & \text{subject to} && g_i(\mathbf{x}) - g_t(\mathbf{x}) = 0. \end{aligned} \quad (3.24)$$

In this optimization, we simplify  $\max_{j \neq t} \{g_j(\mathbf{x})\}$  to  $g_i(\mathbf{x})$  because there are only two classes in our problem. If the index is not  $t$ , then it has to be  $i$ . We also replaced the inequality  $g_i(\mathbf{x}) - g_t(\mathbf{x}) \leq 0$  as an equality, because the projection theorem in Equation (3.17) suggests that the solution must be on the decision boundary.

**Theorem 2** (Minimum  $\ell_2$  Norm Attack for Two-Class Linear Classifier). *The adversarial attack to a two-class linear classifier is the solution of*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 \quad \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0, \quad (3.25)$$

which is given by

$$\mathbf{x}^* = \mathbf{x}_0 - \left( \frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_2} \right) \frac{\mathbf{w}}{\|\mathbf{w}\|_2}. \quad (3.26)$$

*Proof.* The Lagrange multiplier of the constrained optimization is given by

$$\mathcal{L}(\mathbf{x}, \lambda) = \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0).$$

The solution of the optimization is the saddle point  $(\mathbf{x}^*, \lambda^*)$  such that  $\nabla_{\mathbf{x}} \mathcal{L} = 0$  and  $\nabla_{\lambda} \mathcal{L} = 0$ .

Taking derivative with respect to  $\mathbf{x}$  and  $\lambda$  yields

$$\nabla_{\mathbf{x}} \mathcal{L} = \mathbf{x} - \mathbf{x}_0 + \lambda \mathbf{w} = 0,$$

$$\nabla_{\lambda} \mathcal{L} = \mathbf{w}^T \mathbf{x} + w_0 = 0.$$

# Regularization norm attack

## C. Regularization-based Attack

In both minimum norm attack and maximum allowable attack, the optimizations are constrained. For certain simple cases, e.g., binary linear classifiers, closed-form solutions can be derived. However, for advanced classifiers such as deep neural networks, solving an optimization involving constraints are typically very difficult. The regularization-based attack alleviates the problem by considering

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda \left( \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \right).$$

For two-class linear classifier, this is simplified to

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0).$$

Since the regularization-based attack is an unconstrained version of the minimum norm attack and the maximum allowable attack, one should expect to have a very similar solution.

**Theorem 6** (Regularization-based Attack for Two-Class Linear Classifier). *The regularization-based attack for a two-class linear classifier generates the attack by solving*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0), \quad (3.46)$$

*of which the solution is given by*

$$\mathbf{x} = \mathbf{x}_0 - \lambda \mathbf{w}. \quad (3.47)$$

*Proof.* Let  $\varphi(\mathbf{x}) = \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0)$ . Taking the first order derivative and sending to zero yields

$$0 = \nabla \varphi(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_0) + \lambda \mathbf{w}. \quad (3.48)$$

Therefore, we have  $\mathbf{x} = \mathbf{x}_0 - \lambda \mathbf{w}$ . □

**Definition 6** (DeepFool Attack by Moosavi-Dezfooli et al. 2016). *The DeepFool attack for a two-class classification generates the attack by solving the optimization*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 \quad \text{subject to} \quad g(\mathbf{x}) = 0, \quad (3.29)$$

where  $g(\mathbf{x}) = 0$  is the nonlinear decision boundary separating the two classes.

Due to the nonlinearity of  $g(\mathbf{x})$ , it is generally very difficult to derive a closed-form solution. The numerical procedure to compute the solution can be derived by iteratively updating  $\mathbf{x}$ , where each iteration minimizes  $\mathbf{x}$  over the first order approximation of  $g(\mathbf{x})$ :

$$g(\mathbf{x}) \approx g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T(\mathbf{x} - \mathbf{x}^{(k)}),$$

where  $\mathbf{x}^{(k)}$  is the  $k$ -th iterate of the solution. In other words, we are defining a procedure

$$\mathbf{x}^{(k+1)} = \underset{\mathbf{x}}{\text{argmin}} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|^2 \quad \text{subject to} \quad g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T(\mathbf{x} - \mathbf{x}^{(k)}) = 0. \quad (3.30)$$

By identifying  $\mathbf{w}^{(k)} = \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})$  and  $w_0^{(k)} = g(\mathbf{x}^{(k)}) - \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T\mathbf{x}^{(k)}$ , the linearized problem Equation (3.30) becomes

$$\mathbf{x}^{(k+1)} = \underset{\mathbf{x}}{\text{argmin}} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|^2 \quad \text{subject to} \quad (\mathbf{w}^{(k)})^T\mathbf{x} + w_0^{(k)} = 0$$

The solution can thus be found by using Equation (3.26), yielding

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \left( \frac{(\mathbf{w}^{(k)})^T\mathbf{x}^{(k)} + w_0^{(k)}}{\|\mathbf{w}^{(k)}\|^2} \right) \mathbf{w}^{(k)} \\ &= \mathbf{x}^{(k)} - \left( \frac{g(\mathbf{x}^{(k)})}{\|\nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})\|^2} \right) \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)}). \end{aligned}$$

This result is summarized in the following corollary.

**Corollary 1** (DeepFool Algorithm for Two-Class Problem). *An iterative procedure to obtain the DeepFool attack solution is*

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \underset{\mathbf{x}}{\text{argmin}} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|^2 \quad \text{subject to} \quad g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T(\mathbf{x} - \mathbf{x}^{(k)}) = 0 \\ &= \mathbf{x}^{(k)} - \left( \frac{g(\mathbf{x}^{(k)})}{\|\nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})\|^2} \right) \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)}). \end{aligned} \quad (3.31)$$

# Non-linear attacks

# Can we attack with random noise?

## 3.4 Can Random Noise Attack?

Adversarial attack works because the perturbation is carefully designed. What if we perturb the data by pure i.i.d. Gaussian noise?

Recall that when launching attacks for a linear classifier, the perturbation always takes the form of  $\mathbf{x} = \mathbf{x}_0 + \lambda \mathbf{w}$ . That is, we want to move  $\mathbf{x}_0$  along  $\mathbf{w}$  by an amount  $\lambda$  so that  $\mathbf{x}$  is misclassified. Now, if we do not move along  $\mathbf{w}$  but along a random vector  $\mathbf{r}$  such that

$$\mathbf{x} = \mathbf{x}_0 + \sigma_r \mathbf{r}, \quad (3.57)$$

where  $\mathbf{r} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , then we can ask if we can still misclassify the data point  $\mathbf{x}$ . Clearly, this requires us to check whether  $\mathbf{w}^T \mathbf{r} > 0$ . If  $\mathbf{w}^T \mathbf{r} > 0$ , then  $\mathbf{r}$  and  $\mathbf{w}$  will form an acute angle and so for sufficient step size we will be able to move  $\mathbf{x}_0$  to another class. If  $\mathbf{w}^T \mathbf{r} < 0$ , then  $\mathbf{w}$  and  $\mathbf{r}$  are form an obtuse angle and so  $\mathbf{r}$  will move  $\mathbf{x}_0$  to an opposite direction.

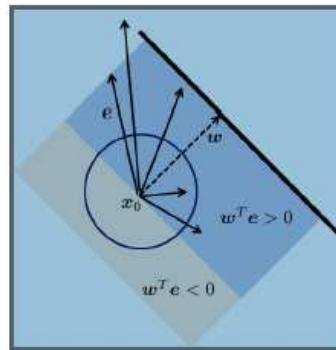


Figure 3.10: Attacking the linear classifier with i.i.d. noise is equivalent to putting an uncertainty circle around  $\mathbf{x}_0$  with radius  $\sigma_r$ . The possible attack directions are those that form acute angle with the normal vector  $\mathbf{w}$ . Therefore, among all the possible  $\mathbf{r}$ 's, we are only interested in those such that  $\mathbf{w}^T \mathbf{r} > 0$ , which occupy half of the space in 2D.

Figure 3.10 provides a pictorial illustration of the noise attack. From this figure, we are tempted to think that i.i.d. noise can achieve 50% attack rate because  $\mathbf{w}^T \mathbf{r} > 0$  occupies half of the space. The problem of such argument is that in high dimension, the probability of  $\mathbf{w}^T \mathbf{r} > 0$  is diminishing very quickly as the dimensionality of  $\mathbf{r}$  grows. This is a well-known

# Can we attack with random noise?

phenomenon called **curse of dimensionality**. To illustrate the idea, let us evaluate the probability of  $\mathbf{w}^T \mathbf{r} \geq \epsilon$  for some  $\epsilon > 0$ . To this end, let us consider

$$\mathbb{P}\left[\frac{1}{d}\mathbf{w}^T \mathbf{r} \geq \epsilon\right] = \mathbb{P}\left[\frac{1}{d}\sum_{j=1}^d w_j r_j \geq \epsilon\right],$$

where  $d$  is the dimensionality of  $\mathbf{w}$ , i.e.,  $\mathbf{w} \in \mathbb{R}^d$ . The tolerance level  $\epsilon$  is a small positive constant that stays away from 0.

**Example.** Before we discuss the theoretical results, it will be useful to do a quick numerical experiment. Consider a special case where  $\mathbf{w} = \mathbf{1}_{d \times 1}$ , i.e., a  $d$ -dimensional all-one vector, and  $\mathbf{r} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . In this case, we define the average as

$$Y \stackrel{\text{def}}{=} \frac{1}{d} \sum_{j=1}^d r_j. \quad (3.58)$$

It follows that  $Y$  is a Gaussian random variable because linear combination of Gaussian remains a Gaussian. The mean and variance are

$$\mathbb{E}[Y] = 0, \quad \text{and} \quad \text{Var}[Y] = \frac{1}{d}. \quad (3.59)$$

Therefore, the probability of the event  $\{Y > \epsilon\}$  is

$$\begin{aligned} \mathbb{P}[Y > \epsilon] &= \int_{\epsilon}^{\infty} \frac{1}{\sqrt{2\pi/d}} \exp\left\{-\frac{t^2}{2/d}\right\} dt \\ &= \int_{\epsilon\sqrt{\frac{d}{2}}}^{\infty} \frac{1}{\sqrt{\pi}} \exp\{-t^2\} dt \\ &= \frac{1}{2} \text{erfc}\left(\epsilon\sqrt{d/2}\right), \end{aligned} \quad (3.60)$$

where erfc is the complementary error function defined as  $\text{erfc}(z) \stackrel{\text{def}}{=} \frac{2}{\sqrt{\pi}} \int_z^{\infty} e^{-t^2} dt$ . As we can see in Figure 3.11, the probability drops rapidly as  $d$  increases.

Let us discuss the general case. One classical result to derive the probability is the tail bound on Gaussian by Feller 1968.<sup>1</sup> If  $Y$  is a Gaussian random variable with  $Y \sim \mathcal{N}(\mu, \sigma)$ , then

$$\mathbb{P}[Y \geq \mu + \sigma\epsilon] \leq \frac{1}{\epsilon} \frac{e^{-\epsilon^2/2}}{\sqrt{2\pi}}. \quad (3.61)$$

# Can we attack with random noise?

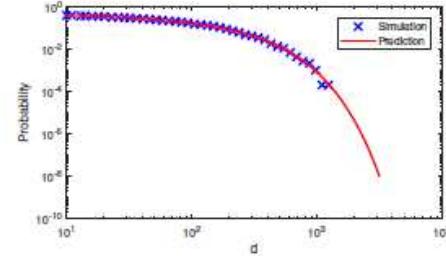


Figure 3.11: Empirical probability and the theoretically predicted value. In this experiment, we let  $Y = \frac{1}{d} \sum_{j=1}^d r_j$  where  $r_j \sim \mathcal{N}(0, 1)$  is an iid Gaussian random variable. We plot the probability  $\mathbb{P}[Y > \epsilon]$  as a function of  $d$ . Note the tightness of the prediction, and the rapid decaying behavior of the probability. The result shows that when  $d$  grows, the value  $Y$  is concentrated at 0.

Now, let  $Y = \frac{1}{d} \sum_{j=1}^d w_j r_j$ . Since a linear combination of Gaussian remains a Gaussian, it holds that  $Y$  is Gaussian and

$$\mu = \mathbb{E}[Y] = 0, \quad \text{and} \quad \sigma^2 = \text{Var}[Y] = \frac{1}{d^2} \sum_{j=1}^d w_j^2 = \frac{\|\mathbf{w}\|^2}{d^2}. \quad (3.62)$$

Therefore, by substituting  $\varepsilon = \sigma\epsilon$  we can show that

$$\mathbb{P}\left[\frac{1}{d} \sum_{j=1}^d w_j r_j \geq \varepsilon\right] = \mathbb{P}[Y \geq \varepsilon] \leq \frac{\sigma e^{-\frac{\varepsilon^2}{2\sigma^2}}}{\varepsilon \sqrt{2\pi}} = \frac{\|\mathbf{w}\|}{\varepsilon d \sqrt{2\pi}} \exp\left\{-d^2 \frac{\varepsilon^2}{2\|\mathbf{w}\|^2}\right\}. \quad (3.63)$$

As  $d \rightarrow \infty$ , it holds that  $\mathbb{P}\left[\frac{1}{d} \sum_{j=1}^d w_j r_j \geq \varepsilon\right] \rightarrow 0$ . That means, the probability of getting a “good attack direction” is diminishing to zero exponentially. Putting everything together we have the following theorem.

**Theorem 7.** Let  $\mathbf{w}^T \mathbf{x} + w_0 = 0$  be the decision boundary of a linear classifier, and let  $\mathbf{x}_0 \in \mathbb{R}^d$  be an input data point. Suppose we attack the classifier by adding i.i.d. Gaussian noise  $\mathbf{r} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  to  $\mathbf{x}_0$ . The probability of a successful attack at a tolerance level  $\varepsilon$  is  $\mathbb{P}\left[\frac{1}{d} \sum_{j=1}^d w_j r_j \geq \varepsilon\right]$ , and such probability is upper bounded by

$$\mathbb{P}\left[\frac{1}{d} \sum_{j=1}^d w_j r_j \geq \varepsilon\right] \leq \frac{\|\mathbf{w}\|}{\varepsilon d \sqrt{2\pi}} \exp\left\{-d^2 \frac{\varepsilon^2}{2\|\mathbf{w}\|^2}\right\}. \quad (3.64)$$

Therefore, as  $d \rightarrow \infty$  it becomes increasingly more difficult for i.i.d. Gaussian noise to succeed in attacking.

# Equivalence robust classifier/maximum margin

**Definition 2.2** (Dual norm). Let  $\|\cdot\|$  be a norm on  $\mathbb{R}^n$ . The associated *dual norm*, denoted  $\|\cdot\|_*$ , is defined as  $\|\delta\|_* = \sup_{\mathbf{x}} \{ |\langle \delta, \mathbf{x} \rangle| \mid \|\mathbf{x}\| \leq 1 \}$ .

**Definition 2.3** (Linear Separability). We say a dataset is linearly separable if there exists  $\mathbf{w}, b$  such that  $y_i(\mathbf{w}^\top \mathbf{x}_i + b) > 0$  for all  $i$ .

**Lemma 2.1** (Maximally Robust Linear Classifier (Ben-Tal et al. [14], §12)). *For linear models and linearly separable data, the following problems are equivalent; i.e., from a solution of one, a solution of the other is readily found.*

$$\text{Maximally robust classifier: } \arg \max_{\mathbf{w}, b} \{\varepsilon \mid y_i(\mathbf{w}^\top (\mathbf{x}_i + \delta) + b) > 0, \forall i, \|\delta\| \leq \varepsilon\}, \quad (4)$$

$$\text{Maximum margin classifier: } \arg \max_{\mathbf{w}, b: \|\mathbf{w}\|_* \leq 1} \{\varepsilon \mid y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq \varepsilon, \forall i\}, \quad (5)$$

$$\text{Minimum norm classifier: } \arg \min_{\mathbf{w}, b} \{\|\mathbf{w}\|_* \mid y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \forall i\}. \quad (6)$$

The expression  $\min_i y_i(\mathbf{w}^\top \mathbf{x}_i + b)/\|\mathbf{w}\|$  is the margin of a classifier  $\mathbf{w}$  that is the distance of the nearest training point to the classification boundary, i.e. the line  $\{v : \mathbf{w}^\top v = -b\}$ .

*Proof.* We first show that the maximally robust classifier is equivalent to a robust counterpart by removing  $\delta$  from the problem,

$$\begin{aligned} & \arg \max_{\mathbf{w}, b} \{\varepsilon \mid y_i(\mathbf{w}^\top (\mathbf{x}_i + \delta) + b) > 0, \forall i, \|\delta\| \leq \varepsilon\} \\ & \quad (\text{homogeneity of } p\text{-norm}) \\ &= \arg \max_{\mathbf{w}, b} \{\varepsilon \mid y_i(\mathbf{w}^\top (\mathbf{x}_i + \varepsilon \delta) + b) > 0, \forall i, \|\delta\| \leq 1\} \\ & \quad (\text{if it is true for all } \delta \text{ it is true for the worst of them}) \\ &= \arg \max_{\mathbf{w}, b} \{\varepsilon \mid \inf_{\|\delta\| \leq 1} y_i(\mathbf{w}^\top (\mathbf{x}_i + \varepsilon \delta) + b) > 0, \forall i\} \\ &= \arg \max_{\mathbf{w}, b} \{\varepsilon \mid y_i(\mathbf{w}^\top \mathbf{x}_i + b) + \varepsilon \inf_{\|\delta\| \leq 1} \mathbf{w}^\top \delta > 0, \forall i\} \\ & \quad (\text{definition of dual norm}) \\ &= \arg \max_{\mathbf{w}, b} \{\varepsilon \mid y_i(\mathbf{w}^\top \mathbf{x}_i + b) > \varepsilon \|\mathbf{w}\|_*, \forall i\} \end{aligned}$$

Assuming  $\mathbf{w} \neq 0$ , which is a result of linear separability assumption, we can divide both sides by  $\|\mathbf{w}\|_*$  and change variables,

$$= \arg \max_{\mathbf{w}, b} \{\varepsilon \mid y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq \varepsilon, \forall i, \|\mathbf{w}\|_* \leq 1\},$$

where we are also allowed to change  $>$  to  $\geq$  because any solution to one problem gives an equivalent solution to the other given  $\mathbf{w} \neq 0$ .

# Proof continues

Now we show that the robust counterpart is equivalent to the minimum norm classification problem by removing  $\varepsilon$ . When the data is linearly separable there exists a solution with  $\varepsilon > 0$ ,

$$\begin{aligned} & \arg \max_{\mathbf{w}, b} \{ \varepsilon \mid y_i(\mathbf{w}^\top \mathbf{x}_i + b) > \varepsilon \|\mathbf{w}\|_*, \forall i \} \\ &= \arg \max_{\mathbf{w}, b} \left\{ \varepsilon \mid y_i \left( \frac{\mathbf{w}^\top}{\varepsilon \|\mathbf{w}\|_*} \mathbf{x}_i + \frac{b}{\varepsilon \|\mathbf{w}\|_*} \right) \geq 1, \forall i \right\} \end{aligned}$$

This problem is invariant to any non-zero scaling of  $(\mathbf{w}, b)$ , so with no loss of generality we set  $\|\mathbf{w}\|_* = 1$ .

$$= \arg \max_{\mathbf{w}, b} \left\{ \varepsilon \mid y_i \left( \frac{\mathbf{w}^\top}{\varepsilon} \mathbf{x}_i + b \right) \geq 1, \forall i, \|\mathbf{w}\|_* = 1 \right\}$$

Let  $\mathbf{w}' = \mathbf{w}/\varepsilon$ , then the solution to the following problem gives a solution for  $\mathbf{w}$ ,

$$\begin{aligned} & \arg \max_{\mathbf{w}', b} \left\{ \frac{1}{\|\mathbf{w}'\|_*} \mid y_i(\mathbf{w}'^\top \mathbf{x}_i + b) \geq 1, \forall i \right\} \\ &= \arg \min_{\mathbf{w}', b} \{ \|\mathbf{w}'\|_* \mid y_i(\mathbf{w}'^\top \mathbf{x}_i + b) \geq 1, \forall i \}. \end{aligned}$$

# Robustness and implicit Bias

**Theorem 3.1** (Implicit Bias of Steepest Descent (Gunasekar et al. [12] (Theorem 5))). *For any separable dataset  $\{x_i, y_i\}$  and any norm  $\|\cdot\|$ , consider the steepest descent updates from (7) for minimizing the empirical risk  $\mathcal{L}(\mathbf{w})$  (defined in Section 2) with the exponential loss,  $\zeta(z) = \exp(-z)$ . For all initializations  $\mathbf{w}_0$ , and all bounded step-sizes satisfying a known upper bound, the iterates  $\mathbf{w}_t$  satisfy*

$$\lim_{t \rightarrow \infty} \min_i \frac{y_i \mathbf{w}_t^\top \mathbf{x}_i}{\|\mathbf{w}_t\|} = \max_{\mathbf{w}: \|\mathbf{w}\| \leq 1} \min_i y_i \mathbf{w}^\top \mathbf{x}_i. \quad (8)$$

*In particular, if a unique maximum margin classifier  $\mathbf{w}_{\|\cdot\|}^* = \arg \max_{\mathbf{w}: \|\mathbf{w}\| \leq 1} \min_i y_i \mathbf{w}^\top \mathbf{x}_i$  exists, the limit direction converges to it, i.e.  $\lim_{t \rightarrow \infty} \frac{\mathbf{w}_t}{\|\mathbf{w}_t\|} = \mathbf{w}_{\|\cdot\|}^*$ .*

# Robustness and implicit Bias

**Corollary 1** (Implicit Robustness of Steepest Descent). *For any linearly separable dataset and any norm  $\|\cdot\|$ , steepest descent iterates minimizing the empirical risk,  $\mathcal{L}(\mathbf{w})$ , satisfying the conditions of Theorem 3.1, converge in direction to a maximally robust classifier,*

$$\arg \max_{\mathbf{w}} \{\varepsilon \mid y_i \mathbf{w}^\top (\mathbf{x}_i + \boldsymbol{\delta}) > 0, \forall i, \|\boldsymbol{\delta}\|_* \leq \varepsilon\}.$$

*In particular, a maximally robust classifier against  $\ell_1$ ,  $\ell_2$ , and  $\ell_\infty$  is reached, respectively, by sign gradient descent, gradient descent, and coordinate descent.*

*Proof.* By Theorem 3.1, the margin of the steepest descent iterates,  $\min_i \frac{y_i \mathbf{w}_t^\top \mathbf{x}_i}{\|\mathbf{w}_t\|}$ , converges as  $t \rightarrow \infty$  to the maximum margin,  $\max_{\mathbf{w}: \|\mathbf{w}\| \leq 1} \min_i y_i \mathbf{w}^\top \mathbf{x}_i$ . By Lemma 2.1, any maximum margin classifier w.r.t.  $\|\cdot\|$  gives a maximally robust classifier w.r.t.  $\|\cdot\|_*$ .  $\square$

# IB Fourier and maximal robustness

**Definition 3.3** (Discrete Fourier Transform).  $\mathcal{F}(\mathbf{w}) \in \mathbb{C}^d$  denotes the Fourier coefficients of  $\mathbf{w}$  where  $[\mathcal{F}(\mathbf{w})]_d = \frac{1}{\sqrt{d}} \sum_{k=0}^{d-1} [\mathbf{w}]_k \exp(-\frac{2\pi j}{d} kd)$  and  $j^2 = -1$ .

**Theorem 3.2** (Implicit Bias towards Fourier Sparsity (Gunasekar et al. [22], Theorem 2, 2.a)). *Consider the family of  $L$ -layer linear convolutional networks and the sequence of gradient descent iterates,  $\mathbf{w}_t$ , minimizing the empirical risk,  $\mathcal{L}(\mathbf{w})$ , with the exponential loss,  $\exp(-z)$ . For almost all linearly separable datasets under known conditions on the step size and convergence of iterates,  $\mathbf{w}_t$  converges in direction to the classifier minimizing the norm of the Fourier coefficients given by*

$$\arg \min_{\mathbf{w}_1, \dots, \mathbf{w}_L} \{ \|\mathcal{F}(\mathbf{w})\|_{2/L} \mid y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1, \forall i \}. \quad (9)$$

*In particular, for two-layer linear convolutional networks the implicit bias is towards the solution with minimum  $\ell_1$  norm of the Fourier coefficients,  $\|\mathcal{F}(\mathbf{w})\|_1$ . For  $L > 2$ , the convergence is to a first-order stationary point.*

# IB Fourier and maximal robustness

**Corollary 2** (Maximally Robust to Perturbations with Bounded Fourier Coefficients). *Consider the family of two-layer linear convolutional networks and the gradient descent iterates,  $w_t$ , minimizing the empirical risk. For almost all linearly separable datasets under conditions of Theorem 3.2,  $w_t$  converges in direction to a maximally robust classifier,*

$$\arg \max_{\mathbf{w}_1, \dots, \mathbf{w}_L} \{\varepsilon \mid y_i \varphi_{conv}(\mathbf{x}_i + \delta; \{\mathbf{w}_l\}_{l=1}^L) > 0, \forall i, \|\mathcal{F}(\delta)\|_\infty \leq \varepsilon\}.$$

# Fourier, proof

## B.2 Proof of Maximally Robust to Perturbations Bounded in Fourier Domain (Corollary 2)

The proof mostly follows from the equivalence for linear models in Appendix B.1 by substituting the dual norm of Fourier- $\ell_1$ . Here,  $\mathbf{A}^*$  denotes the complex conjugate transpose,  $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^\top \mathbf{v}^*$  is the complex inner product,  $[\mathbf{F}]_{ik} = \frac{1}{\sqrt{D}} \omega_D^{ik}$  the DFT matrix where  $\omega_D = e^{-j2\pi/D}$ ,  $j = \sqrt{-1}$ .

Let  $\|\cdot\|$  be a norm on  $\mathbb{C}^n$  and  $\langle \cdot, \cdot \rangle$  be the complex inner product. Similar to  $\mathbb{R}^n$ , the associated dual norm is defined as  $\|\delta\|_* = \sup_{\mathbf{x}} \{ |\langle \delta, \mathbf{x} \rangle| \mid \|\mathbf{x}\| \leq 1 \}$ .

$$\begin{aligned} & \|\mathcal{F}(w)\|_1 \\ &= \sup_{\|\delta\|_\infty \leq 1} |\langle \mathcal{F}(w), \delta \rangle| \\ &\quad (\text{Expressing DFT as a linear transformation.}) \\ &= \sup_{\|\delta\|_\infty \leq 1} |\langle Fw, \delta \rangle| \\ &= \sup_{\|\delta\|_\infty \leq 1} |\langle w, F^* \delta \rangle| \\ &\quad (\text{Change of variables and } F^{-1} = F^*.) \\ &= \sup_{\|F\delta\|_\infty \leq 1} |\langle w, \delta \rangle| \\ &= \sup_{\|\mathcal{F}(\delta)\|_\infty \leq 1} |\langle w, \delta \rangle|. \end{aligned}$$

# An aspect of the implicit bias: important Fourier features. A step towards testing for interpretability

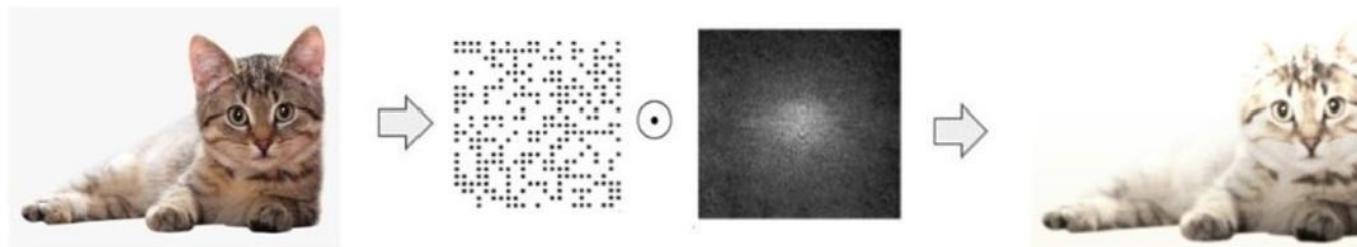
## Question

Which Fourier features are essential to the network to solve the task?

$x$

$M_\Phi \odot \mathcal{F}x$

$\mathcal{F}^{-1}(M_\Phi \odot \mathcal{F}x)$



$$M_\Phi(\lambda, p) = \operatorname{argmin}_{M_\Phi} \sum_{x \in \mathcal{X}_V} e^{[\mathcal{L}(\Phi(\bar{x}), y) - \mathcal{L}(\Phi(x), y)]^2} + \lambda \|M_\Phi\|_p, \quad \lambda \in \mathbb{R}_+,$$

**Result:**  $M_\Phi$  highlights the important Fourier features.

The study hasn't been done in the context of neuroscience but the technique can be applied.

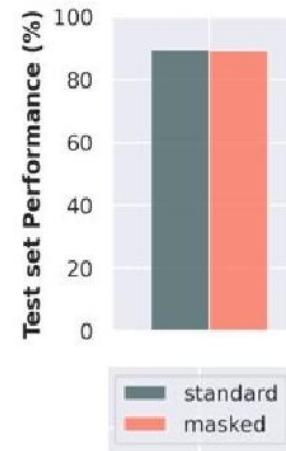
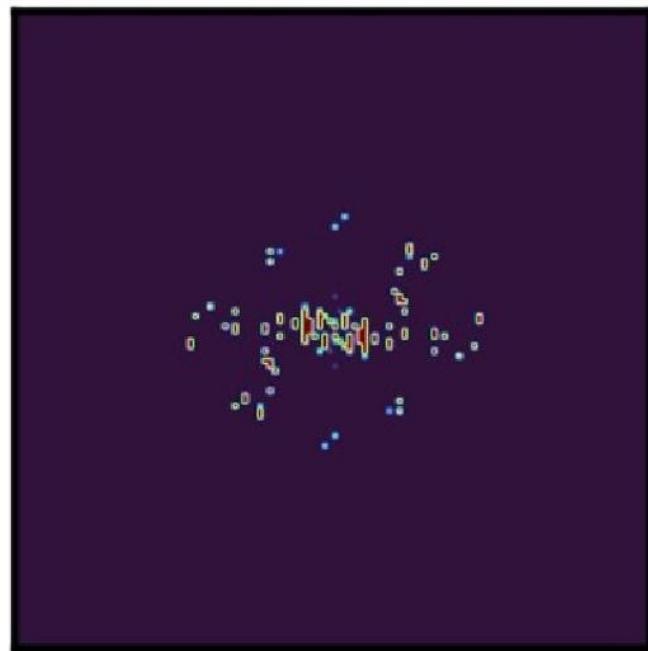
# Algorithm

- Train the model
- Freeze the weights and add a mask layer
- Train the mask to
  - ① preserve the score
  - ② maximize the frequency sparsity

## Why Fourier? Many successful stories

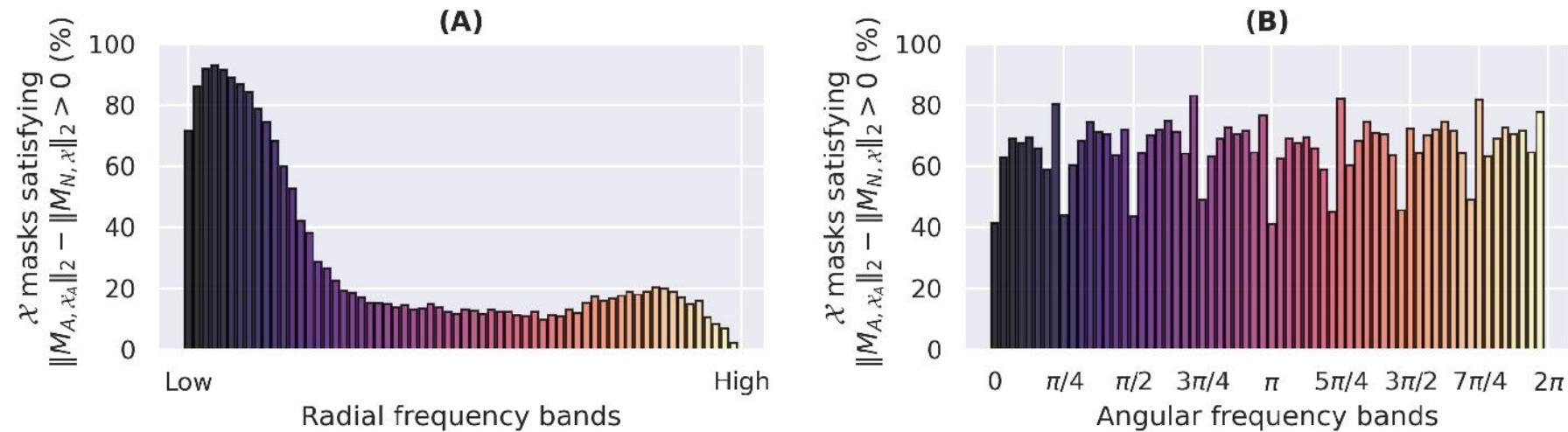
- Abello et al. (2021). Dissecting the high-frequency bias in convolutional neural In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.
- Caro et al. (2020). Local convolutions cause an implicit bias towards high frequency adversarial examples.
- Christian et al. (2021). Ringing relus: Harmonic distortion analysis of nonlinear feedforward networks. ICLR
- Geirhos, R. et al. Imagenet trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. ICLR. Li, Z. et al. (2022). Robust deep learning object recognition models rely on low frequency information in natural images.
- Ortiz-Jimenez, G. at al. (2020). Neural anisotropy directions.
- Sharma, Y., Ding, G. W., and Brubaker, M. A. (2019). On the Effectiveness of Low Frequency Perturbations. IJCAI

# Single image Fourier masks



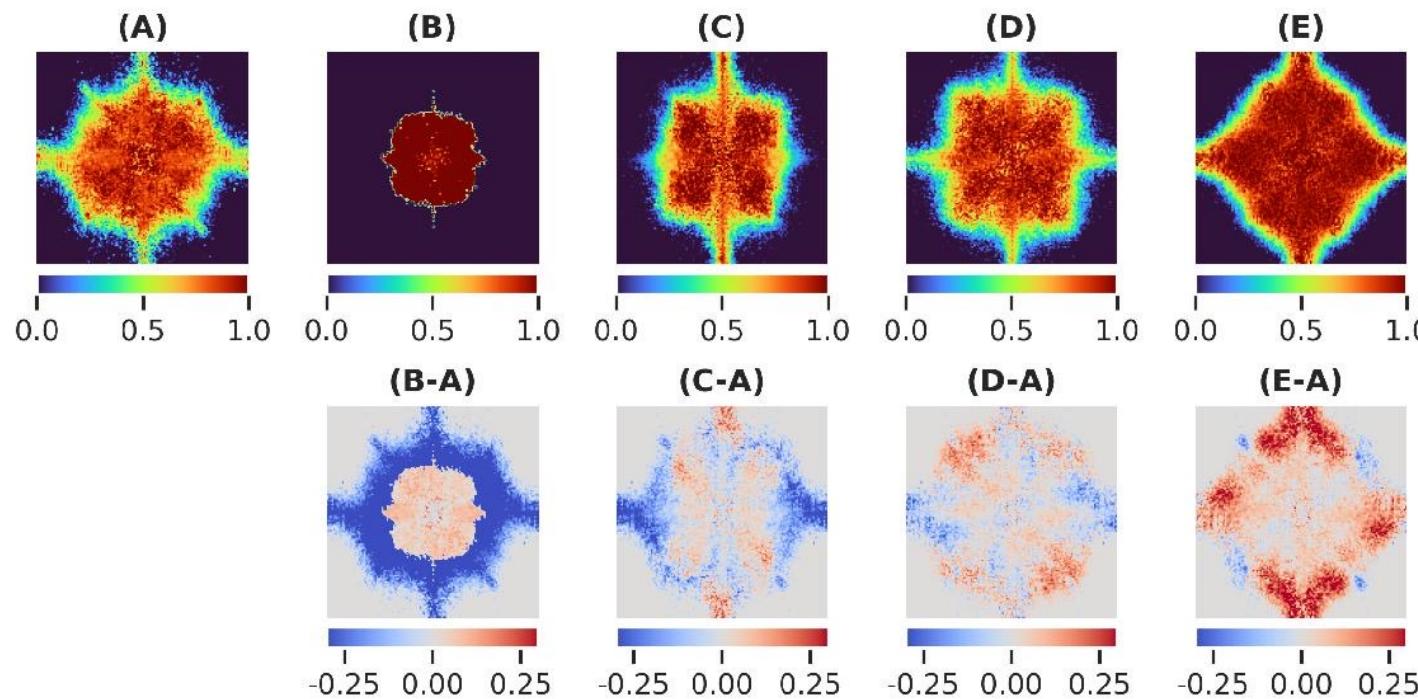
- Important frequencies are very sparse.
- The performance is preserved with the mask
- The performance drops 50% with complementary mask

# Energy distribution



Lower frequencies are more important in the adversarially trained net w.r.t. the naive net.

# Data augmentation (full dataset): testing the low-frequency bias



**A:** Naive

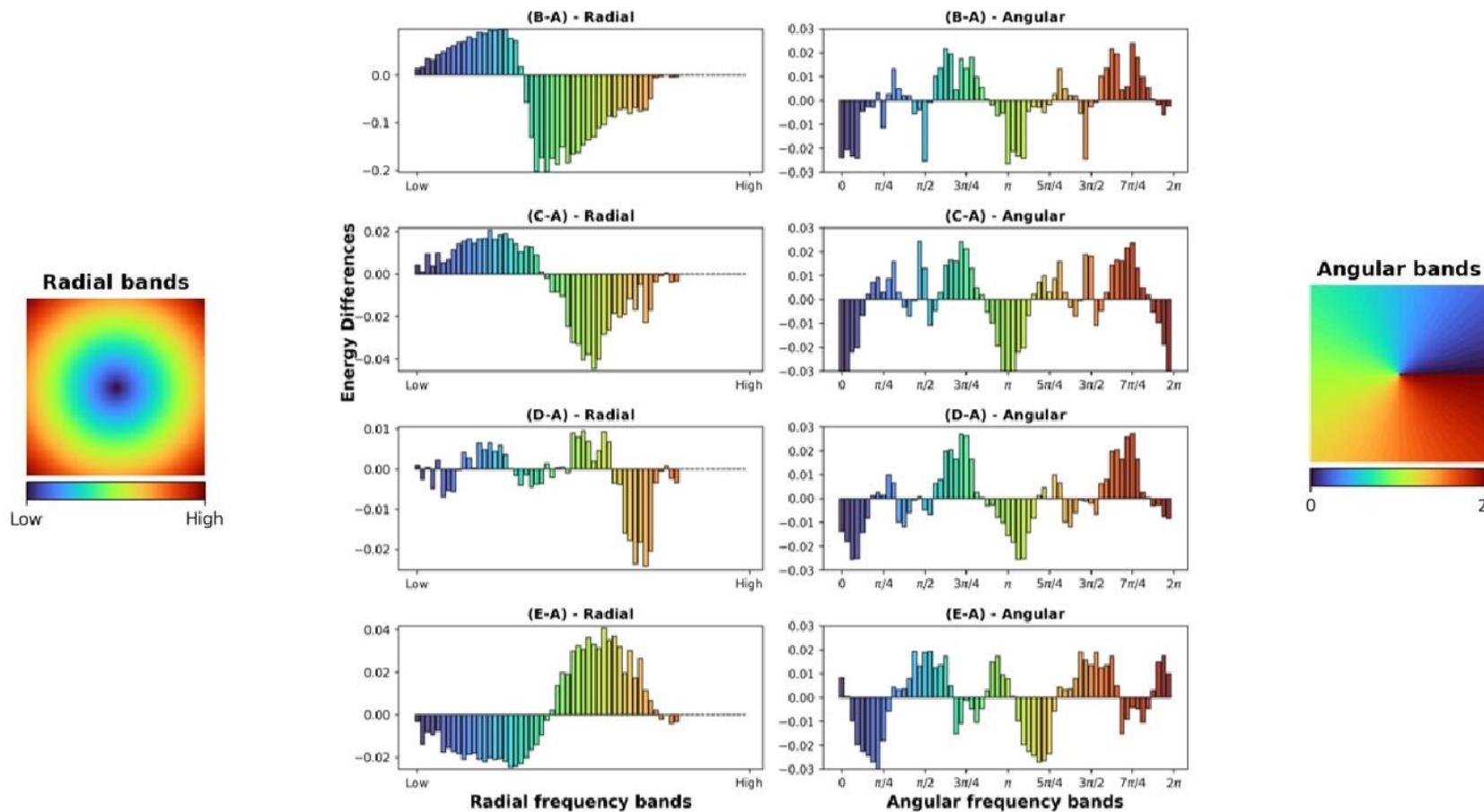
**B:** Adversarially trained

**C:** Scale augmentation

**D:** Translation augmentation

**E:** Rotation augmentation

# Data augmentation (full dataset): testing the low-frequency bias



\*Adversarial training and scale are biasing towards lower frequencies.  
\*The bias is highly direction dependent.

## Where does this leave us?

- Although reaching good performance ANNs are relying on features that do not seem to be those the brain is relying on.
- Test “brain-like” Anns:
  - ① Train an ANN to a recognition task that an animal also performs well.
  - ② Find the essential frequencies and produce filtered images.
  - ③ Show back to the animal the filtered images and test for the animal performance.
- Generate “brain-like” Anns:
  - ① Online learning of a frequency filter that maximize real neurons activity (how?)
  - ② Learn a network to do e.g. recognition on mask-filtered images.