



Pacific Ballroom #26

Loss Landscapes of Regularized Linear Autoencoders



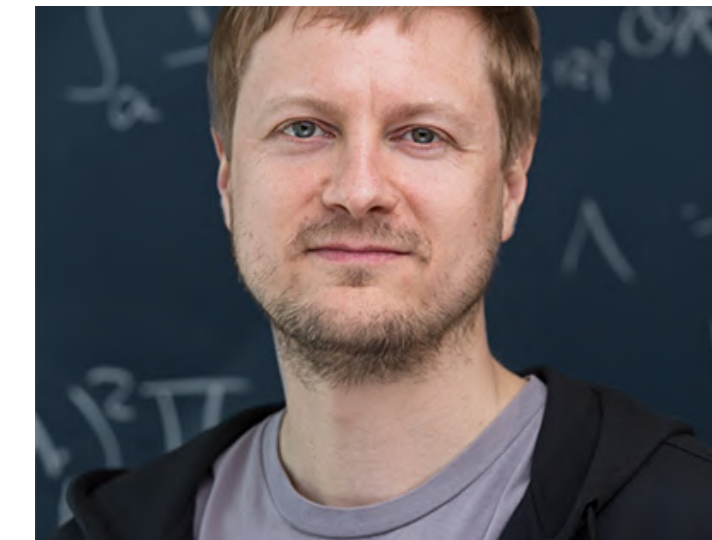
Daniel Kunin



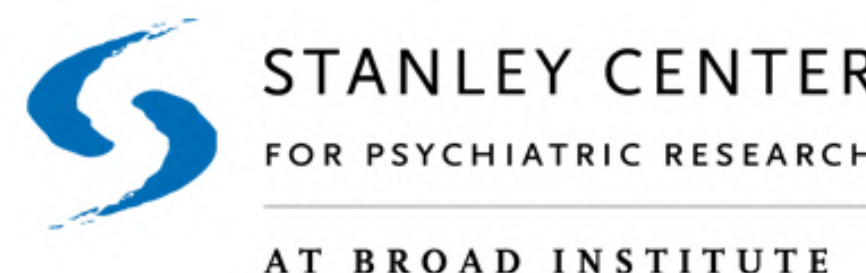
Jonathan M. Bloom



Aleksandrina Goeva



Cotton Seed

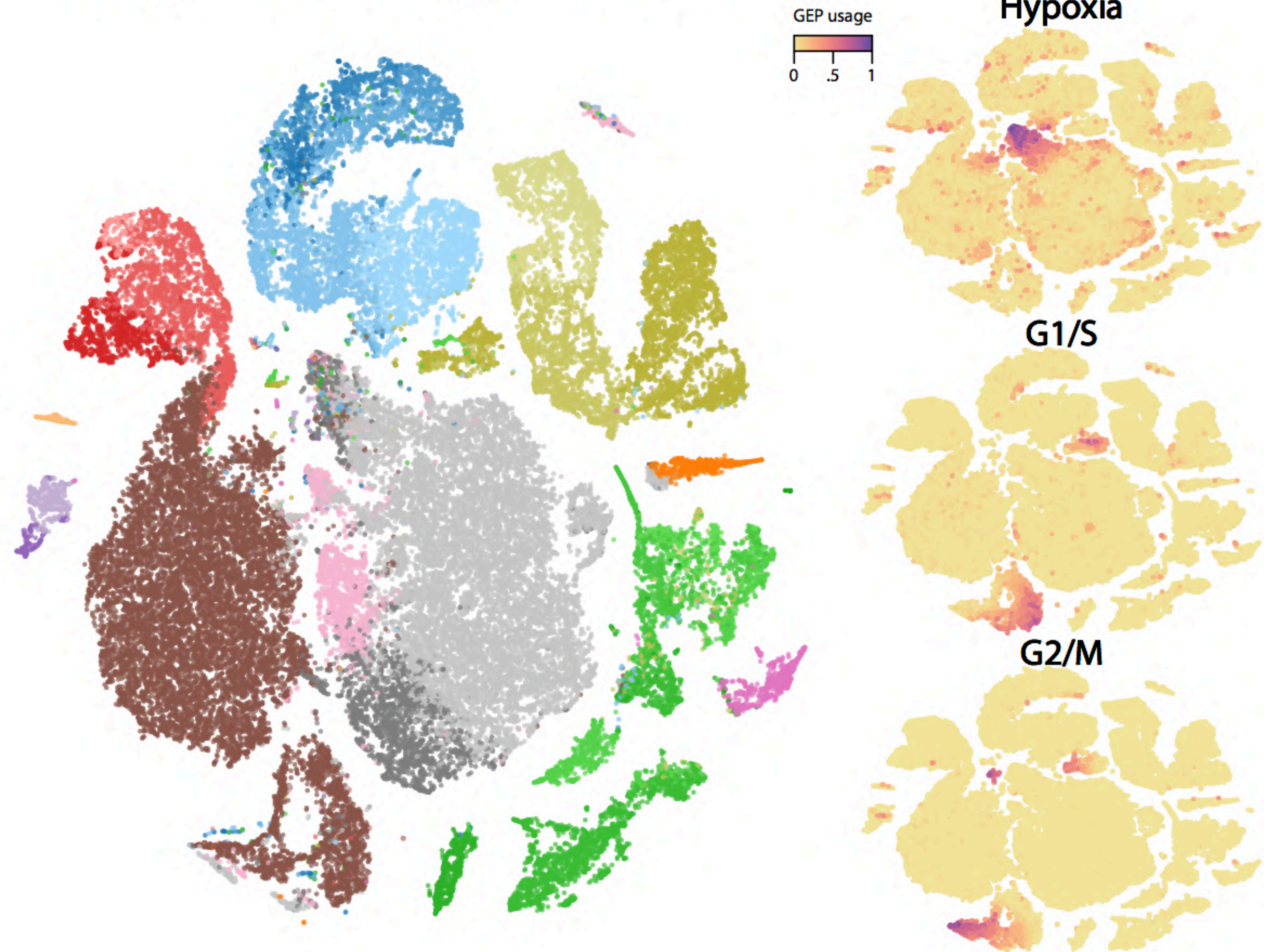


Can we use **autoencoders** to learn meaningful representations of **neurons** from their gene expression?

Organoid cell-types and activities

Identity GEP

- | | |
|-----------|-------------|
| ● Astro-1 | ● FB-3 |
| ● Astro-2 | ● Dop-1 |
| ● Astro-3 | ● Dop-2 |
| ● Astro-4 | ● NE-1 |
| ● Astro-5 | ● NE-2 |
| ● Astro-6 | ● Stem-like |
| ● Ret-1 | ● PP |
| ● Ret-2 | ● Musc-T1 |
| ● Ret-3 | ● Musc-Im |
| ● Ret-4 | ● Musc-T2 |
| ● Ret-5 | ● C6-1 |
| ● Ret-6 | ● C6-2 |
| ● FB-1 | ● C7 |
| ● FB-2 | ● C8 |

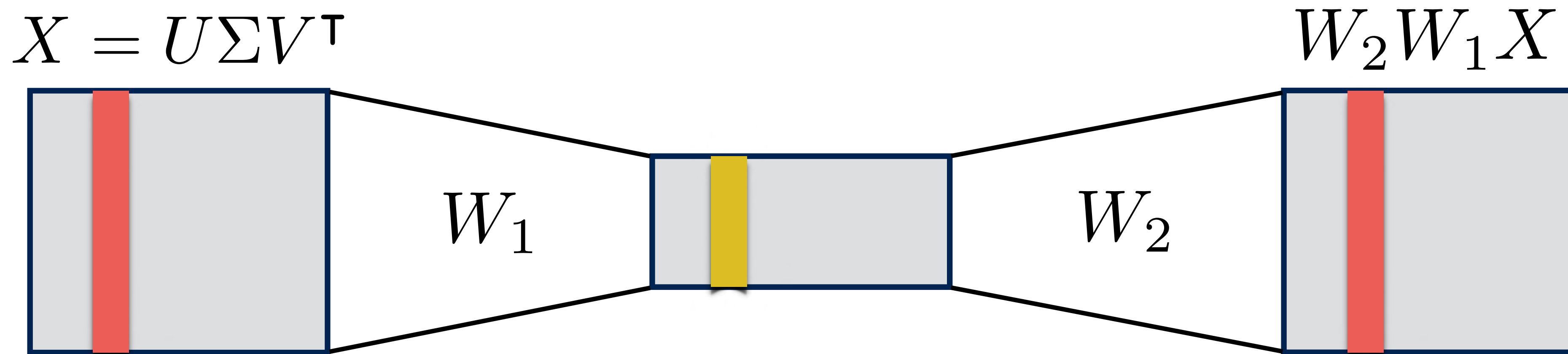


What **does** a linear autoencoder learn?

&

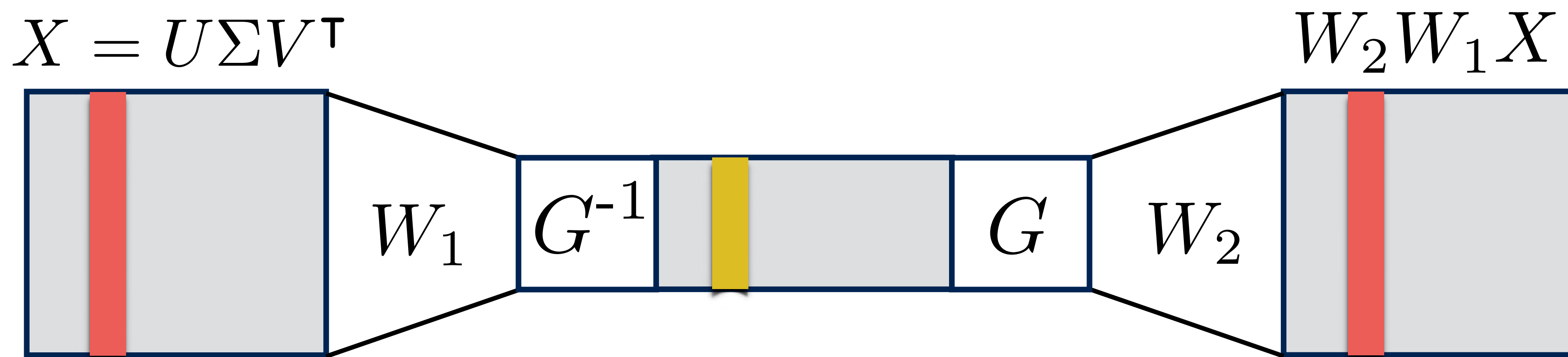
What **doesn't** a linear autoencoder learn?

Does learn the **principal subspace**

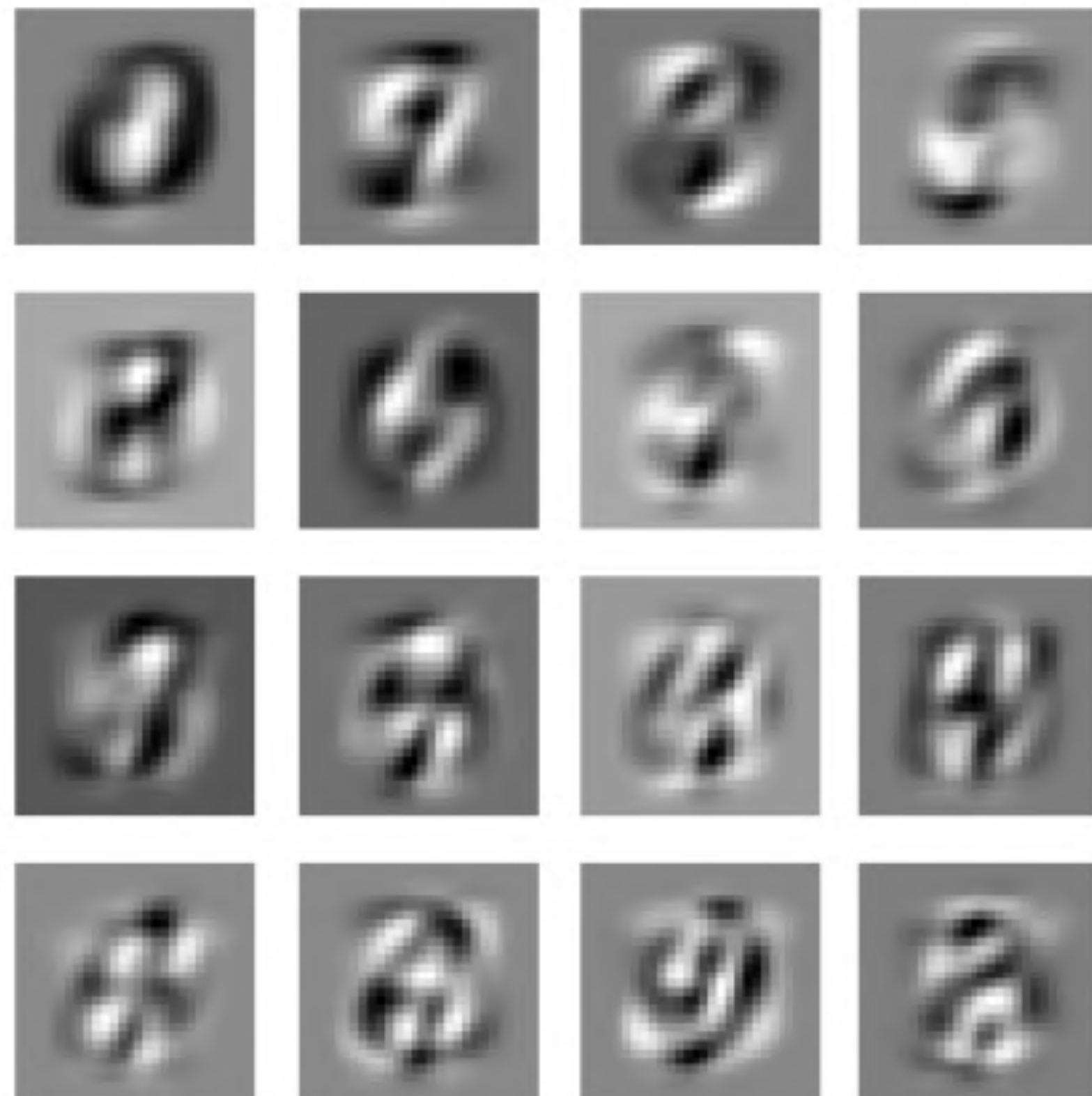


$$\mathcal{L}(W_1, W_2) = \|X - \underbrace{W_2W_1X}_{U_k U_k^\top}\|^2$$

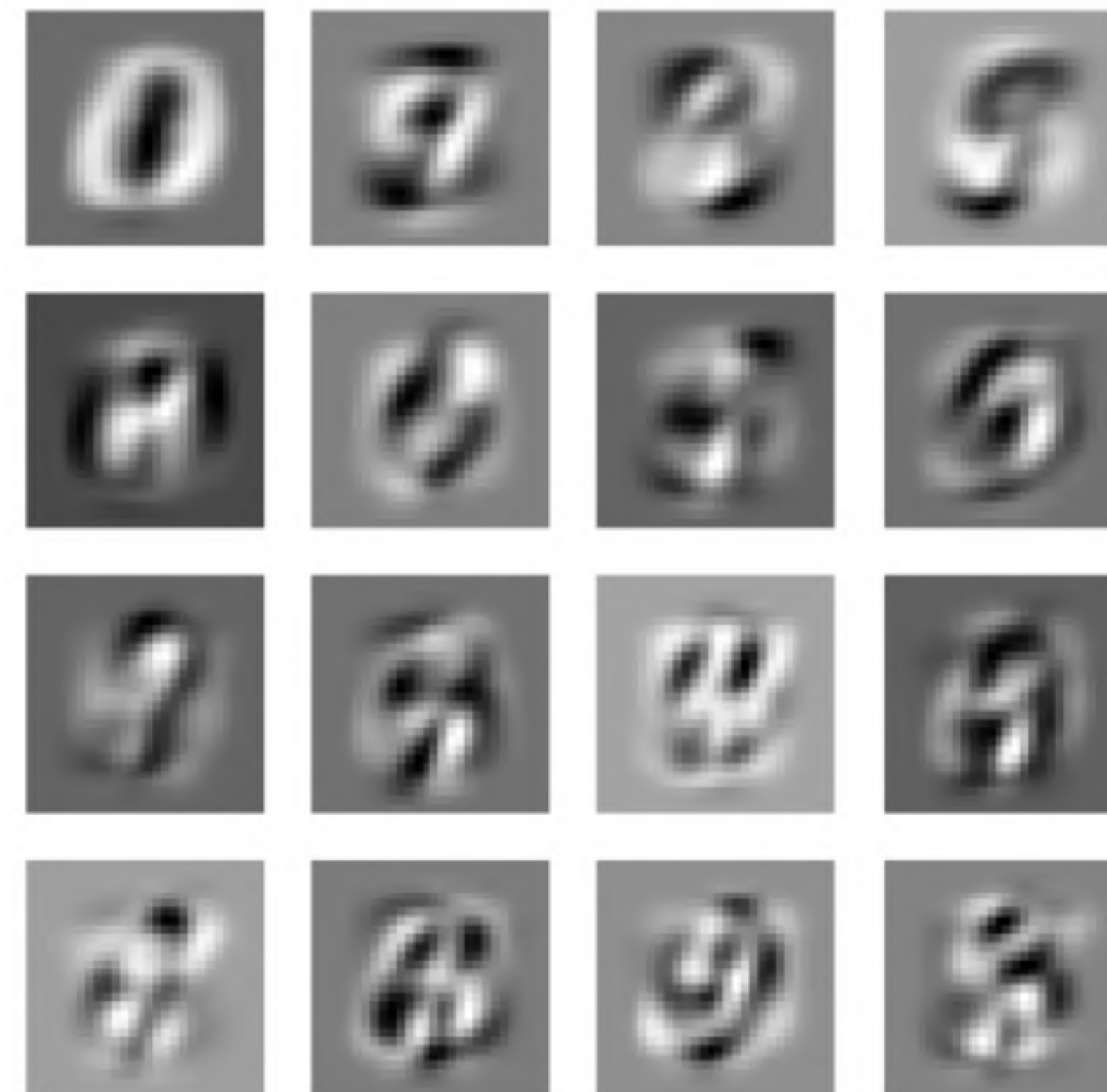
Doesn't learn the **principal directions** or **eigenvalues**



$$\mathcal{L}(W_1, W_2) = \left\| X - \underbrace{W_2W_1X}_{U_k G G^{-1} U_k^T} \right\|^2$$



Principal Directions of X



Left Singular Vectors of W_2

Adding Regularization

$$\mathcal{L}(W_1, W_2) = \|X - W_2 W_1 X\|_F^2$$



$$\mathcal{L}_\sigma(W_1, W_2) = \mathcal{L}(W_1, W_2) + \lambda(\|W_1\|_F^2 + \|W_2\|_F^2)$$

Regularization and Orthogonality

1. Orthogonal matrices are the volume-preserving matrices of minimal Frobenius norm.

$$\min_A \|A\|_F^2 \quad \text{s.t.} \quad \det(A)^2 = 1$$



2. Orthogonal matrices are the inverse matrices of minimal Frobenius norm.

$$\min_{A,B} \|A\|_F^2 + \|B\|_F^2 \quad \text{s.t.} \quad AB = I$$

In particular $A = B^T$ at all minima.

Scalar Case



$$(x - w_2 w_1 x)^2$$

Critical points \longrightarrow

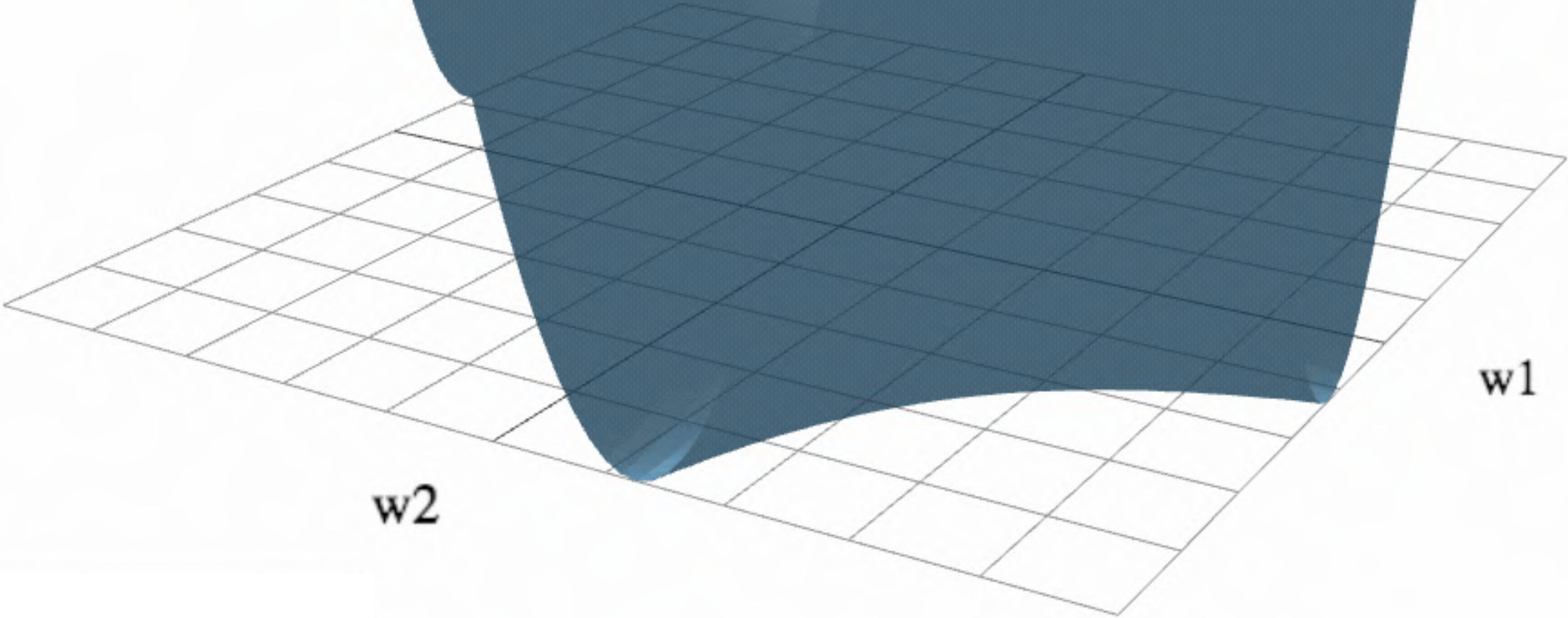
$$\begin{cases} w_1 = w_2 = 0 \\ w_2 w_1 = 1 \end{cases}$$

$$(x - w_2 w_1 x)^2 + \lambda(w_1^2 + w_2^2)$$

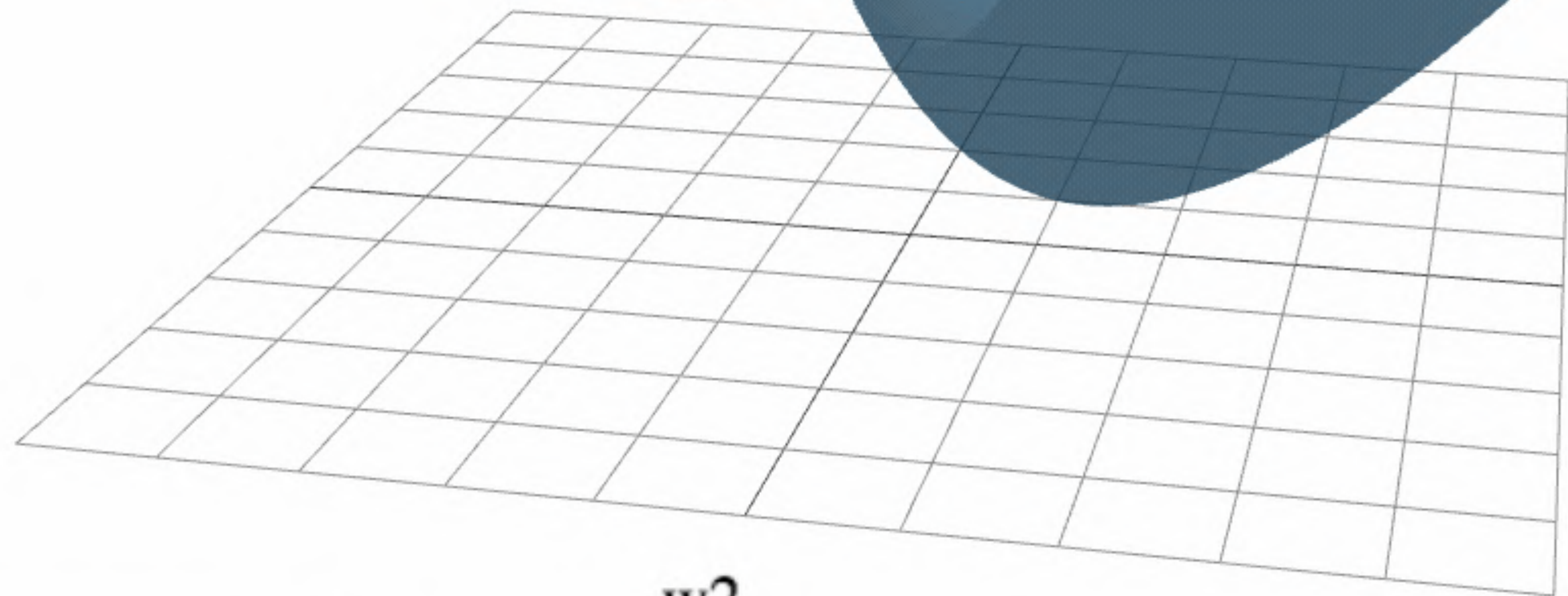
Critical points \longrightarrow

$$\begin{cases} w_1 = w_2 = 0 \\ w_2 w_1 = (1 - \lambda x^{-2}) \\ w_1 = w_2 \end{cases}$$

$$\lambda = 0$$



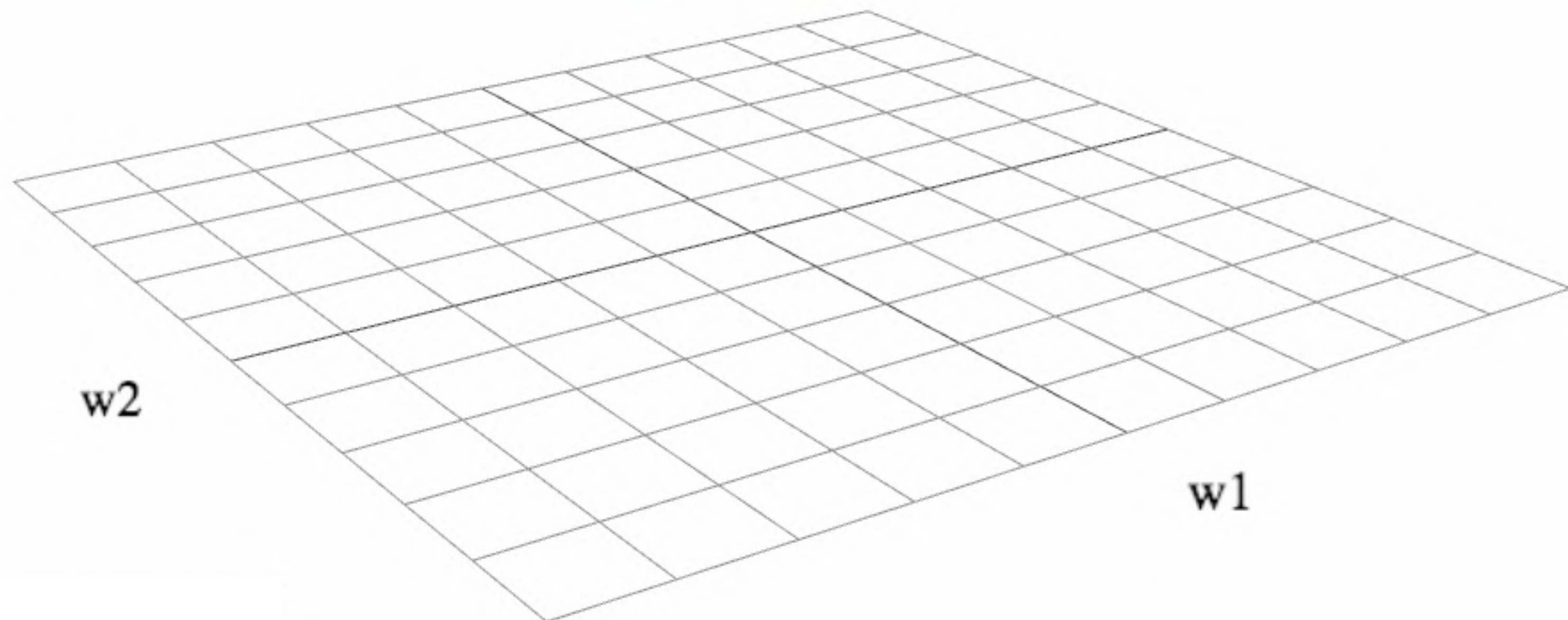
$$x^2 > \lambda > 0$$



w1

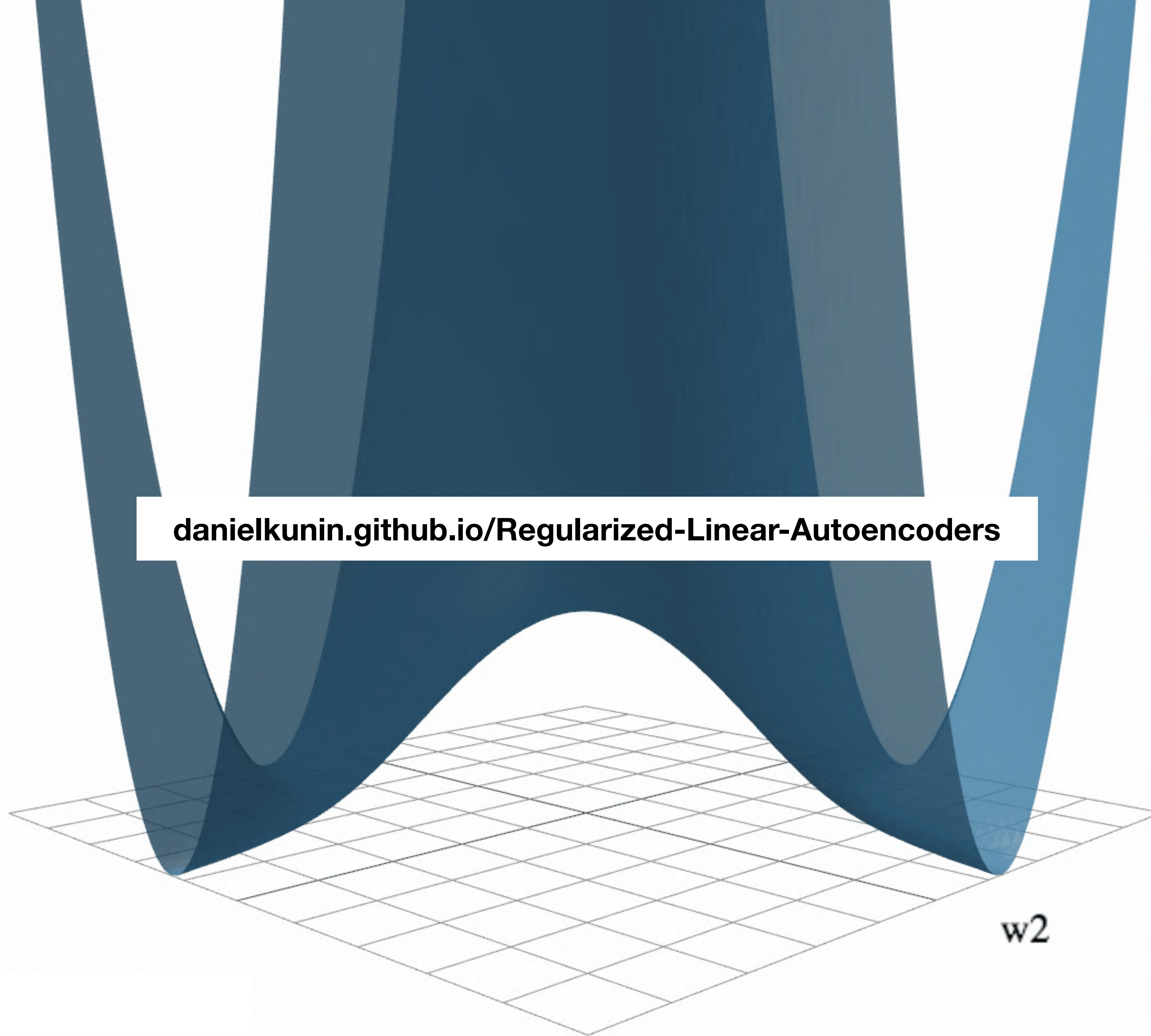
w2

$$\lambda > x^2$$



danielkunin.github.io/Regularized-Linear-Autoencoders

w_2



Loss Functions

Unregularized:

$$\mathcal{L}(W_1, W_2) = \|X - W_2 W_1 X\|_F^2$$

Product Regularized:

$$\mathcal{L}_\pi(W_1, W_2) = \mathcal{L}(W_1, W_2) + \lambda \|W_2 W_1\|_F^2$$

Sum Regularized:

$$\mathcal{L}_\sigma(W_1, W_2) = \mathcal{L}(W_1, W_2) + \lambda (\|W_1\|_F^2 + \|W_2\|_F^2)$$

Theorem 4.2 (Landscape Theorem).

The critical landscape is diffeomorphic to the space of pairs (\mathcal{I}, G) or (\mathcal{I}, O) with

- $\mathcal{I} \subset \{1, \dots, m\}$ of size $0 \leq l \leq k$,
- $G \in \mathbb{R}^{k \times l}$ with independent columns,
- $O \in \mathbb{R}^{k \times l}$ with orthonormal columns.

	W_2	W_1
\mathcal{L}	$U_{\mathcal{I}} G^+$	$G U_{\mathcal{I}}^{\top}$
\mathcal{L}_{π}	$U_{\mathcal{I}} (I_{\ell} + \lambda \Sigma_{\mathcal{I}}^{-2})^{-\frac{1}{2}} G^+$	$G (I_{\ell} + \lambda \Sigma_{\mathcal{I}}^{-2})^{-\frac{1}{2}} U_{\mathcal{I}}^{\top}$
\mathcal{L}_{σ}	$U_{\mathcal{I}} (I_{\ell} - \lambda \Sigma_{\mathcal{I}}^{-2})^{\frac{1}{2}} O^{\top}$	$O (I_{\ell} - \lambda \Sigma_{\mathcal{I}}^{-2})^{\frac{1}{2}} U_{\mathcal{I}}^{\top}$

Theorem 4.2 (Landscape Theorem).

The critical landscape is diffeomorphic to the space of pairs (\mathcal{I}, G) or (\mathcal{I}, O) with

- $\mathcal{I} \subset \{1, \dots, m\}$ of size $0 \leq l \leq k$,
- $G \in \mathbb{R}^{k \times l}$ with independent columns,
- $O \in \mathbb{R}^{k \times l}$ with orthonormal columns.

	W_2	W_1
\mathcal{L}	$U_{\mathcal{I}} G^+$	$G U_{\mathcal{I}}^{\top}$
\mathcal{L}_{π}	$U_{\mathcal{I}} (I_{\ell} + \lambda \Sigma_{\mathcal{I}}^{-2})^{-\frac{1}{2}} G^+$	$G (I_{\ell} + \lambda \Sigma_{\mathcal{I}}^{-2})^{-\frac{1}{2}} U_{\mathcal{I}}^{\top}$
\mathcal{L}_{σ}	$U_{\mathcal{I}} (I_{\ell} - \lambda \Sigma_{\mathcal{I}}^{-2})^{\frac{1}{2}} O^{\top}$	$O (I_{\ell} - \lambda \Sigma_{\mathcal{I}}^{-2})^{\frac{1}{2}} U_{\mathcal{I}}^{\top}$

Theorem 4.2 (Landscape Theorem).

The critical landscape is diffeomorphic to the space of pairs (\mathcal{I}, G) or (\mathcal{I}, O) with

- $\mathcal{I} \subset \{1, \dots, m\}$ of size $0 \leq l \leq k$,
- $G \in \mathbb{R}^{k \times l}$ with independent columns,
- $O \in \mathbb{R}^{k \times l}$ with orthonormal columns.

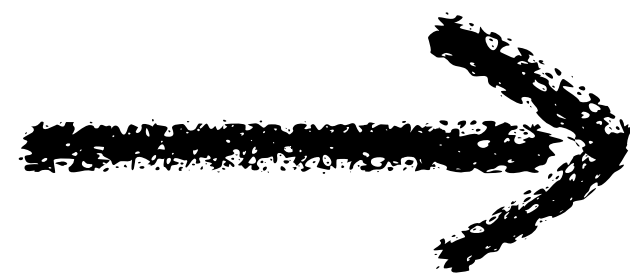
	W_2	W_1
\mathcal{L}	$U_{\mathcal{I}} G^+$	$G U_{\mathcal{I}}^{\top}$
\mathcal{L}_{π}	$U_{\mathcal{I}} (I_{\ell} + \lambda \Sigma_{\mathcal{I}}^{-2})^{-\frac{1}{2}} G^+$	$G (I_{\ell} + \lambda \Sigma_{\mathcal{I}}^{-2})^{-\frac{1}{2}} U_{\mathcal{I}}^{\top}$
\mathcal{L}_{σ}	$U_{\mathcal{I}} (I_{\ell} - \lambda \Sigma_{\mathcal{I}}^{-2})^{\frac{1}{2}} O^{\top}$	$O (I_{\ell} - \lambda \Sigma_{\mathcal{I}}^{-2})^{\frac{1}{2}} U_{\mathcal{I}}^{\top}$

Algebraic Topology of PCA

Regularization reduces the symmetry from linear to orthogonal.

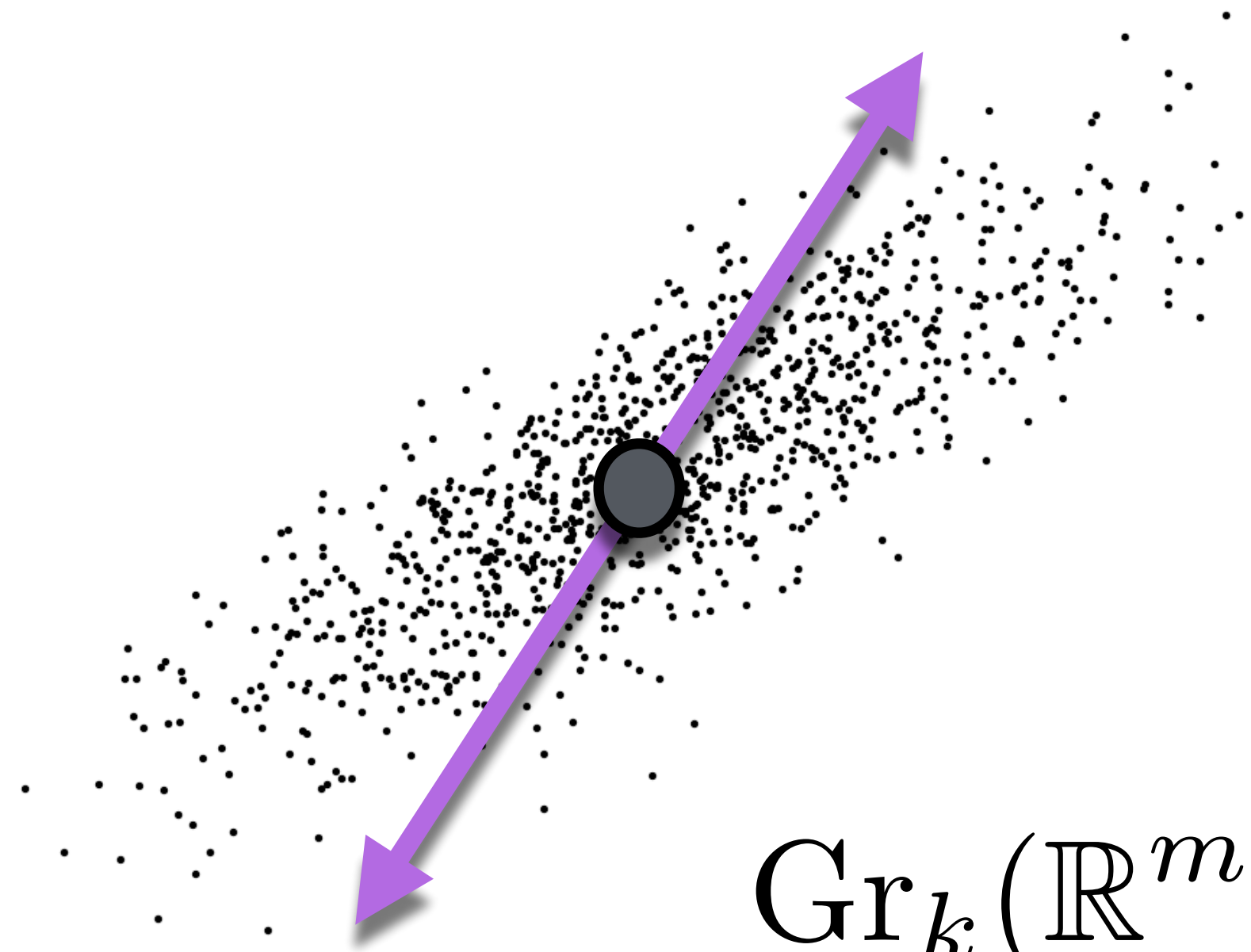
But LAEs are still over-parameterized, obscuring the gradient dynamics.

Issue: $\mathbb{R}^{k \times m} \times \mathbb{R}^{m \times k}$ is not the natural domain of PCA.



Algebraic Topology of PCA

Find the k -subspace nearest a point cloud in \mathbb{R}^m



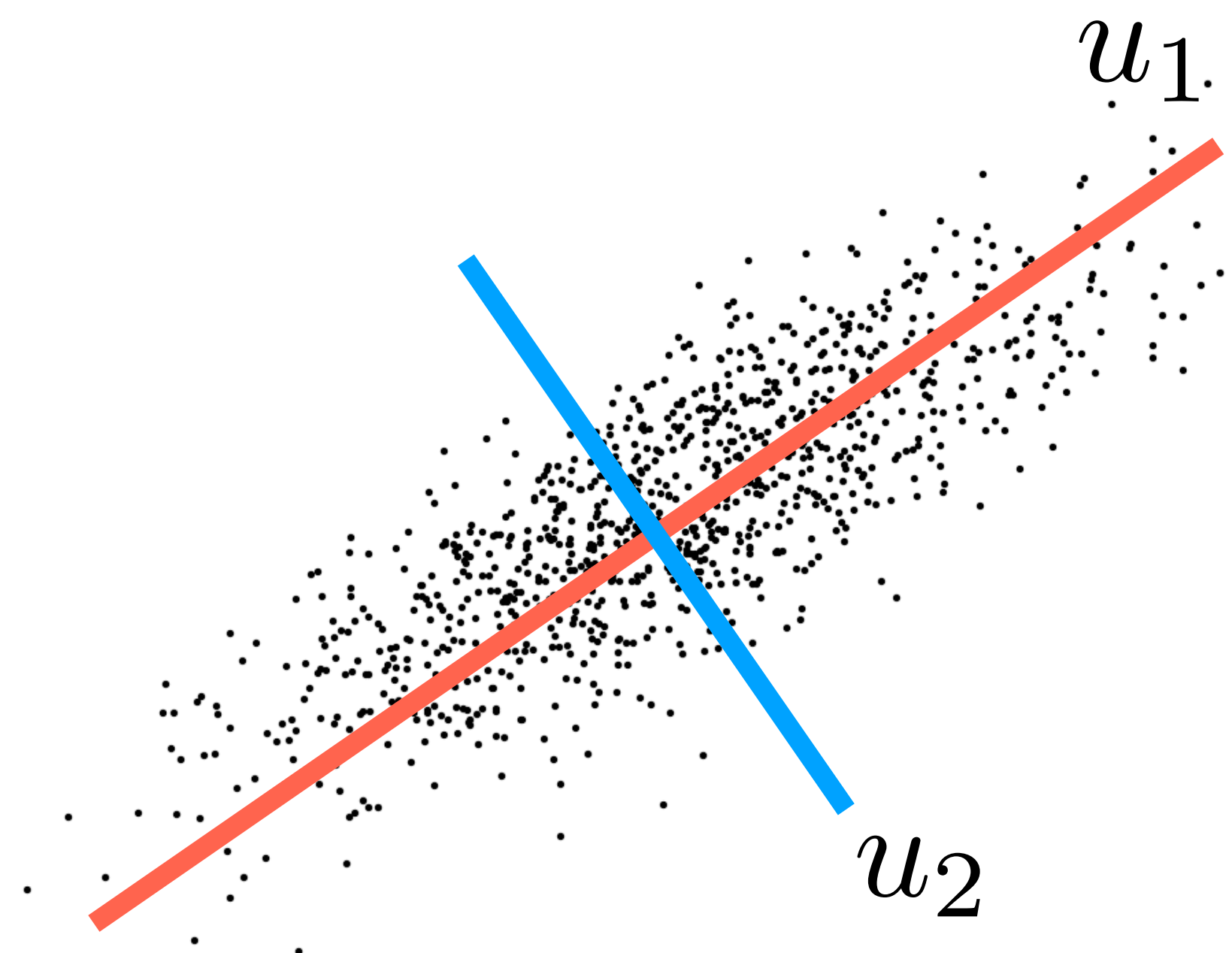
Natural domain is the *Grassmannian manifold*.

Points are subspaces. Distance is loss function.

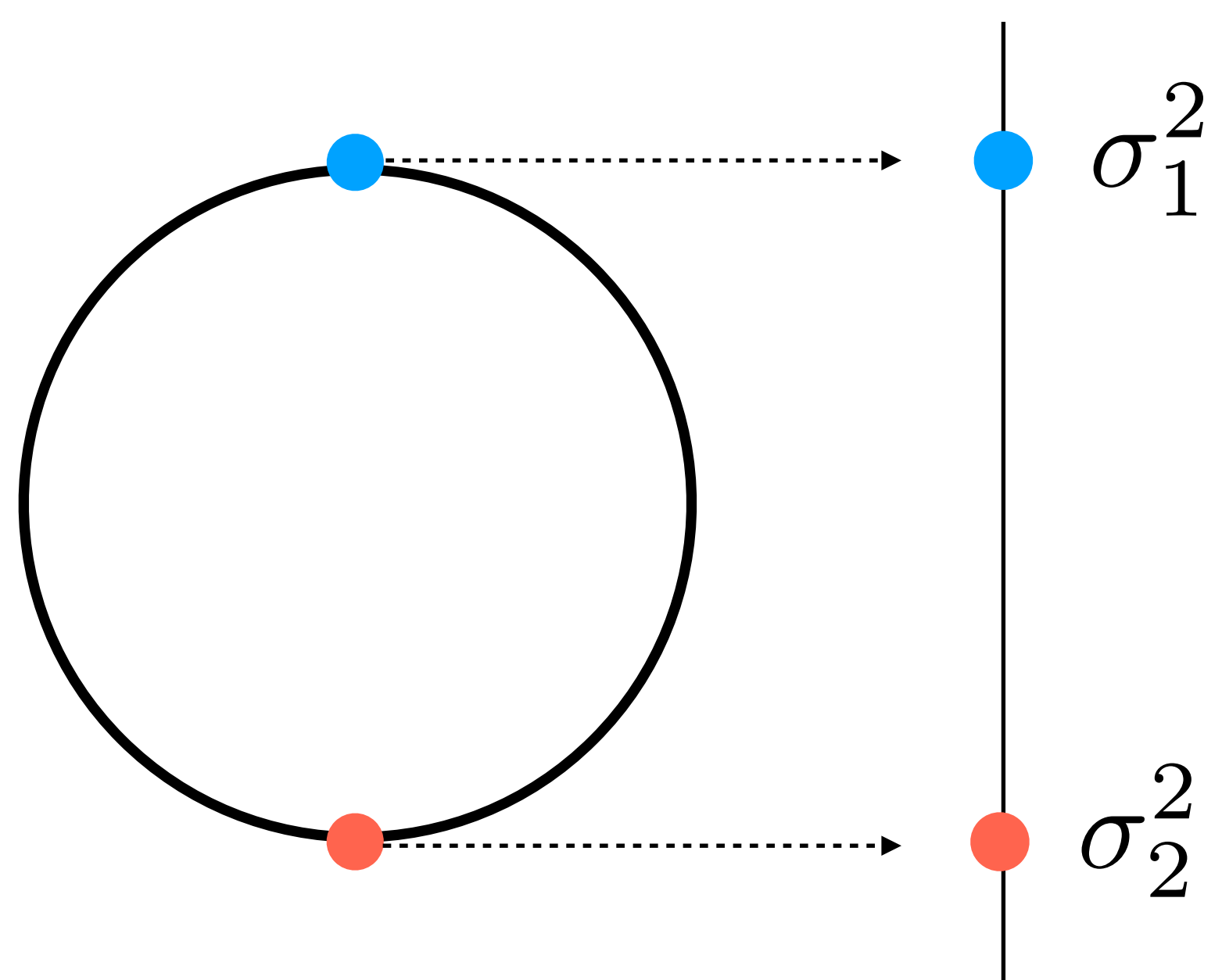
$$\text{Gr}_k(\mathbb{R}^m) \cong \{P = P^2, P = P^\top, \text{tr } P = k\} \subset \mathbb{R}^{m \times m}$$

$$\dim \text{Gr}_k(\mathbb{R}^m) = k(m - k) \qquad \mathcal{L}_X(P) = \|X - PX\|^2$$

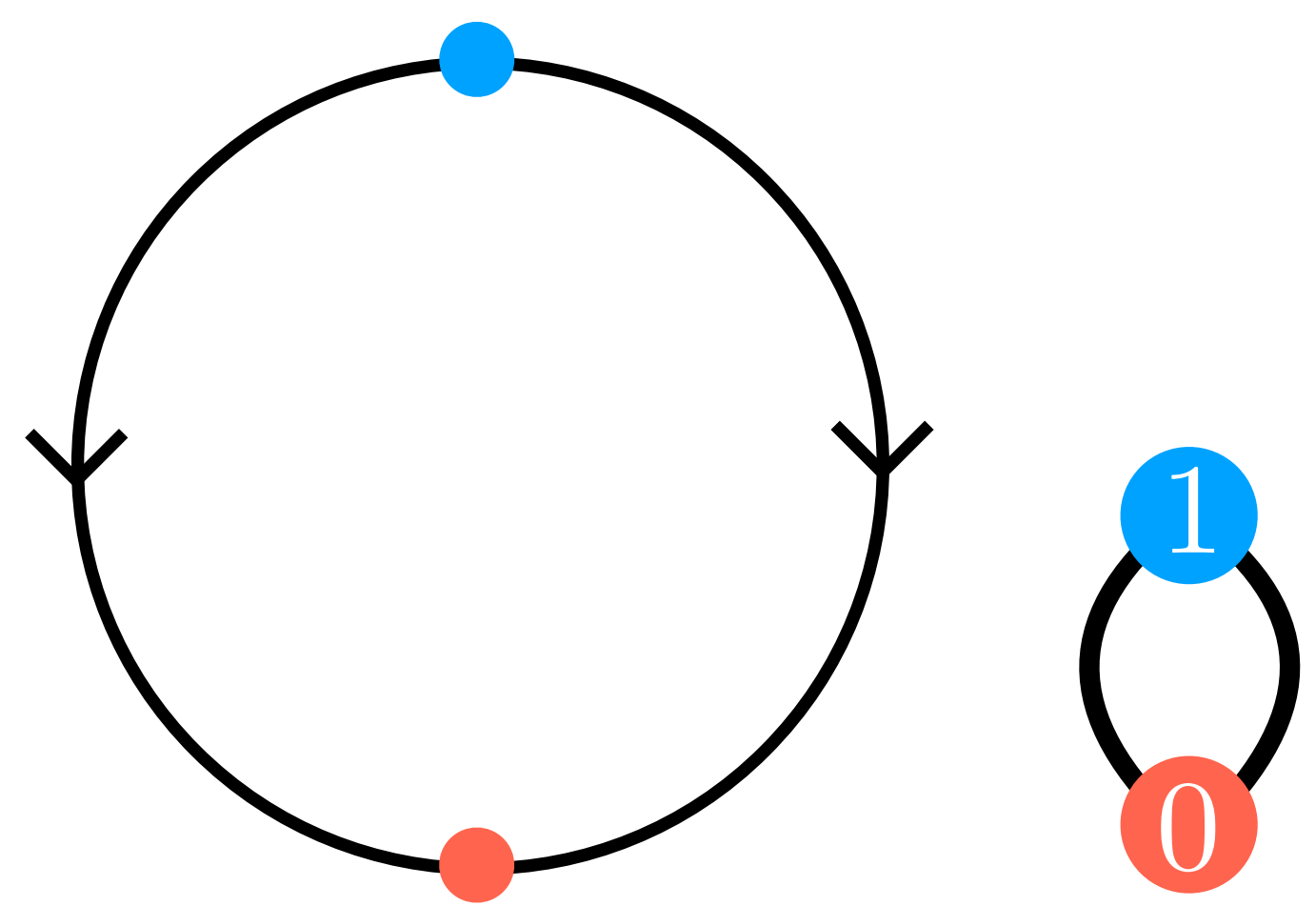
Algebraic Topology of PCA



$X \subset \mathbb{R}^2$

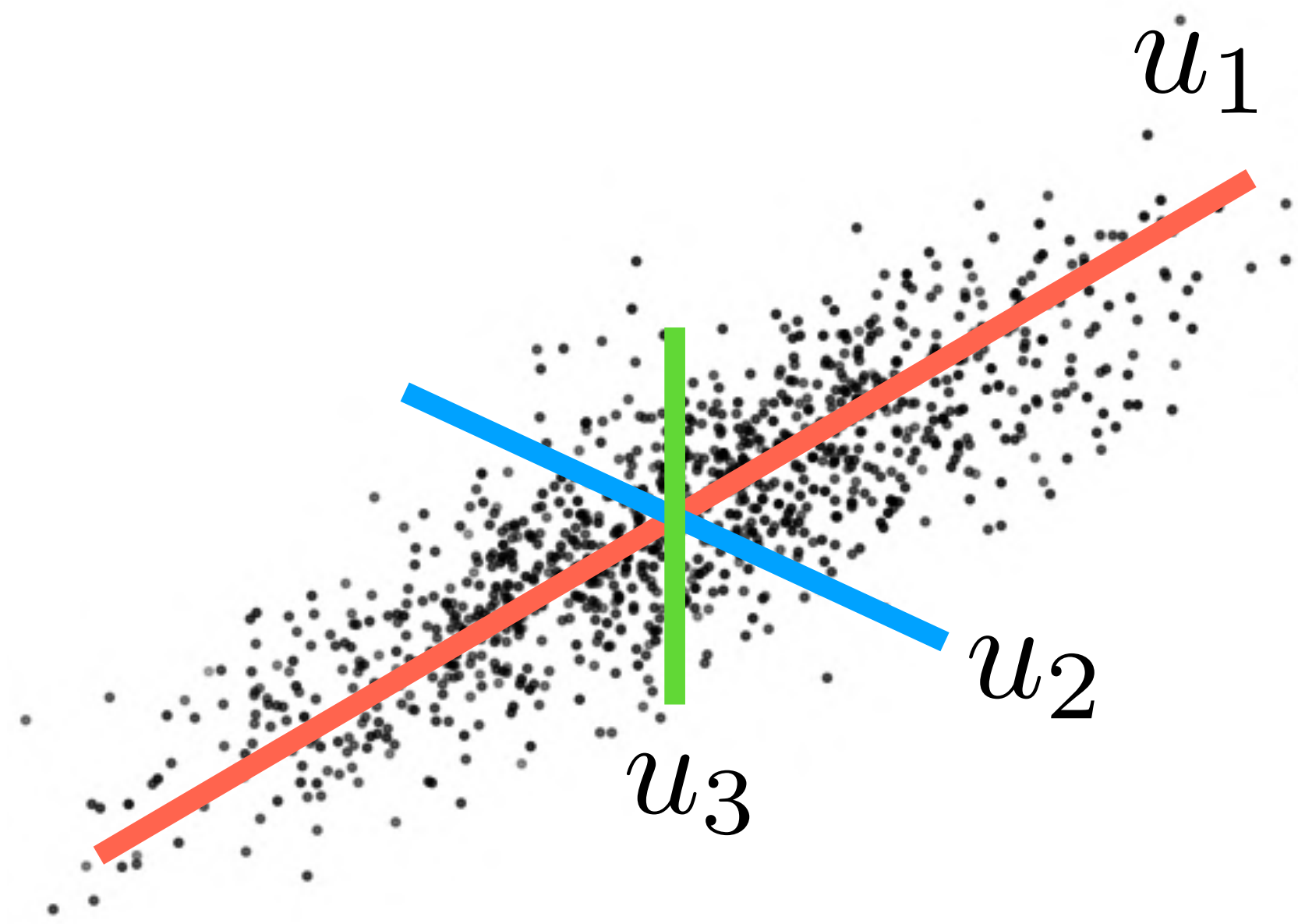


$\text{Gr}_1(\mathbb{R}^2) \xrightarrow{\mathcal{L}_X} \mathbb{R}$

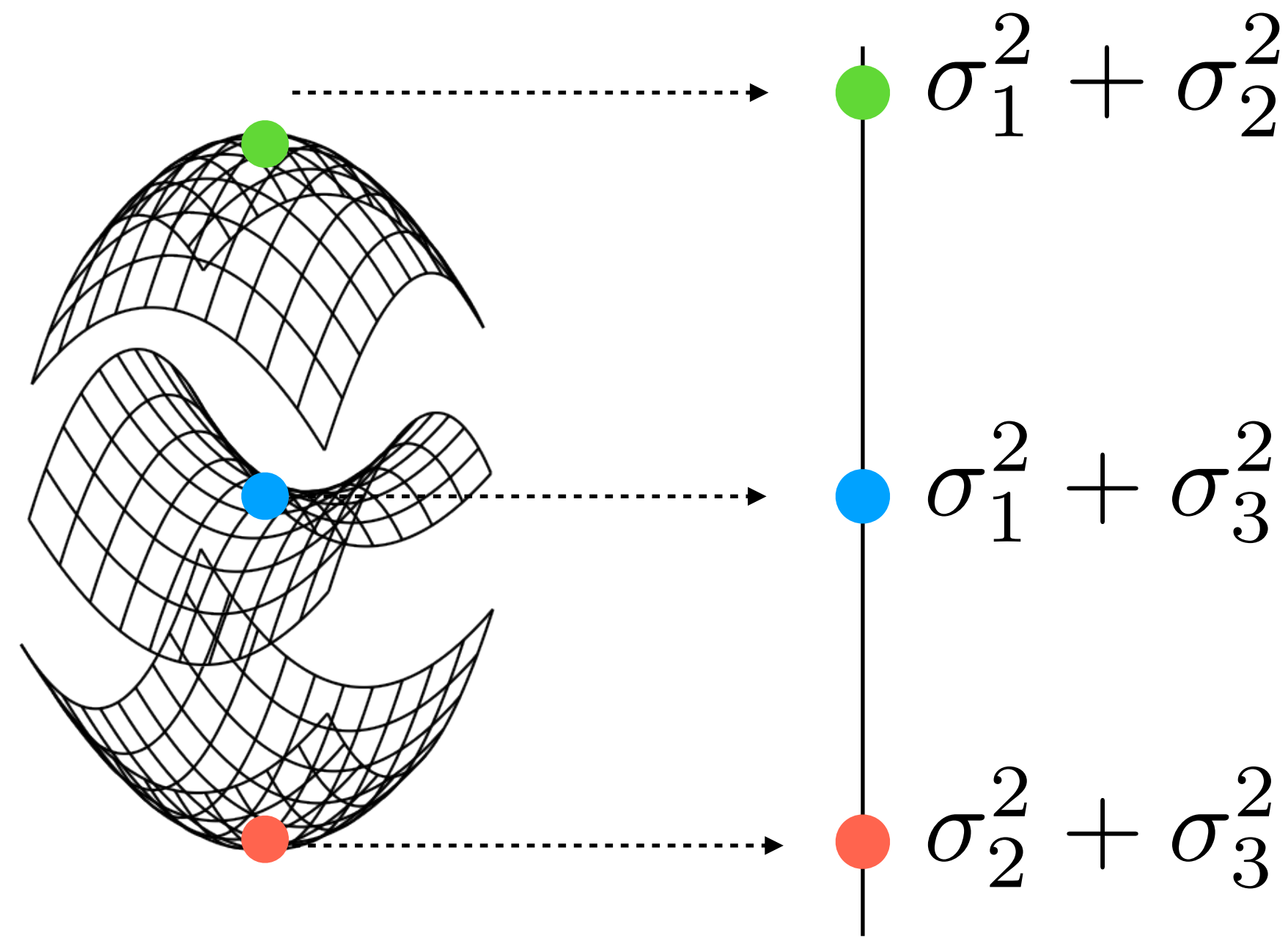


$-\nabla \mathcal{L}_X \quad H_*(\cdot)$

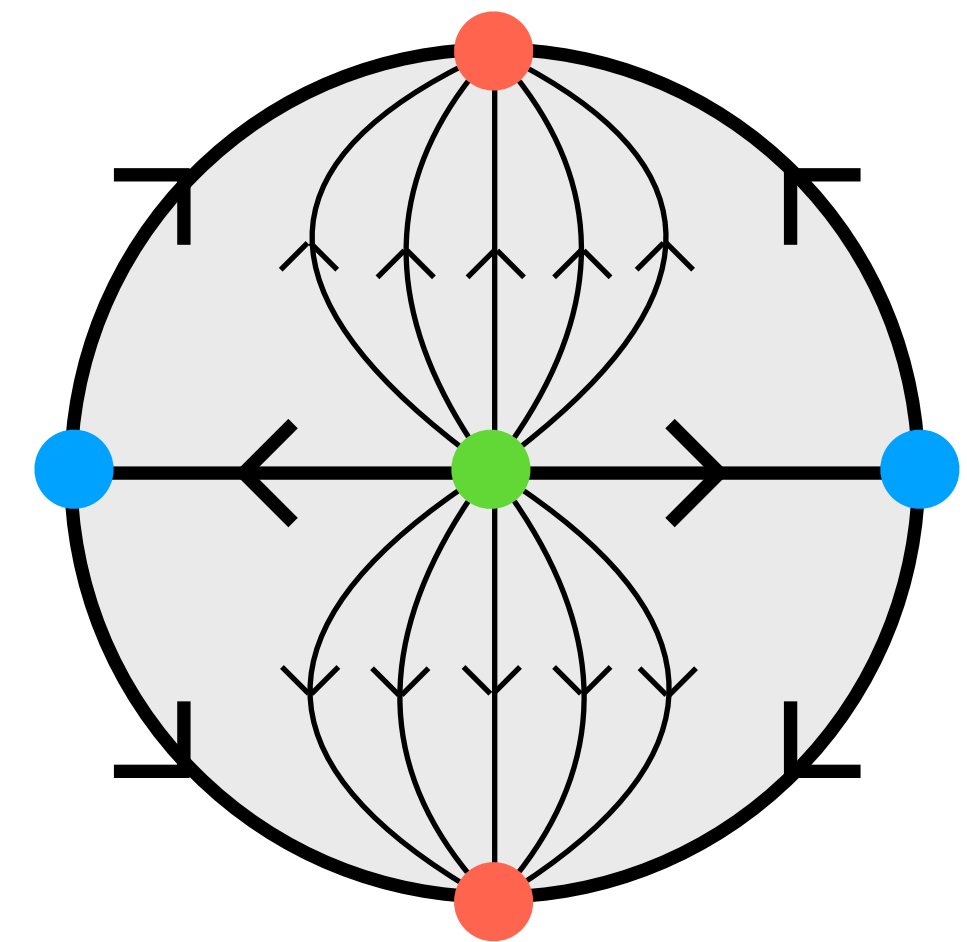
Algebraic Topology of PCA



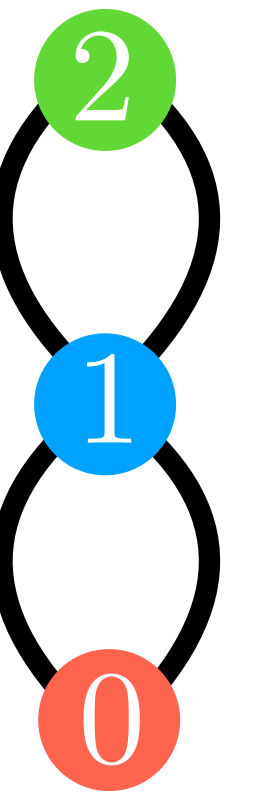
$$X \subset \mathbb{R}^3$$



$$\text{Gr}_1(\mathbb{R}^3) \xrightarrow{\mathcal{L}_X} \mathbb{R}$$



$$-\nabla \mathcal{L}_X$$

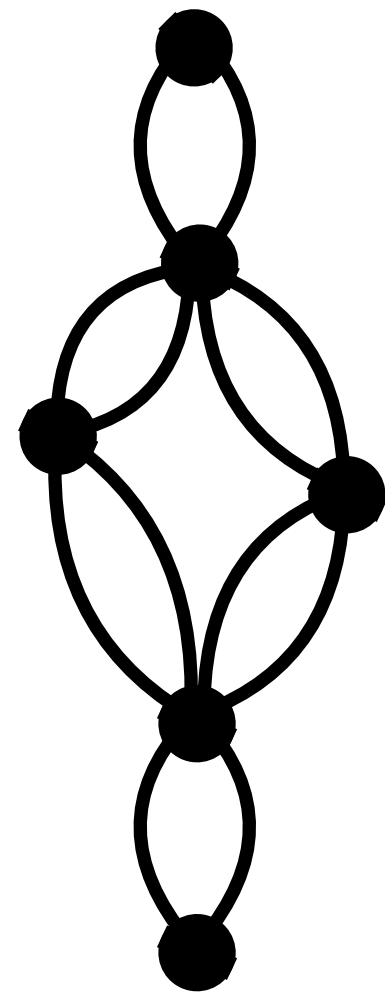


$$H_*(\cdot)$$

Algebraic Topology of PCA

$$\mathcal{L}_X : \text{Gr}_2(\mathbb{R}^4) \rightarrow \mathbb{R}$$

d	u_1	u_2	u_3	u_4
4			•	•
3		•		•
2		•	•	
2	•			•
1	•		•	
0	•	•		



- $\binom{m}{k}$ critical points are principal subspaces
- Critical values are sums of eigenvalues
- Connecting gradient trajectories are rotations
- Loss is std saddle (\mathbb{F}_2 -perfect Morse function)
- LAE loss is degenerate std saddle (*Morse-Bott*)
- Suggests principals, algorithms for deep learning

PCA Algorithms

PCA is a **two-step optimization**:

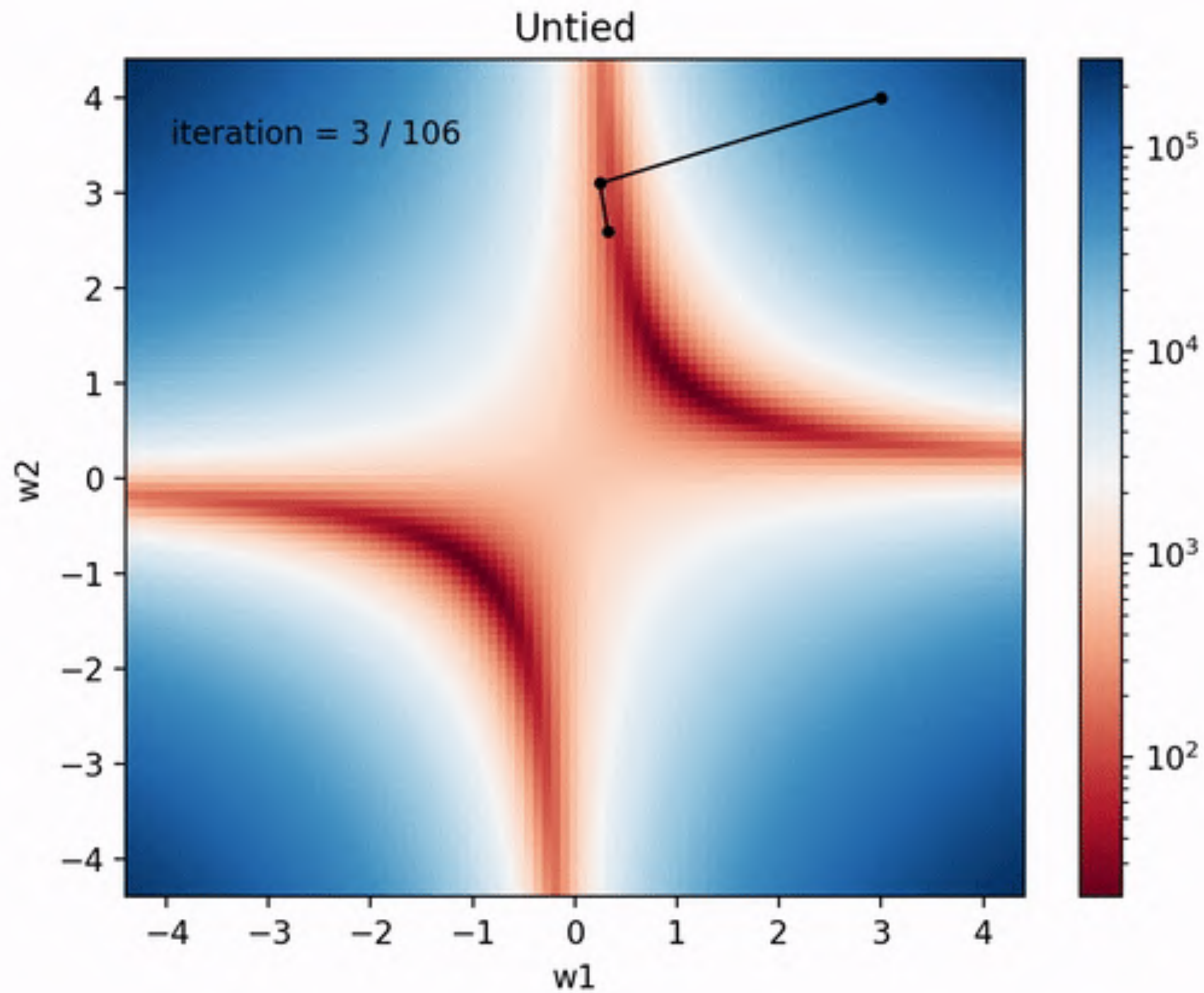
1. Train an L²-regularized LAE on $X \subset \mathbb{R}^{m \times n}$

An Adventure!

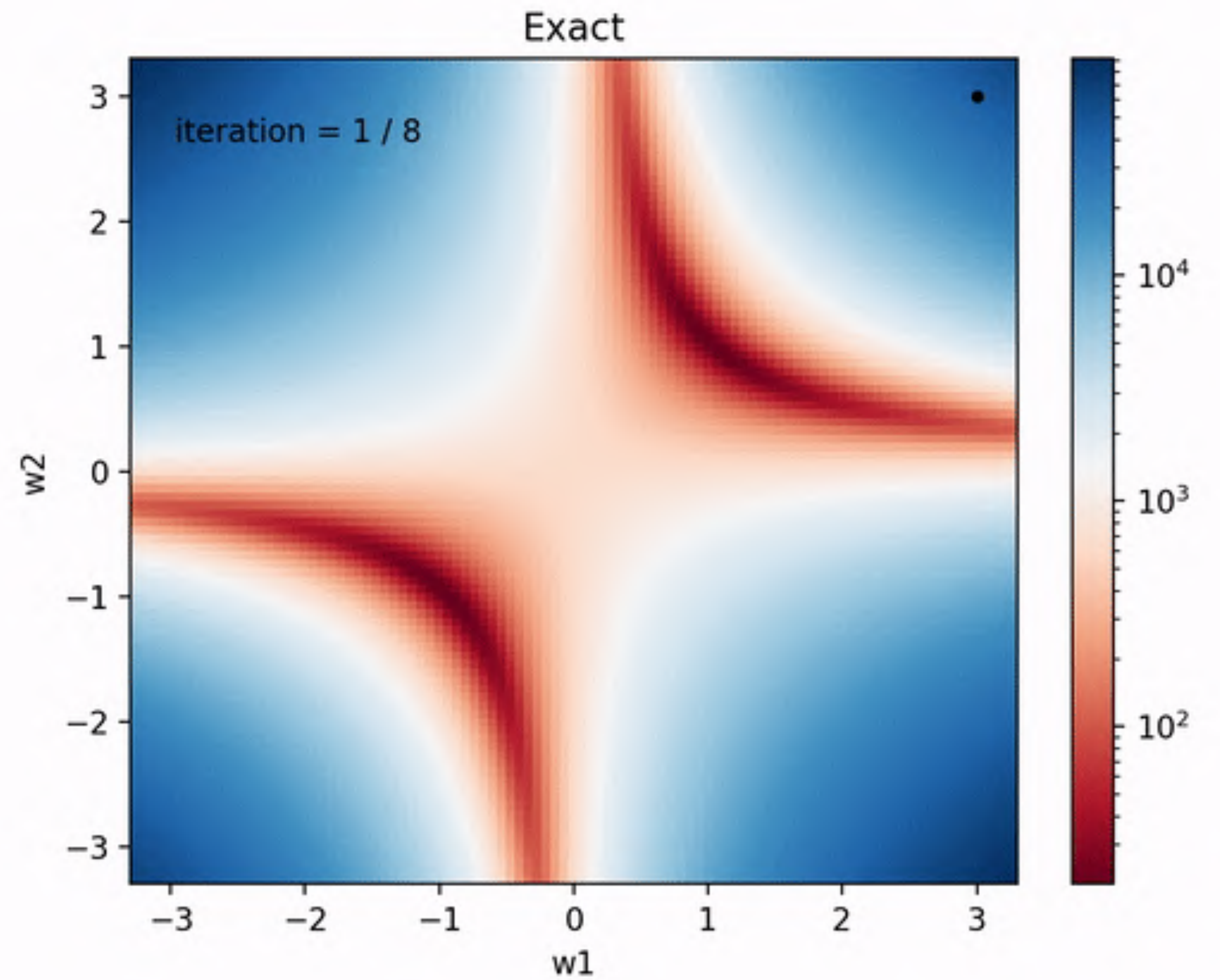
2. Apply SVD to the decoder $W_2 \subset \mathbb{R}^{m \times k}$

Quick!

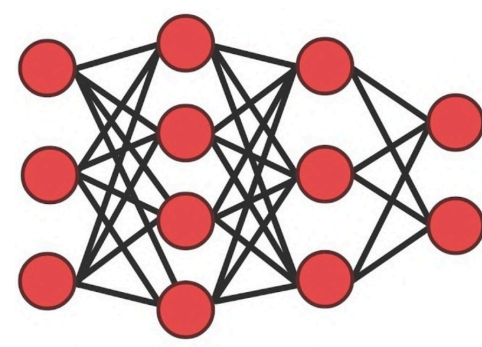
PCA Algorithms



Gradient descent



Solve for W_2 , set $W_1 = W_2^T$



Prediction in artificial neural networks is inspired by the brain.
 Is *learning* in the brain inspired by artificial neural networks?



COGNITIVE SCIENCE **11**, 23–63 (1987)

Competitive Learning: From Interactive Activation to Adaptive Resonance

STEPHEN GROSSBERG
Boston University

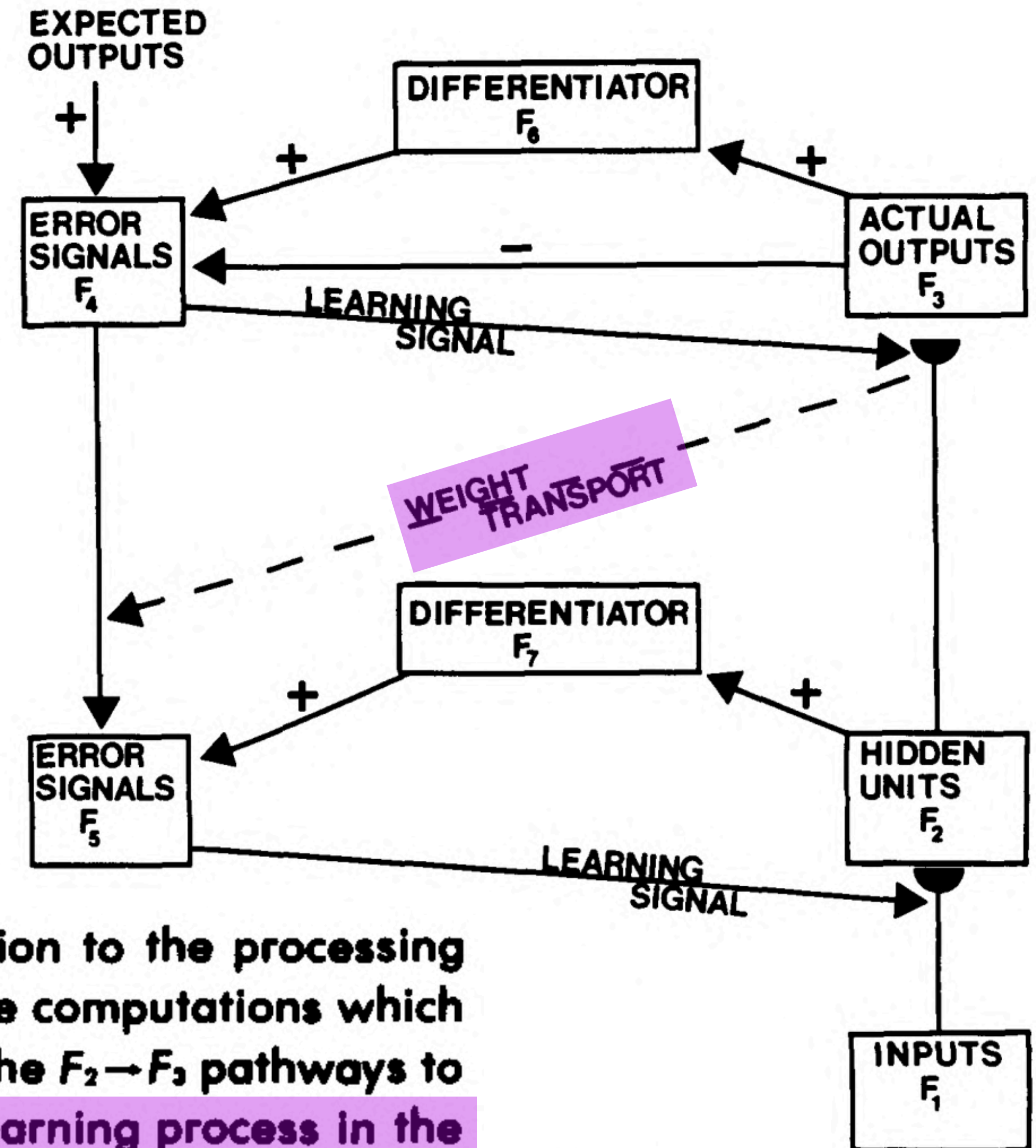
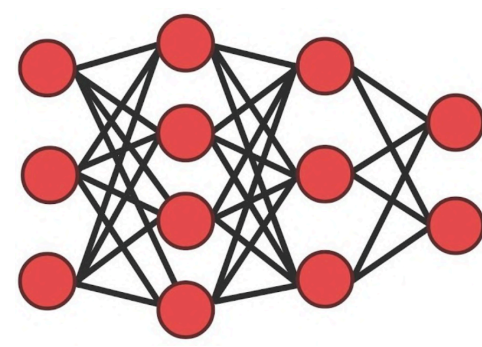


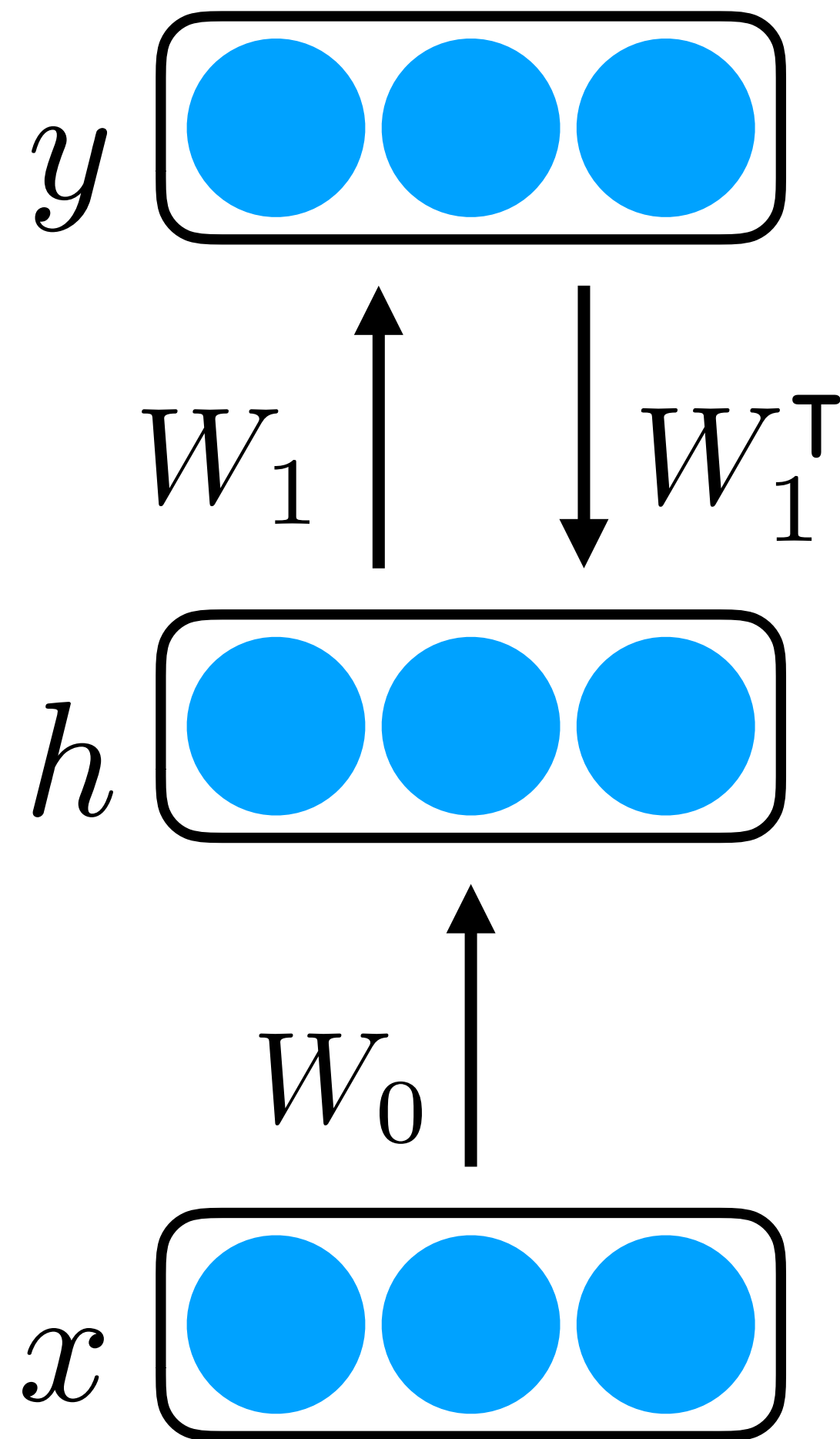
Figure 8. Circuit diagram of the back propagation model: In addition to the processing levels F_1 , F_2 , F_3 , there are also levels F_4 , F_5 , F_6 , and F_7 to carry out the computations which control the learning process. The transport of learned weights from the $F_2 \rightarrow F_3$ pathways to the $F_4 \rightarrow F_5$ pathways shows that this algorithm cannot represent a learning process in the brain.



Prediction in artificial neural networks is inspired by the brain.

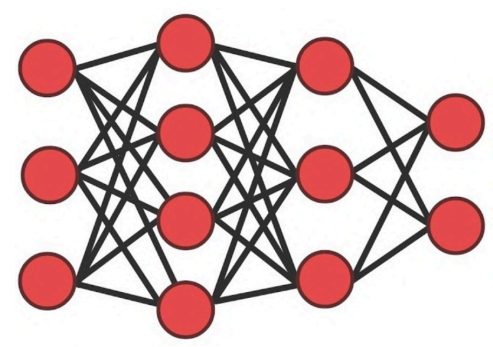


Is *learning* in the brain inspired by artificial neural networks?



weight
transport
problem

$$\Delta W_0 \propto W_1^T e x^T$$



Prediction in artificial neural networks is inspired by the brain.
Is *learning* in the brain inspired by artificial neural networks?

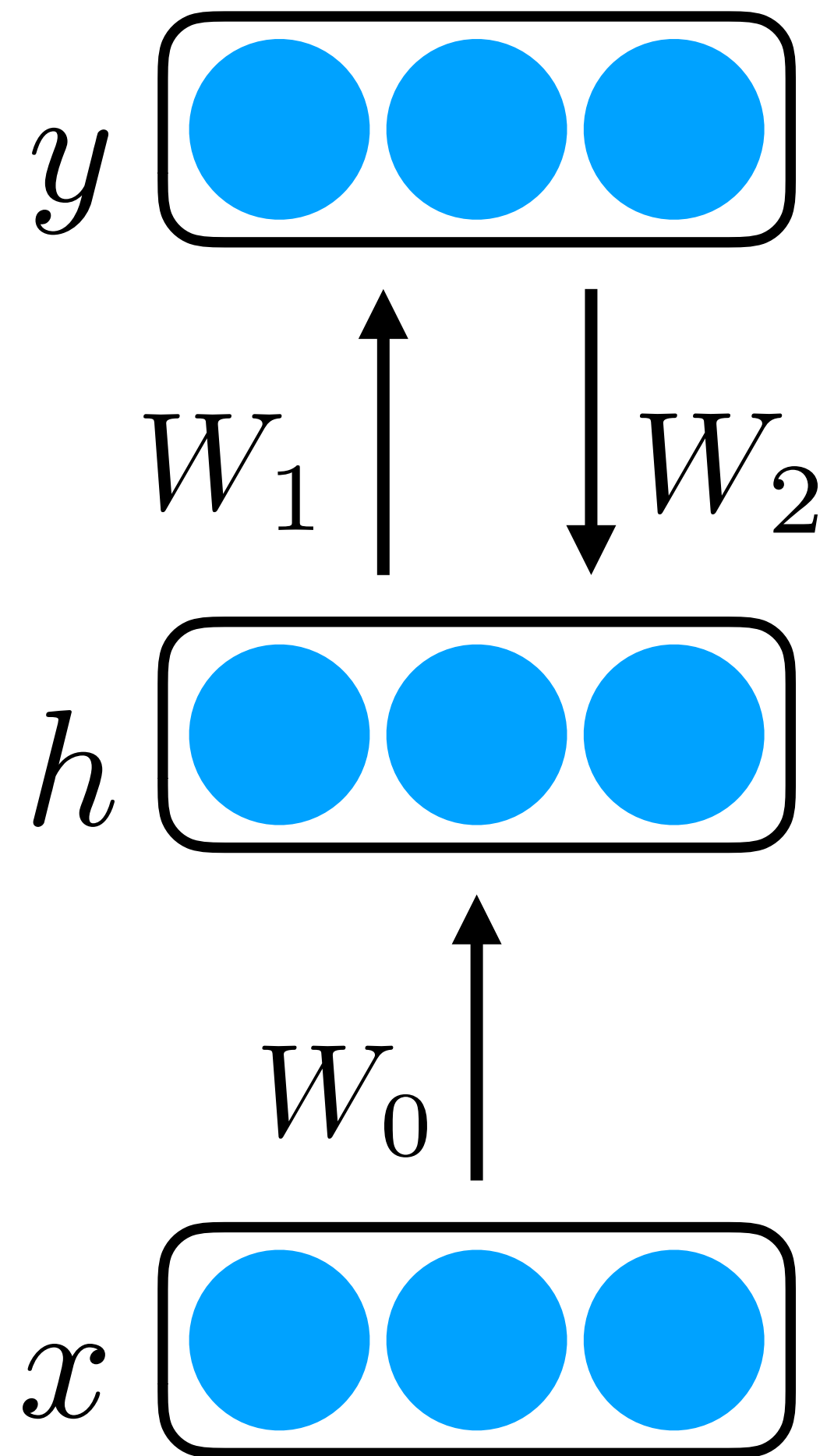


Thm: L_2 -regularized LAEs *symmetrize*.

W_2 *dynamically* aligns to W_1^T by:

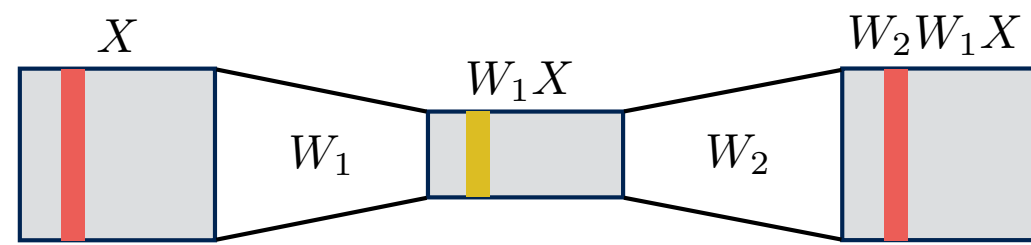
- maximizing flow of information
- minimizing energy

$$\mathcal{L}_{\text{IA}} = \mathcal{L}_{\text{pred}} + \mathcal{L}_{\text{info}} + \mathcal{L}_{\text{reg}}$$



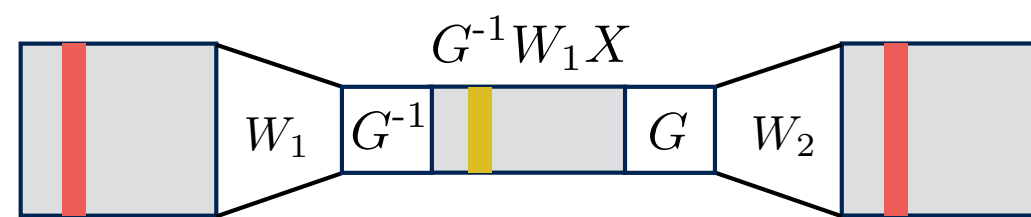
Background

A **linear autoencoder** maps $\mathbb{R}^m \rightarrow \mathbb{R}^k \rightarrow \mathbb{R}^m$.



$$\mathcal{L}(W_1, W_2) = \|X - W_2W_1X\|^2$$

LAEs learn the top principal *subspace* but **not** the principal *directions* or eigenvalues. The optimal latent representation is only defined up to a **linear map** $G \in GL_k(\mathbb{R}^m)$.



LAEs are **pseudoinverses** at all critical pts.

Regularization

We prove that L²-regularized LAEs are **transposes** at all critical points and learn the principal directions as the left singular vectors of the decoder. Define \mathcal{L}_σ by

$$\|X - W_2W_1X\|^2 + \lambda(\|W_1\|^2 + \|W_2\|^2)$$

The minima of \mathcal{L}_σ are defined up to an **orthogonal map** $O \in O_k(\mathbb{R}^m)$ by

$$W_2 = U_k(I - \lambda\Sigma_k^2)^{\frac{1}{2}}O = W_1^\top$$

where $X = U\Sigma V^\top$ and $\sigma_1^2 > \dots > \sigma_k^2 > \lambda$.

$$W_2W_1 = U_k(I - \lambda\Sigma_k^2)U_k^\top$$

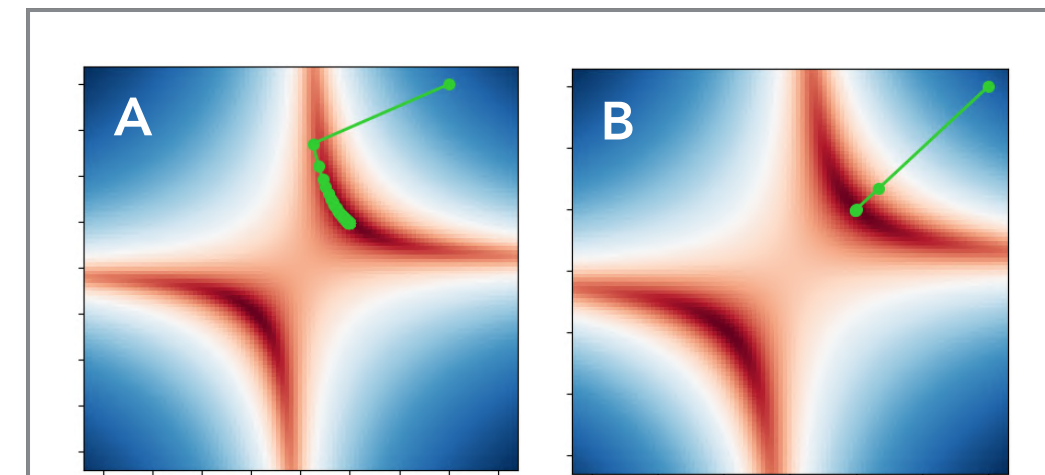
PCA Algorithms

Hence PCA is a **two-step optimization**:

1. Train L²-regularized LAE on $X \subset \mathbb{R}^{m \times n}$.
2. Apply SVD to the decoder $W_2 \subset \mathbb{R}^{m \times k}$.

Step 2 is quick. Step 1 options include:

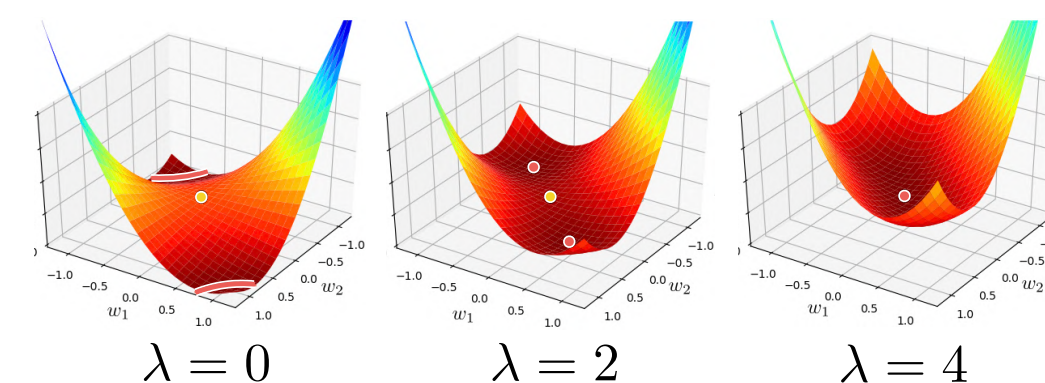
- A. Gradient descent (below).
- B. Solve for W_2 , set $W_1 = W_2^\top$, iterate.



input $X \in \mathbb{R}^{m \times n}$; $k \leq m$; $\lambda, \alpha > 0$
initialize $W_1, W_2^\top \in \mathbb{R}^{k \times m}$
while not converged
 $W_1 \leftarrow \alpha(W_2^\top(W_2W_1 - I)XX^\top + \lambda W_1)$
 $W_2 \leftarrow \alpha((W_2W_1 - I)XX^\top W_1^\top + \lambda W_2)$
 $U, \Sigma, _ = \text{SVD}(W_2)$
return $U, \lambda(I - \Sigma^2)^{-1}$

Posterior Collapse

Principal directions with eigenvalues below λ collapse as in **probabilistic PCA**.



Example of collapse for $X = [2]$.

Symmetry and Backprop

L²-reg LAEs are **symmetric** at all critical pts.

Theorem 2.1 (Transpose Theorem). *All critical points of \mathcal{L}_σ satisfy $W_1 = W_2^\top$.*

Proof. Critical points of \mathcal{L}_σ satisfy:

$$\frac{\partial \mathcal{L}_\sigma}{\partial W_1} = 2W_2^\top(W_2W_1 - I)XX^\top + 2\lambda W_1 = 0,$$

$$\frac{\partial \mathcal{L}_\sigma}{\partial W_2} = 2(W_2W_1 - I)XX^\top W_1^\top + 2\lambda W_2 = 0.$$

We first prove that the matrix

$$C = (I - W_2W_1)XX^\top$$

is positive semi-definite⁸. Rearranging $\frac{\partial \mathcal{L}_\sigma}{\partial W_2} W_2^\top$ gives

$$XX^\top(W_2W_1)^\top = (W_2W_1)XX^\top(W_2W_1)^\top + \lambda W_2W_2^\top.$$

Both terms on the right are positive semi-definite, so their sum on the left is as well and therefore

$$XX^\top(W_2W_1)^\top \succeq (W_2W_1)XX^\top(W_2W_1)^\top.$$

Cancelling $(W_2W_1)^\top$ via Lemma B.1 gives $C \succeq 0$.

We now show the difference $A = W_1 - W_2^\top$ is zero. Rearranging terms using the symmetry of C gives

$$0 = \frac{\partial \mathcal{L}_\sigma}{\partial W_1} - \frac{\partial \mathcal{L}_\sigma}{\partial W_2}^\top = 2A(C + \lambda I).$$

Since $C \succeq 0$ and $\lambda > 0$ imply $C + \lambda I \succ 0$, we conclude from

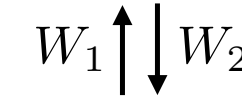
$$A(C + \lambda I)A^\top = 0$$

that $A = 0$. \square

Resolution to **weight transport problem**:



Backprop in  lacks W_1^\top because neurons go one way.



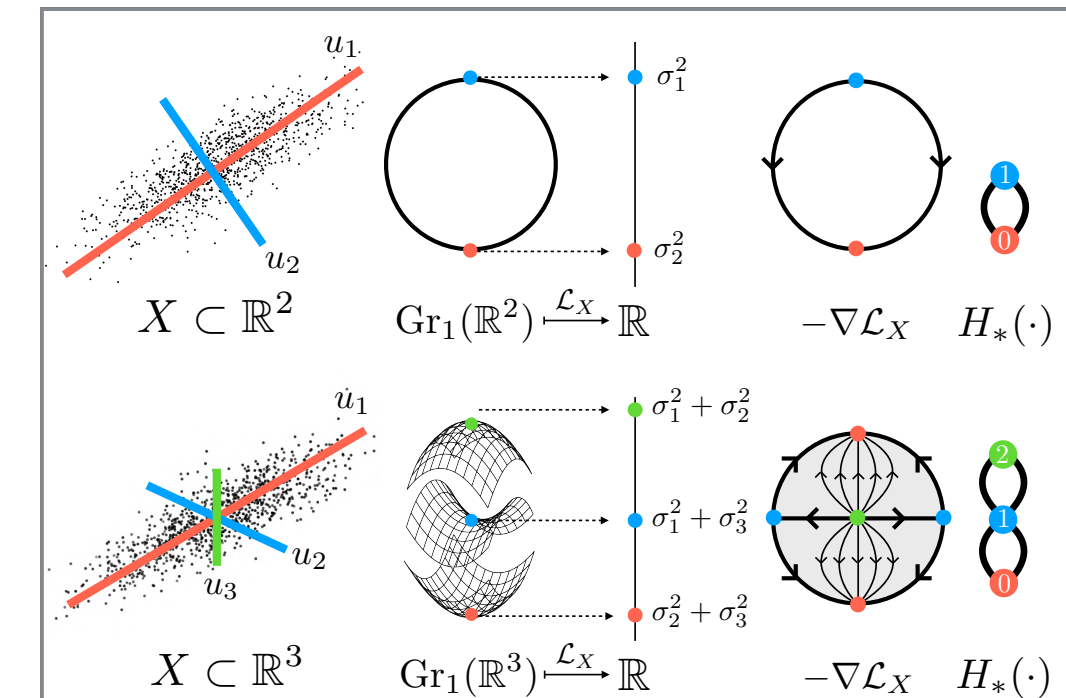
Learn as W_2 by maximizing flow of info and minimizing energy.

Algebraic Topology

We **smoothly parameterize** the critical manifolds of LAEs with several forms of regularization via one elementary proof.

We factor the loss as a **Morse function** on the Grassmannian to reveal the dynamics near and between critical manifolds.

Morse homology suggests principles and algorithms for deep learning.



Theorem 4.4 (Curvature Theorem). *In local coordinates near any point on the critical manifold indexed by \mathcal{I} , all three losses take the form of a standard degenerate saddle with $d_{\mathcal{I}} + (k - \ell)(m - \ell)$ descending directions.*

- \mathcal{L} and \mathcal{L}_π have $k\ell$ flat directions.
- \mathcal{L}_σ has $k\ell - \binom{\ell+1}{2}$ flat directions.

The remaining directions are ascending.

Theorem E.1. \mathcal{L}_X is an \mathbb{F}_2 -perfect Morse function. Its critical points are the rank- k principal subspaces.

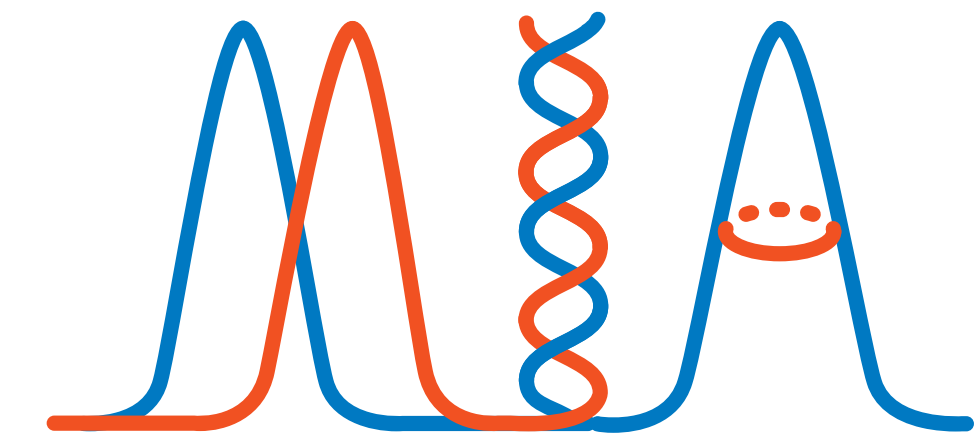
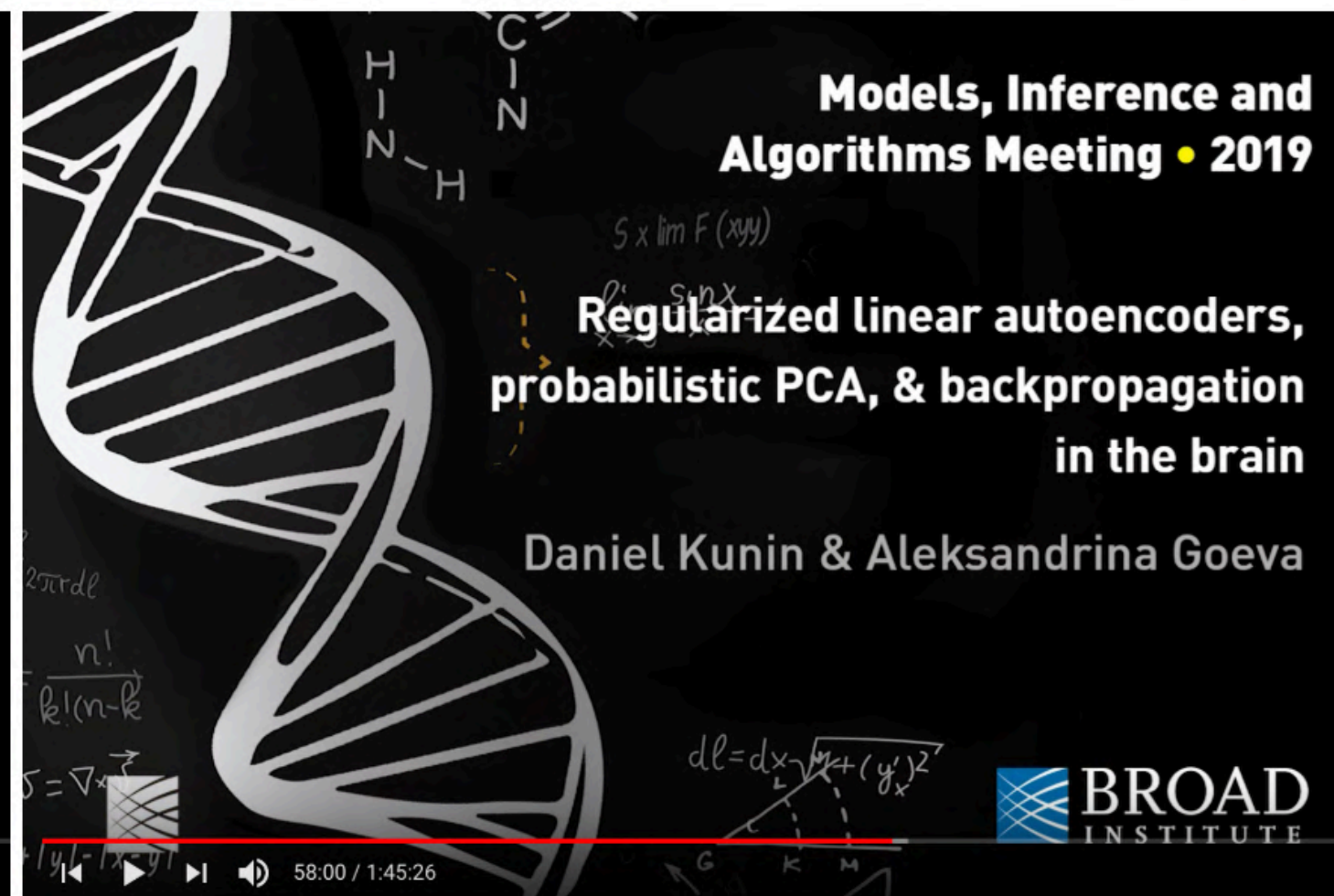
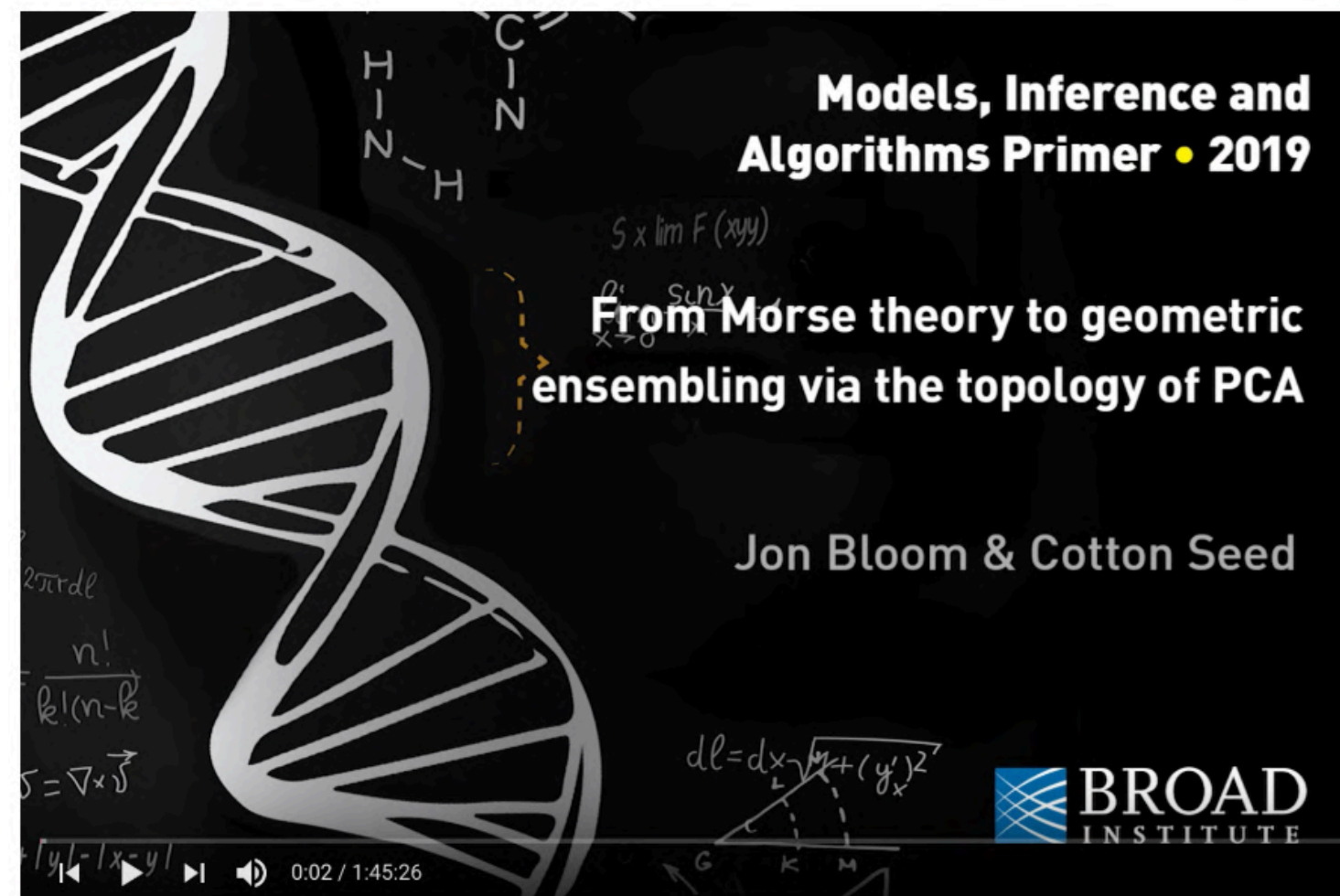
Proof. Consider the commutative diagram

$$\begin{array}{ccc} V_k(\mathbb{R}^m) & \xrightarrow{\pi: O \mapsto \text{Im}(OO^\top)} & \text{Gr}_k(\mathbb{R}^m) \\ \downarrow \iota: O \mapsto (O^\top, O) & & \downarrow \mathcal{L}_X \\ \mathbb{R}^{k \times m} \times \mathbb{R}^{m \times k} & \xrightarrow{\mathcal{L}} & \mathbb{R} \end{array} \quad (10)$$

In the last decade, **biology** has been transformed by the ability to perturb and measure biological systems at massive scale.

However, new ideas in **ML** are needed to translate biomedical data into a *mechanistic understanding* of biology and disease.

Biology and **ML** are poised to powerfully advance one another.



broadinstitute.org/mia