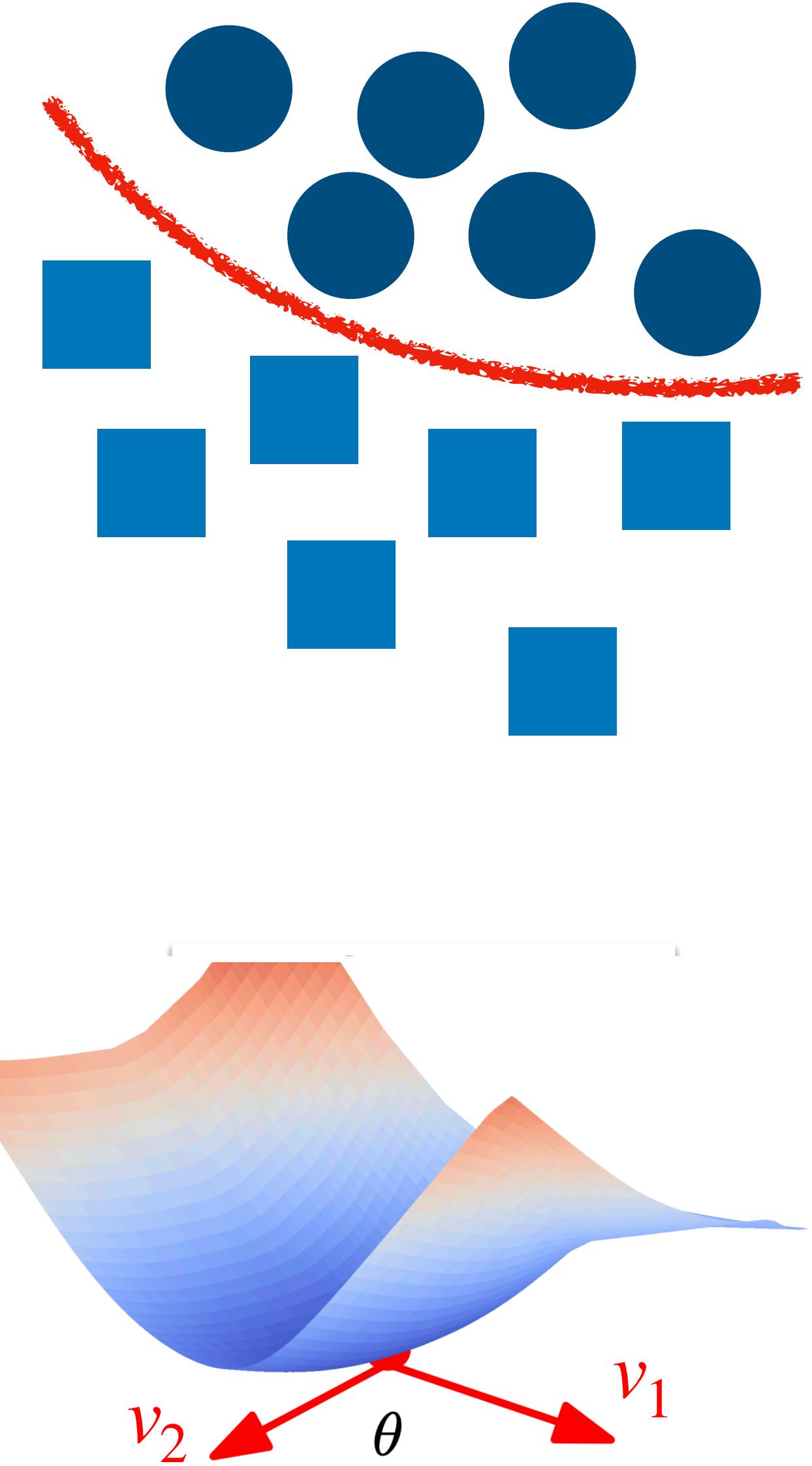


Unveiling the Hessian's Connection to the Decision Boundary

Anna Dawid



Machine learning and statistical physics back together again, Aug 9, 2023



Amazing collaborators

Les Houches Summer School open problem



Mahalakshmi Sabanayagam
TU Munich



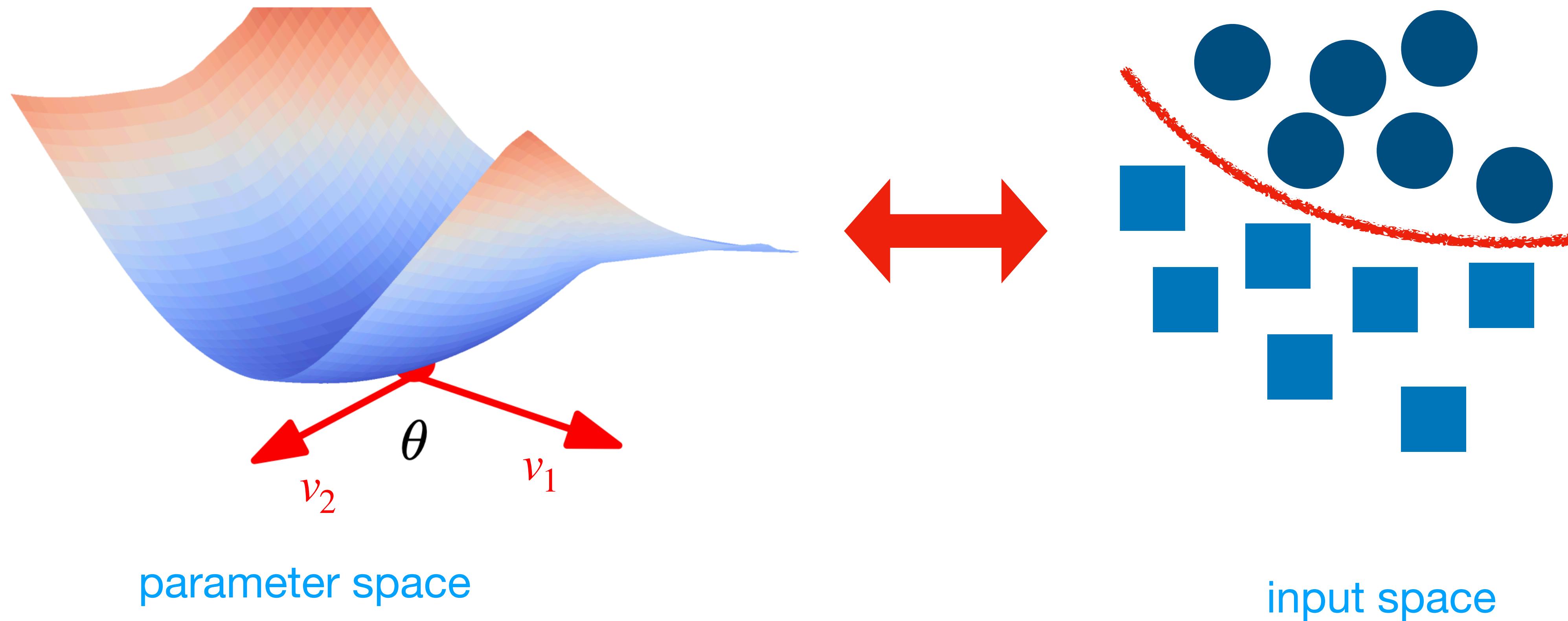
Freya Behrens
EPFL



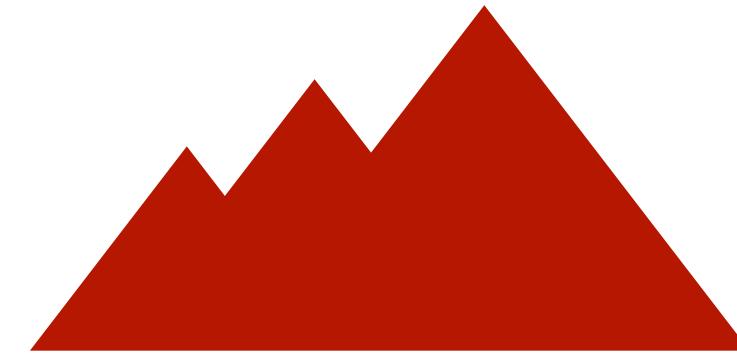
Urte Adomaityte
King's College London

Take-home message

The top Hessian eigenvectors encode the decision boundary



Outline



Empirical mess
in the Land of the Loss



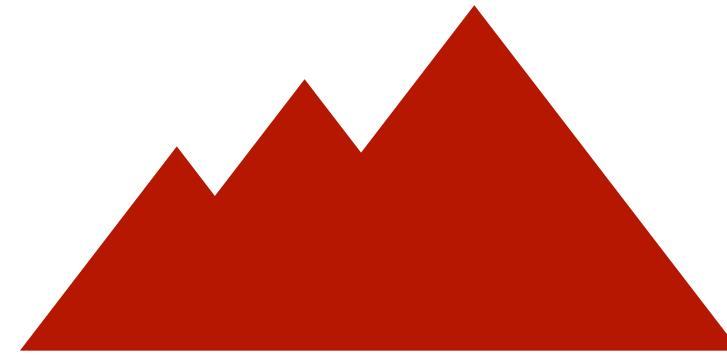
Hessian and the decision
boundary: toy example



Selected interesting
consequences

- measure of the decision boundary complexity
- margin estimation
- perspectives for new generalization measure

Outline



Empirical mess
in the Land of the Loss

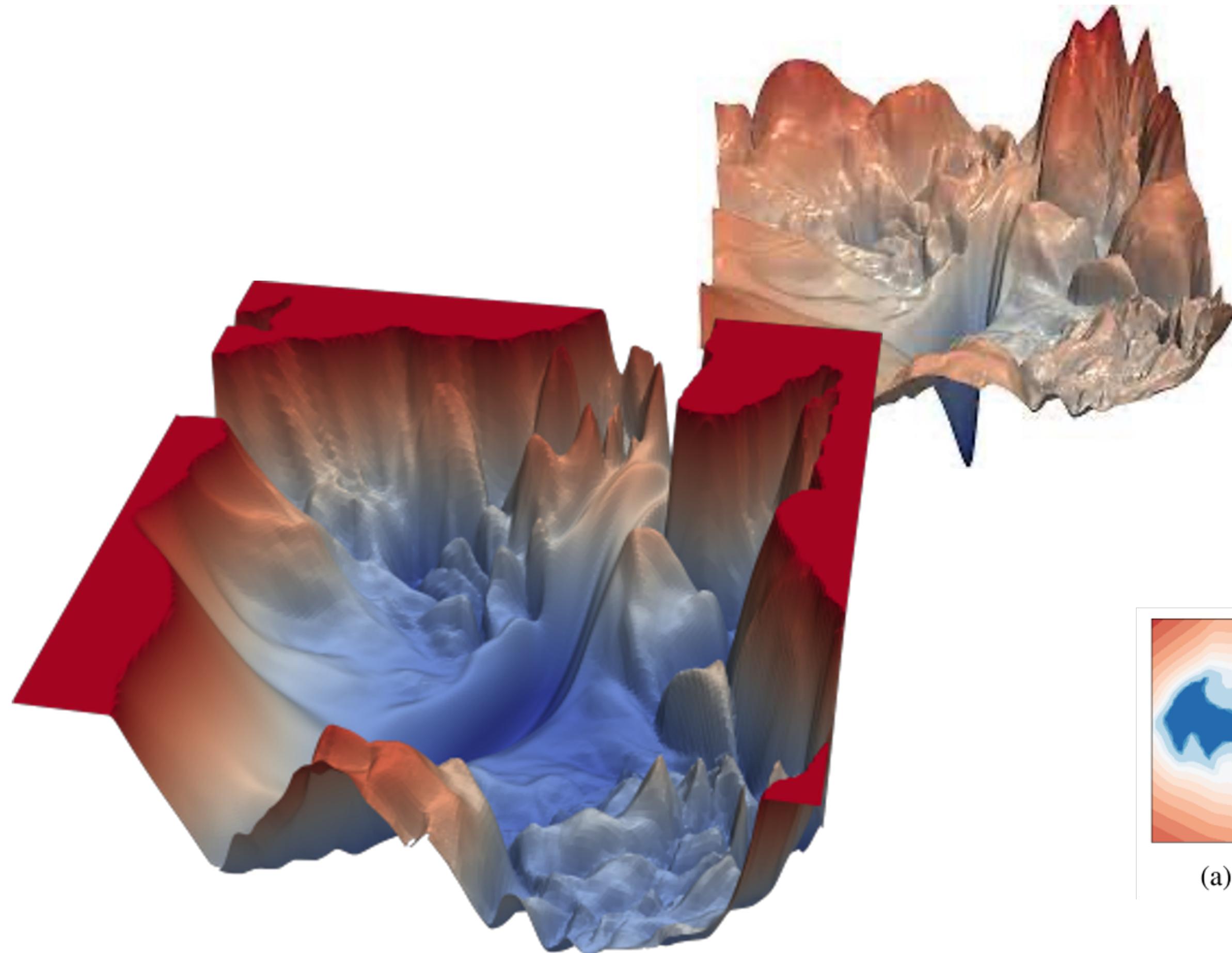


Hessian and the decision
boundary: toy example

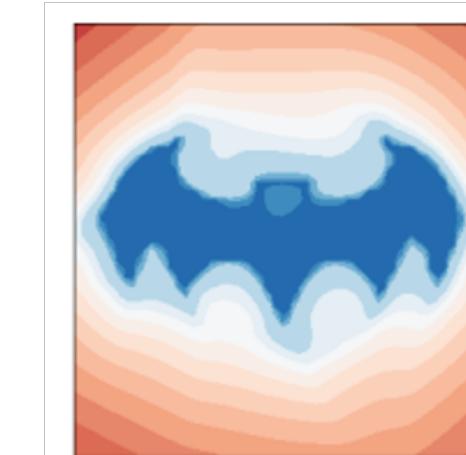


Selected interesting
consequences

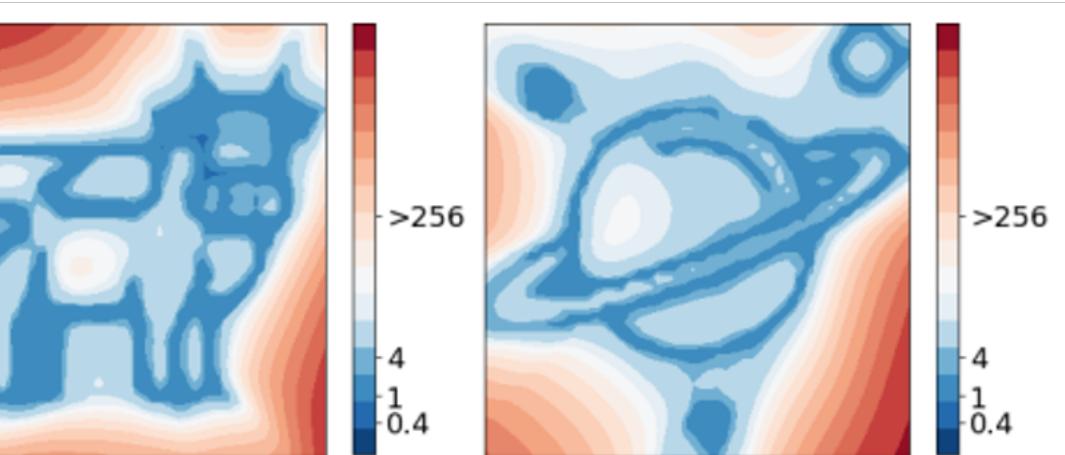
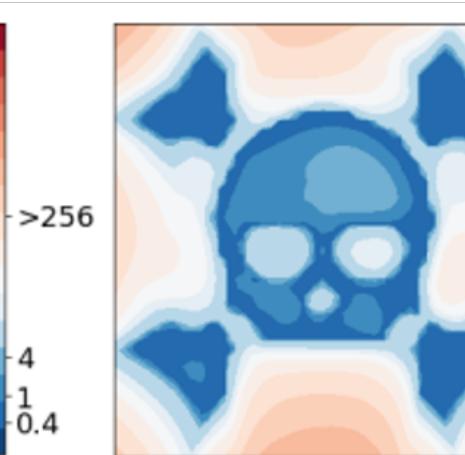
Non-convex loss landscapes of deep networks?



one can empirically find arbitrary 2D binary patterns
inside loss surfaces of popular neural networks!

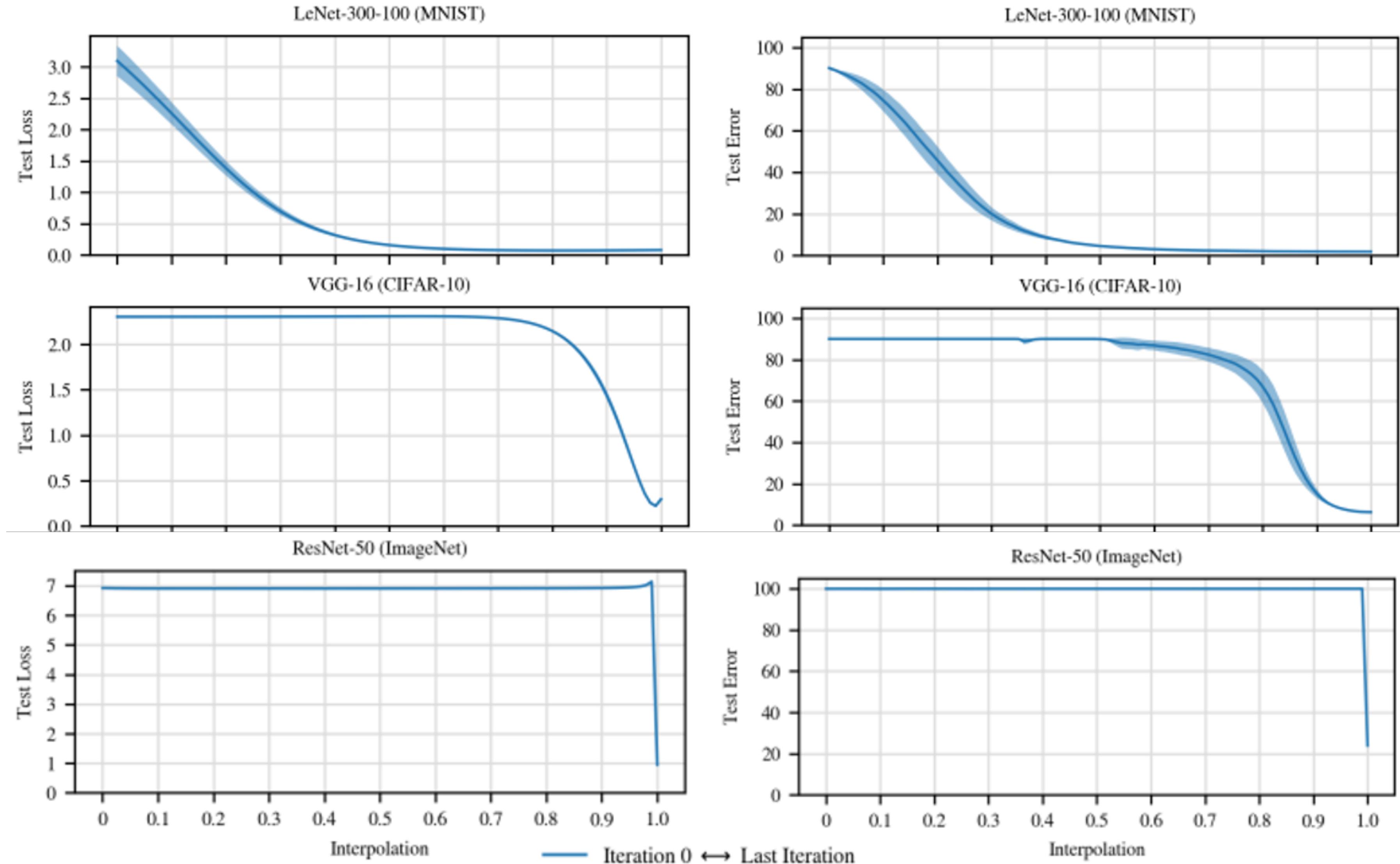


(a) Loss surface on FashionMNIST dataset

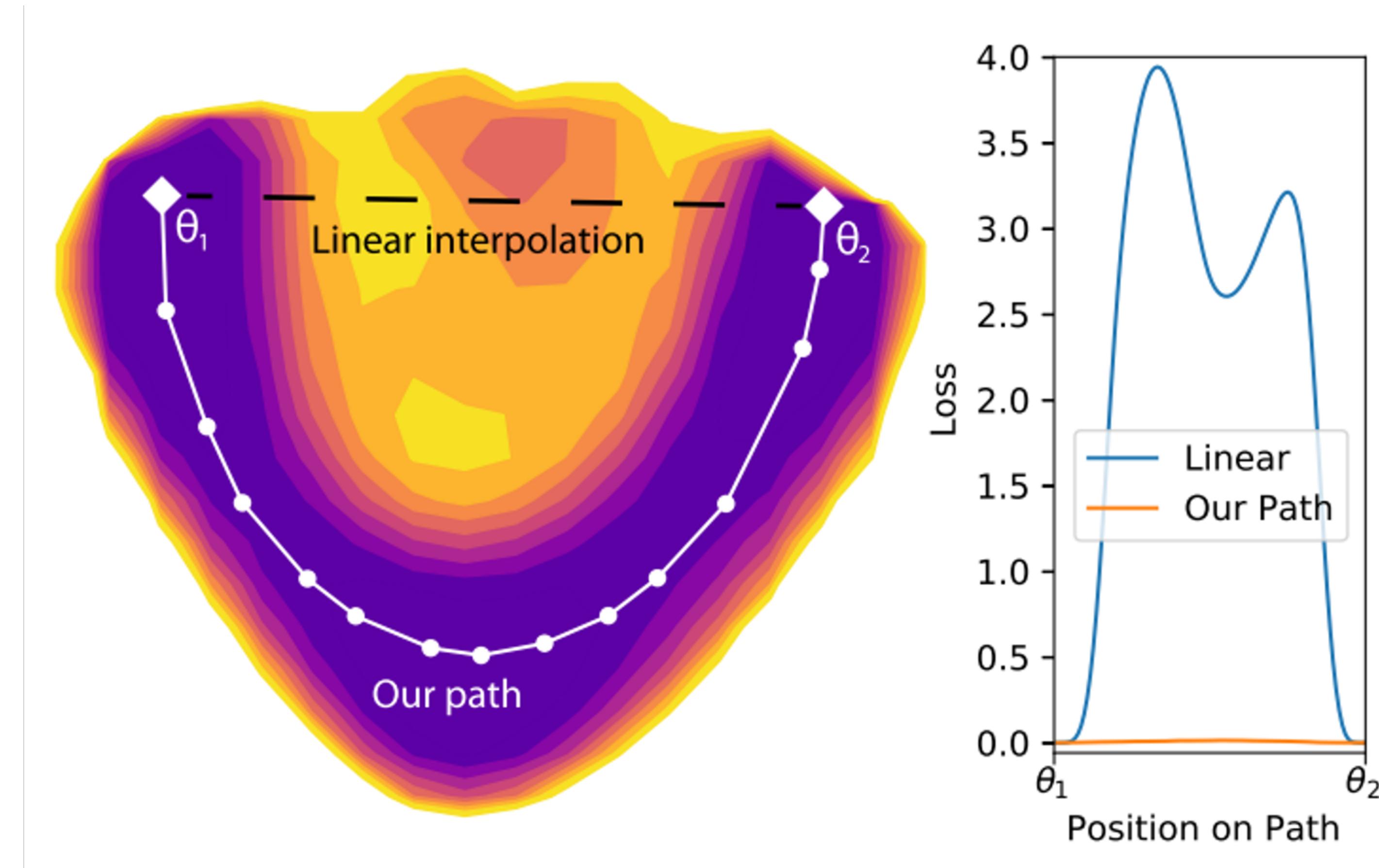


(b) Loss surface on CIFAR10 dataset

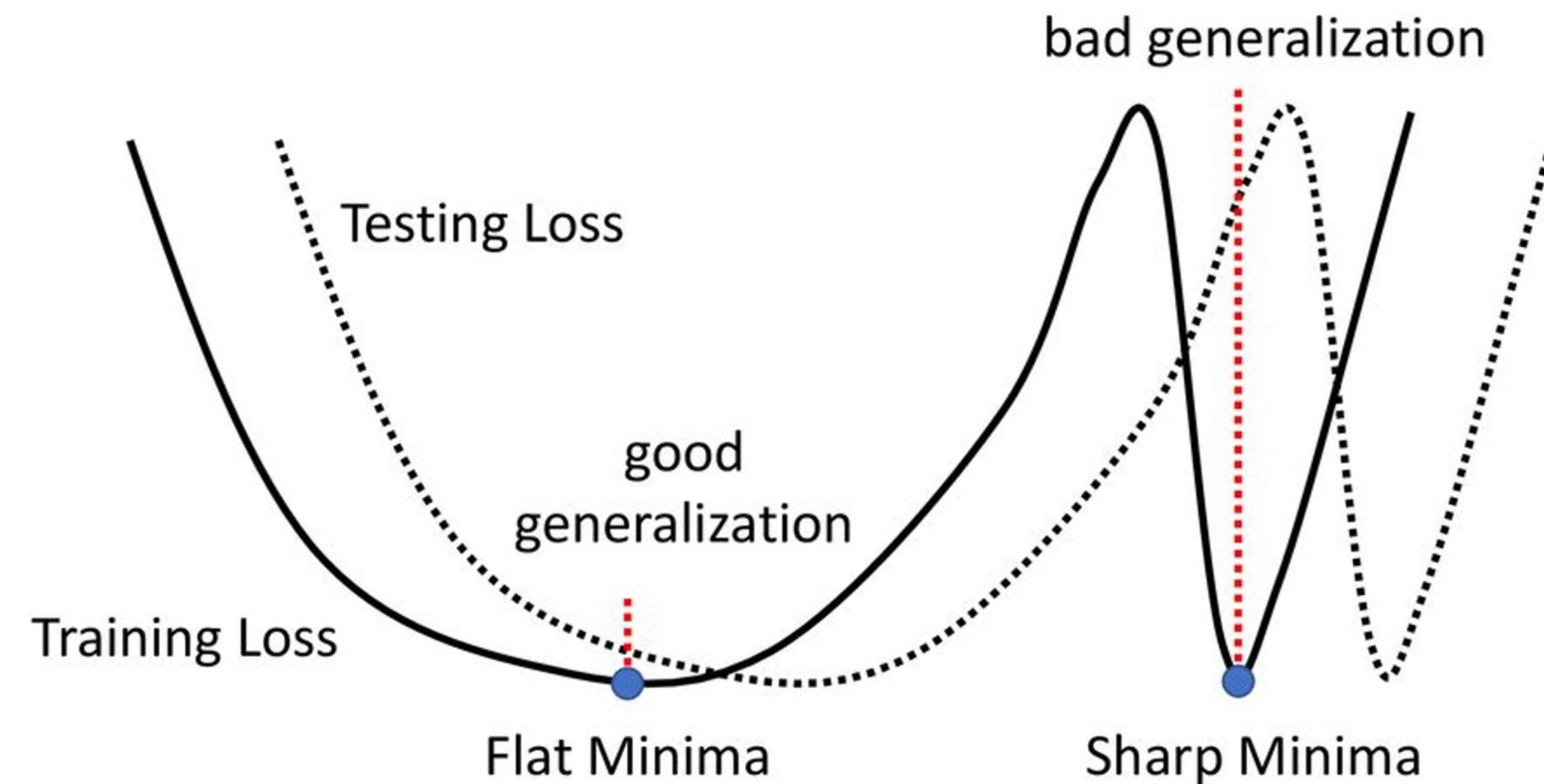
But... Path from a proper initialization to a minimum



But... Minima are connected



Optimizers reach flat minima that generalize better?



Optimizers reach flat minima that generalize better?



Many works: **Yes! Flatter minima generalize better**

Optimizers reach flat minima that generalize better?



Many works: **Yes! Flatter minima generalize better**

Also many works:



- standard (Hessian trace, spectra norm) flatness measures do **not** correlate well with generalization
- they are **not** reparametrization invariant!

$$\text{ReLU}(\alpha x) = \alpha \text{ReLU}(x) \quad (\text{non-negative homogeneous})$$

$$\text{ReLU}(x \cdot (\alpha \theta_1)) \theta_2 = \text{ReLU}(x \cdot \theta_1) \cdot (\alpha \theta_2)$$

1st-layer $\alpha \theta_1$
and 2nd-layer θ_2

1st-layer θ_1
and 2nd-layer $\alpha \theta_2$

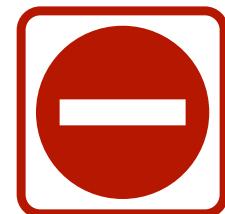
arbitrarily sharp minima with
effectively identical predictions!

Optimizers reach flat minima that generalize better?



Many works: **Yes! Flatter minima generalize better**

Also many works:



- standard (Hessian trace, spectra norm) flatness measures do **not** correlate well with generalization
- they are **not** reparametrization invariant!



Also many works:

Let's fix the flatness measure!

$$\text{ReLU}(\alpha x) = \alpha \text{ReLU}(x) \quad (\text{non-negative homogeneous})$$

$$\text{ReLU}(x \cdot (\alpha \theta_1)) \theta_2 = \text{ReLU}(x \cdot \theta_1) \cdot (\alpha \theta_2)$$

1st-layer $\alpha \theta_1$
and 2nd-layer θ_2

1st-layer θ_1
and 2nd-layer $\alpha \theta_2$

arbitrarily sharp minima with
effectively identical predictions!

Optimizers reach flat minima that generalize better?



Many works: **Yes! Flatter minima generalize better**

Also many works:



- standard (Hessian trace, spectra norm) flatness measures do **not** correlate well with generalization
- they are **not** reparametrization invariant!



Also many works:

Let's **fix the flatness measure!**

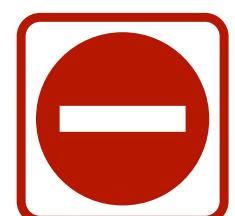
$$\text{ReLU}(\alpha x) = \alpha \text{ReLU}(x) \quad (\text{non-negative homogeneous})$$

$$\text{ReLU}(x \cdot (\alpha \theta_1)) \theta_2 = \text{ReLU}(x \cdot \theta_1) \cdot (\alpha \theta_2)$$

1st-layer $\alpha \theta_1$
and 2nd-layer θ_2

1st-layer θ_1
and 2nd-layer $\alpha \theta_2$

arbitrarily sharp minima with
effectively identical predictions!



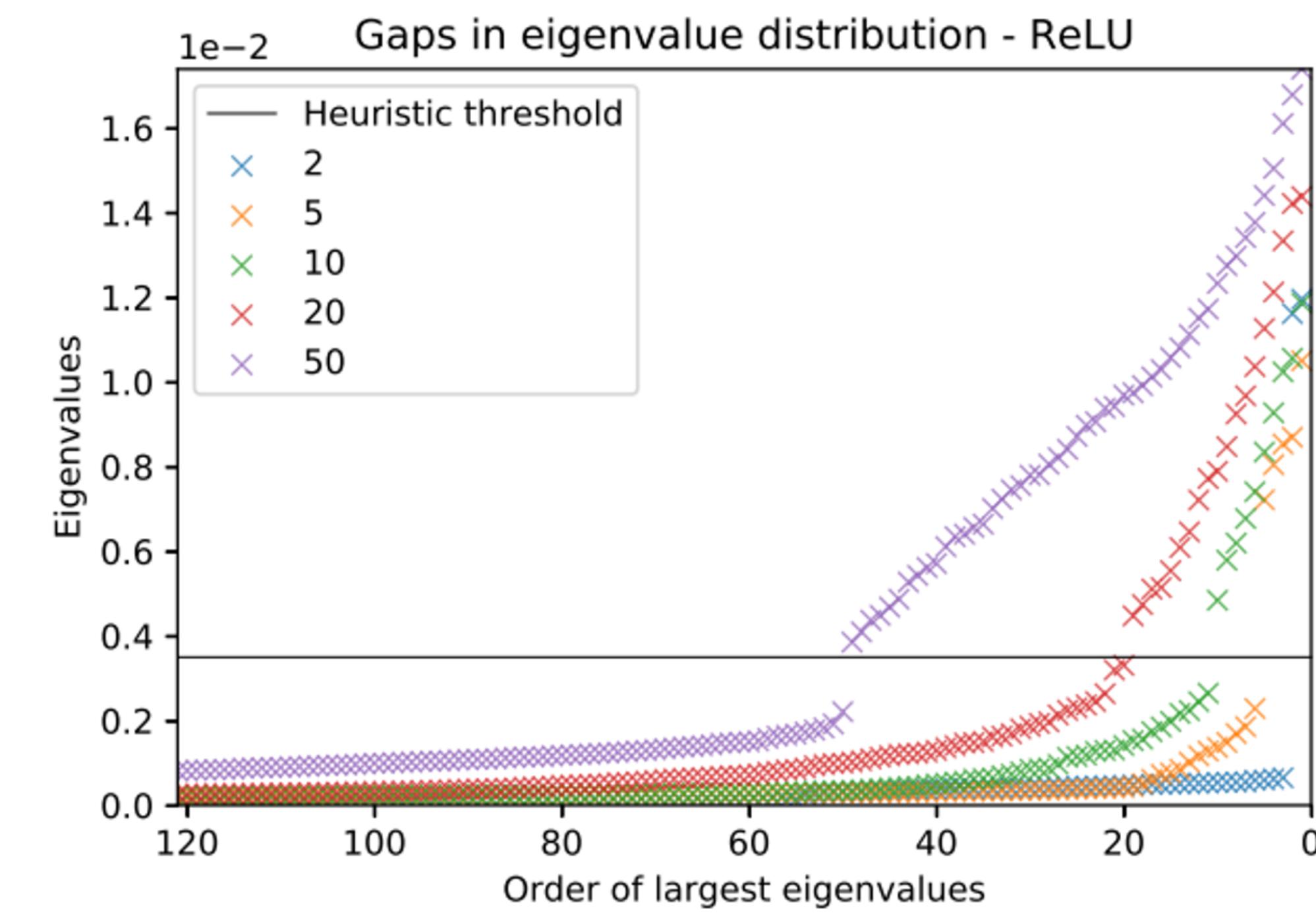
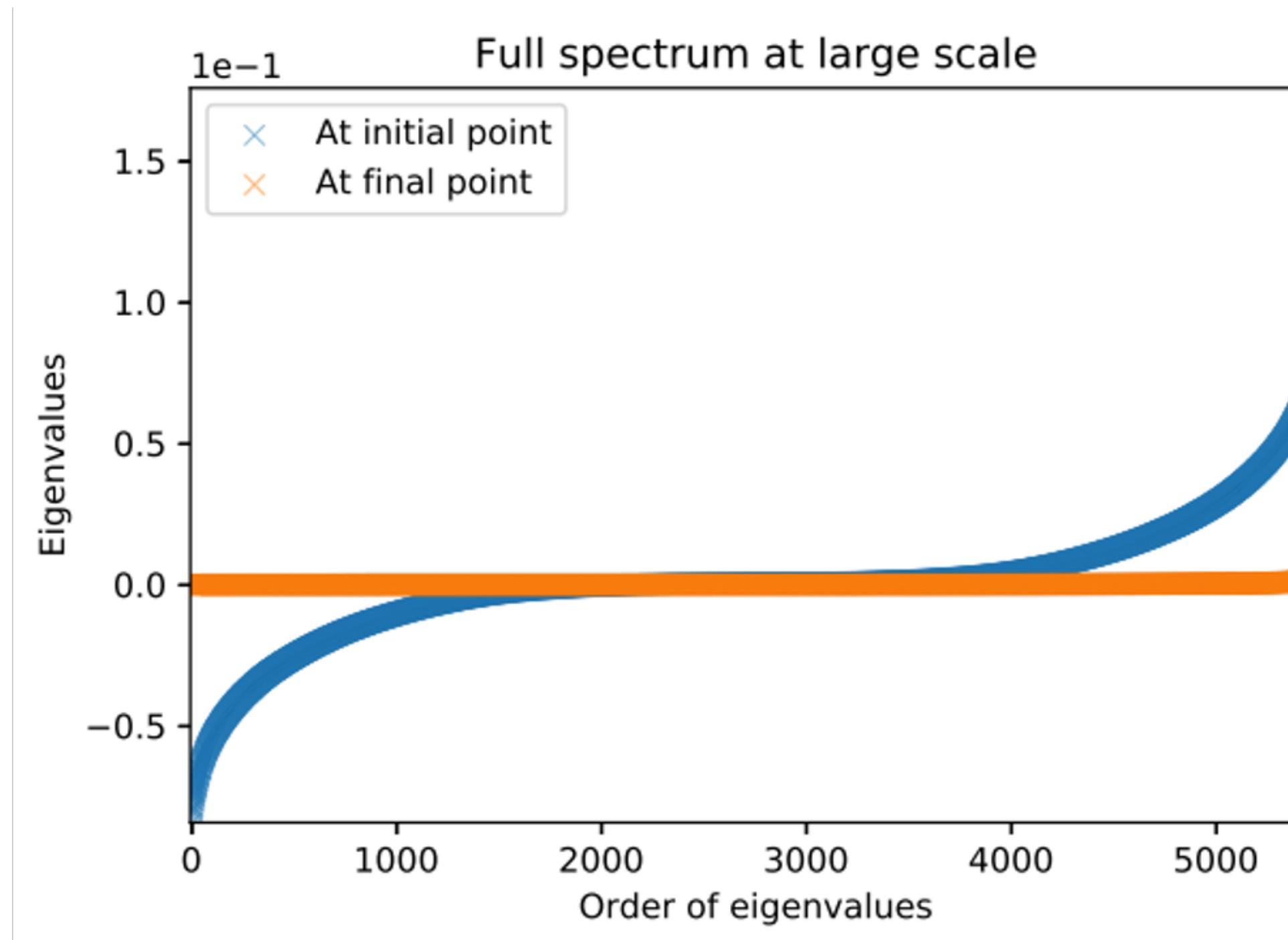
Also many works:

No, they also don't work!

They **correlate rather with hyperparameters** than the generalization itself.

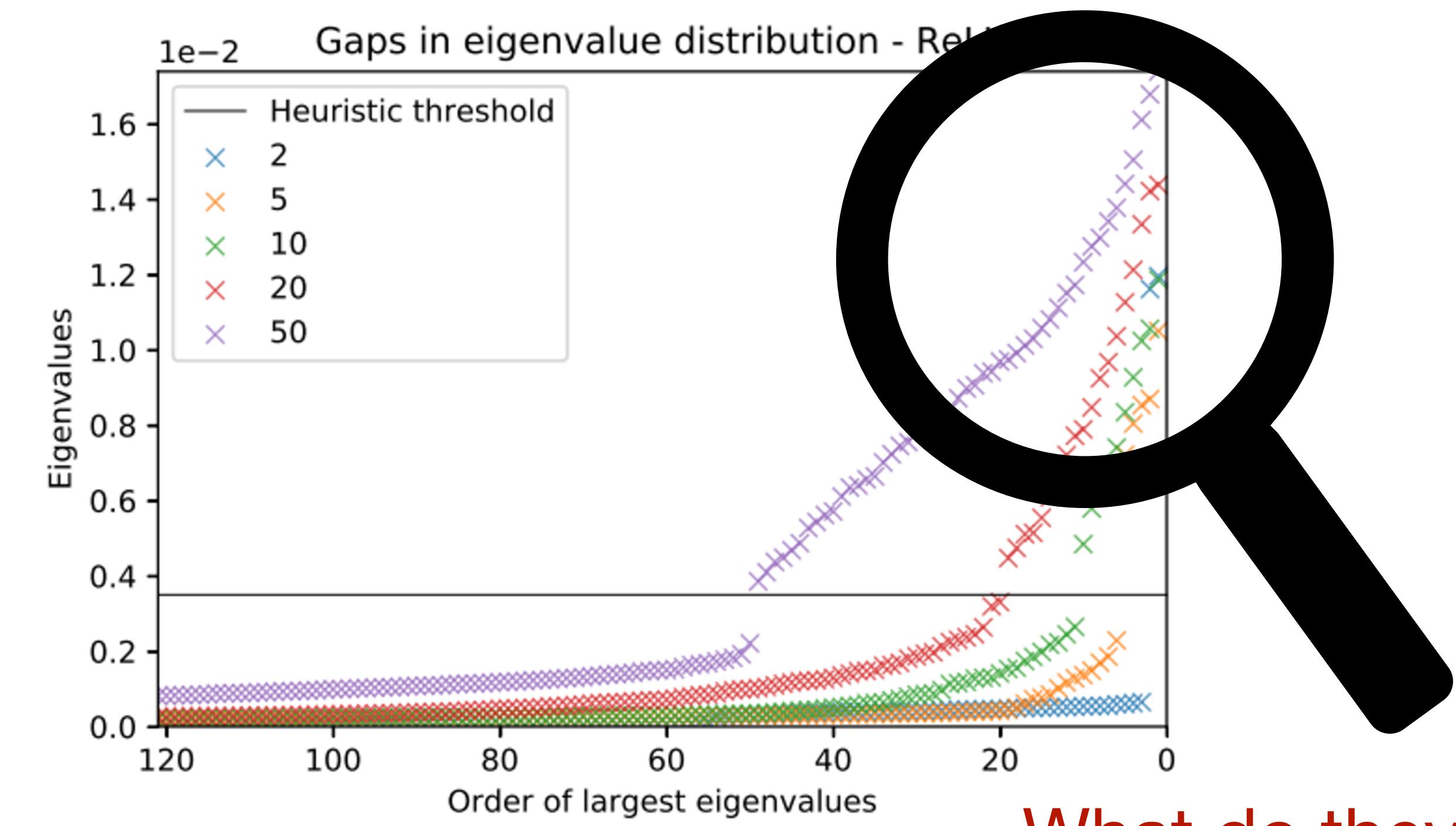
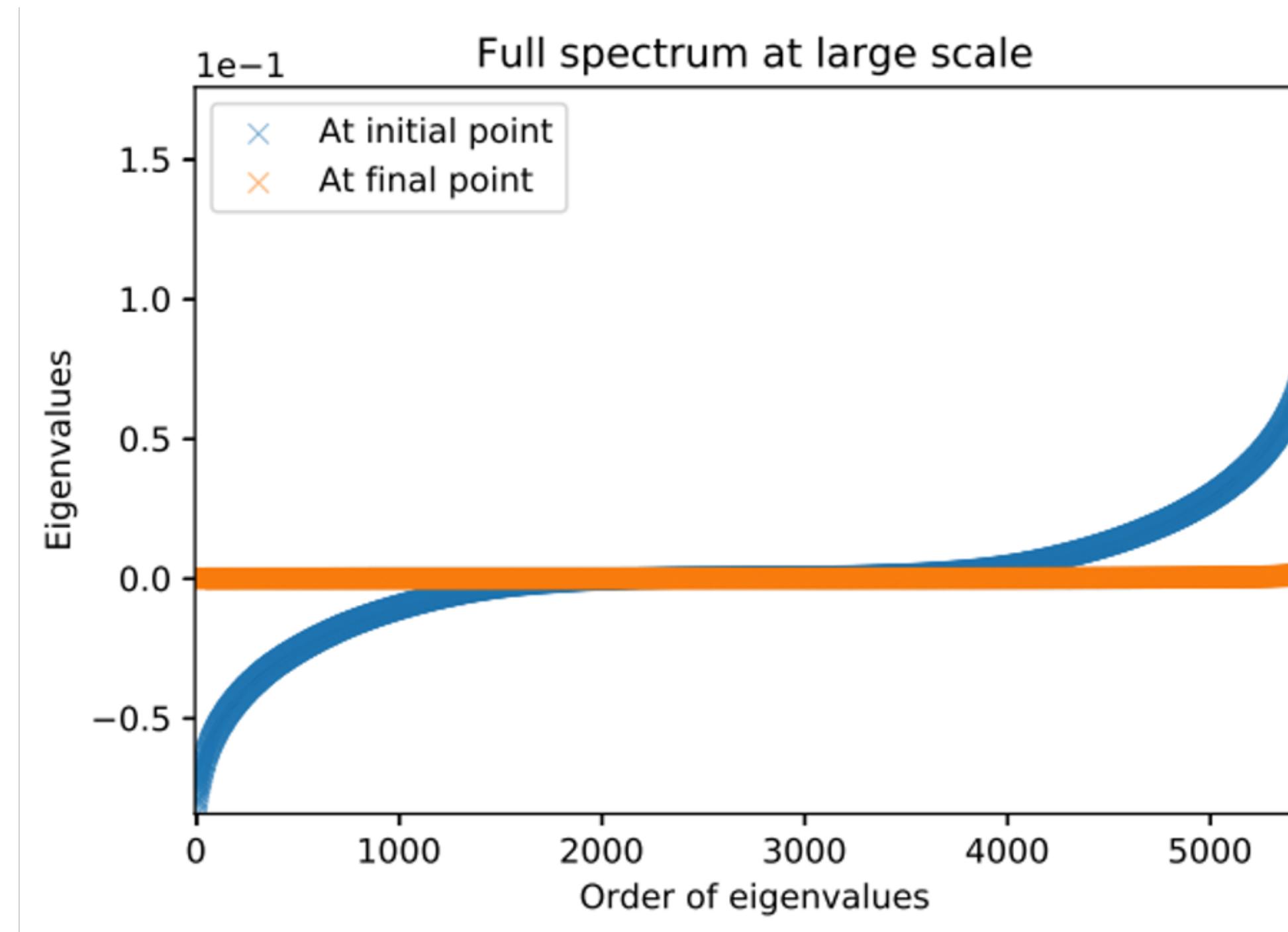
Hessian has universal properties across deep learning setups a.k.a. the community does agree on something!

$$H_{ij} = \left. \frac{\partial^2 \mathcal{L}(\mathcal{D}, \theta)}{\partial \theta_i \partial \theta_j} \right|_{\theta=\theta^*}$$



Hessian has universal properties across deep learning setups a.k.a. the community does agree on something!

$$H_{ij} = \left. \frac{\partial^2 \mathcal{L}(\mathcal{D}, \theta)}{\partial \theta_i \partial \theta_j} \right|_{\theta=\theta^*}$$

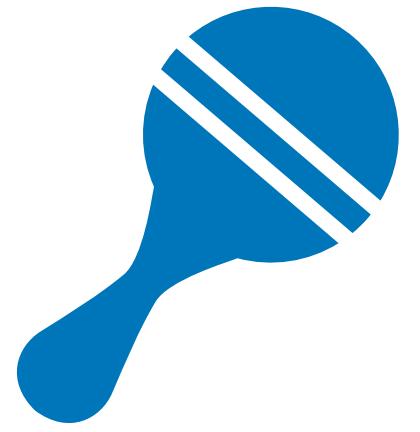


What do they
encode?

Outline



Empirical mess
in the Land of the Loss

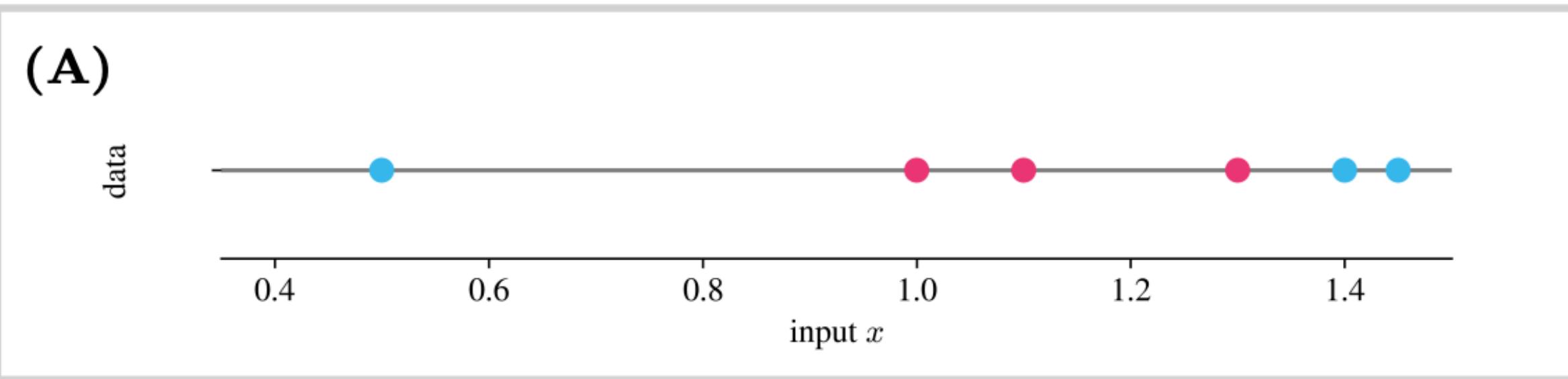


Hessian and the decision
boundary: toy example

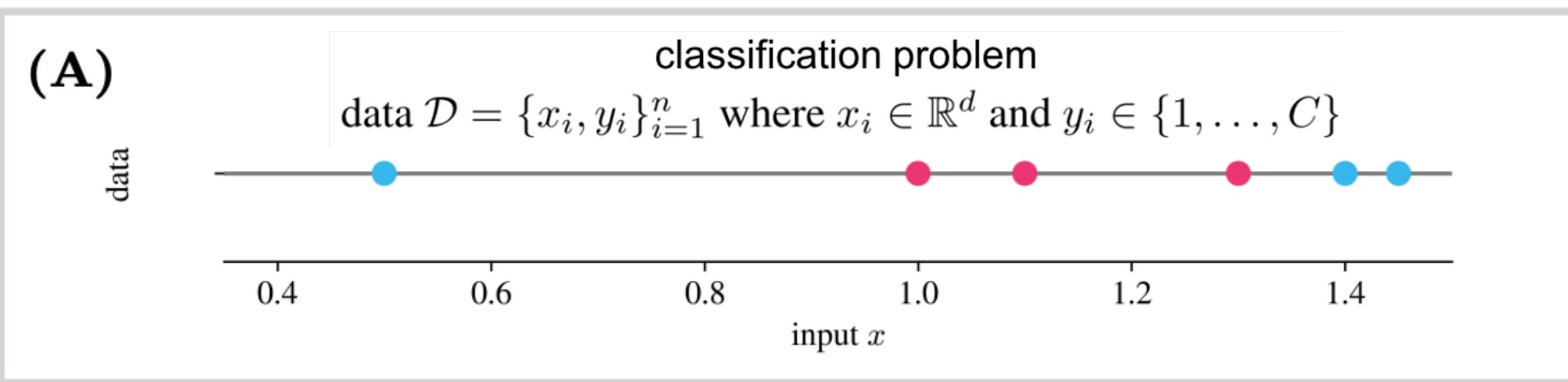


Selected interesting
consequences

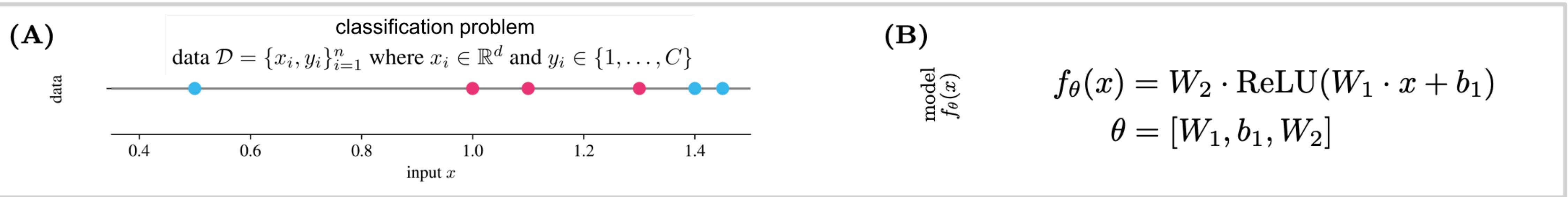
Toy example



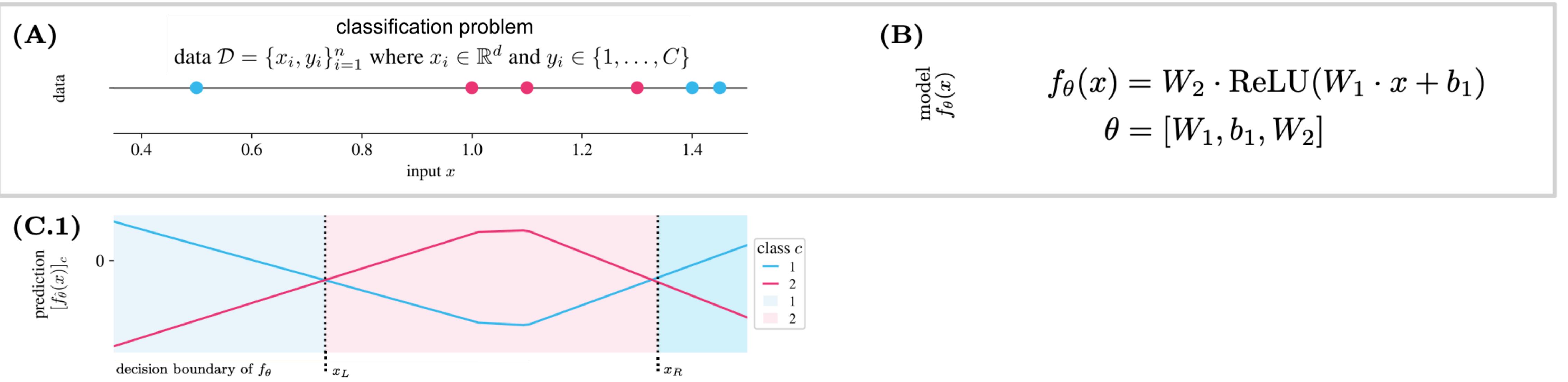
Toy example



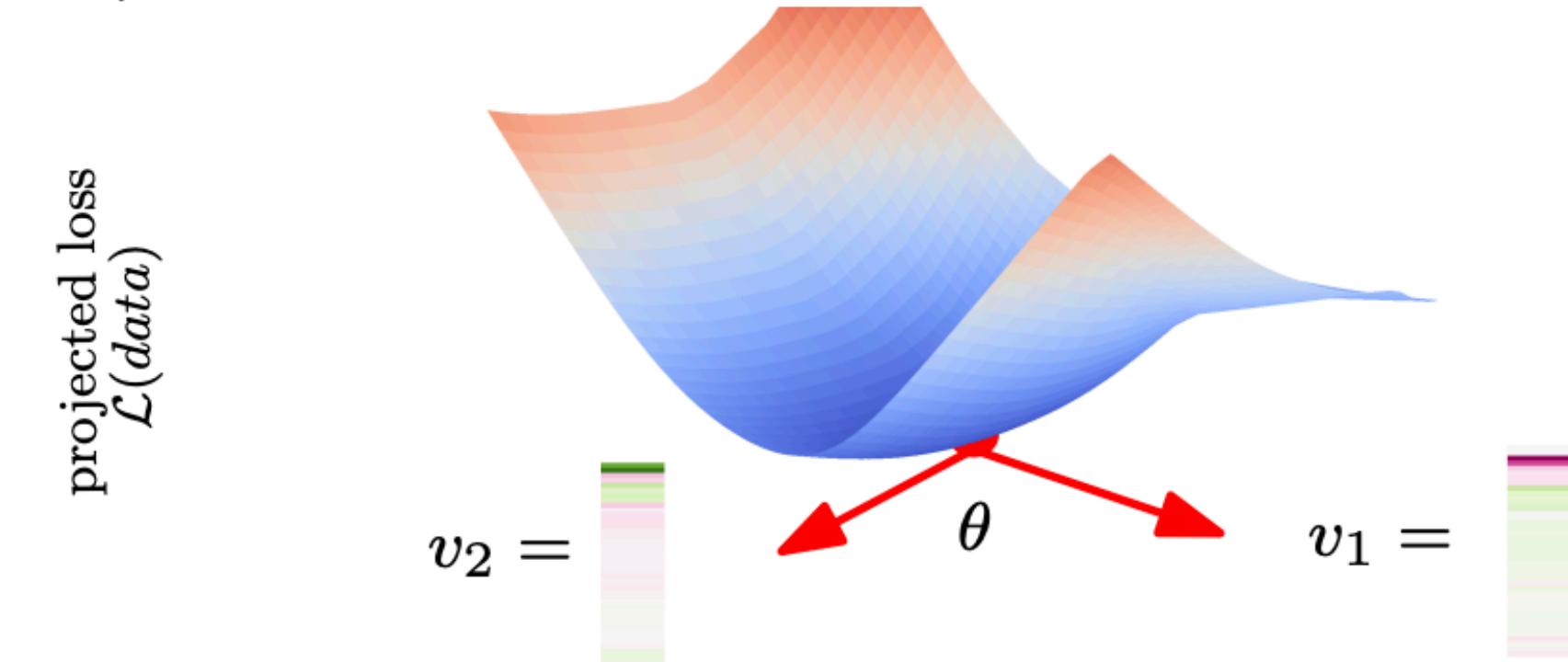
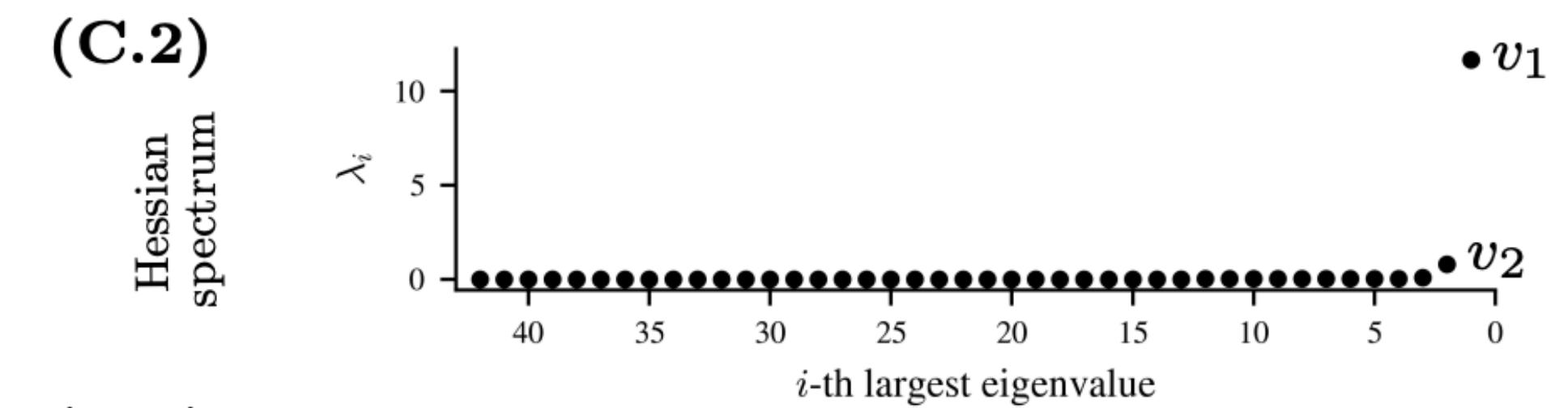
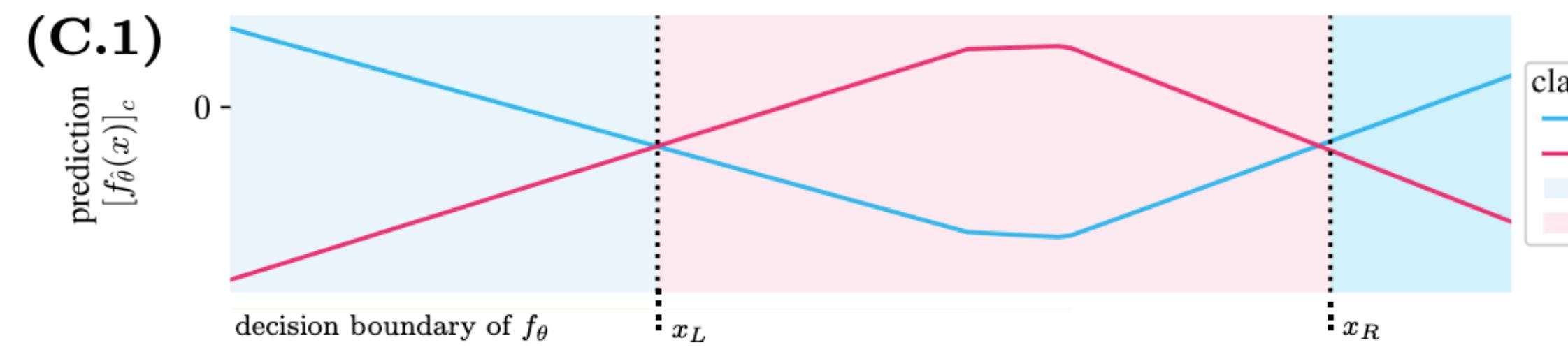
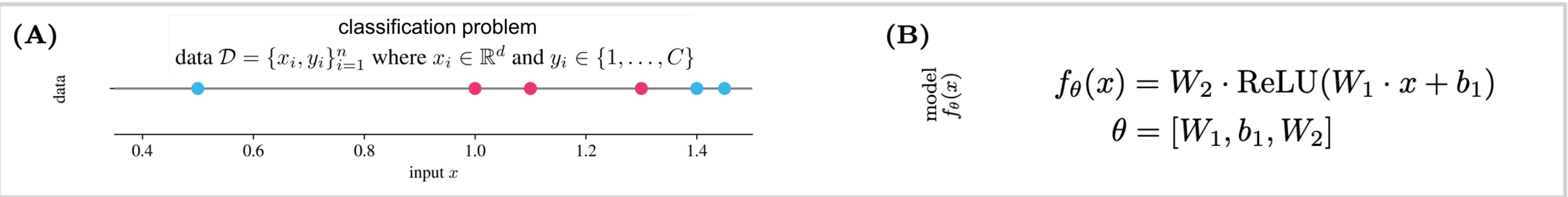
Toy example



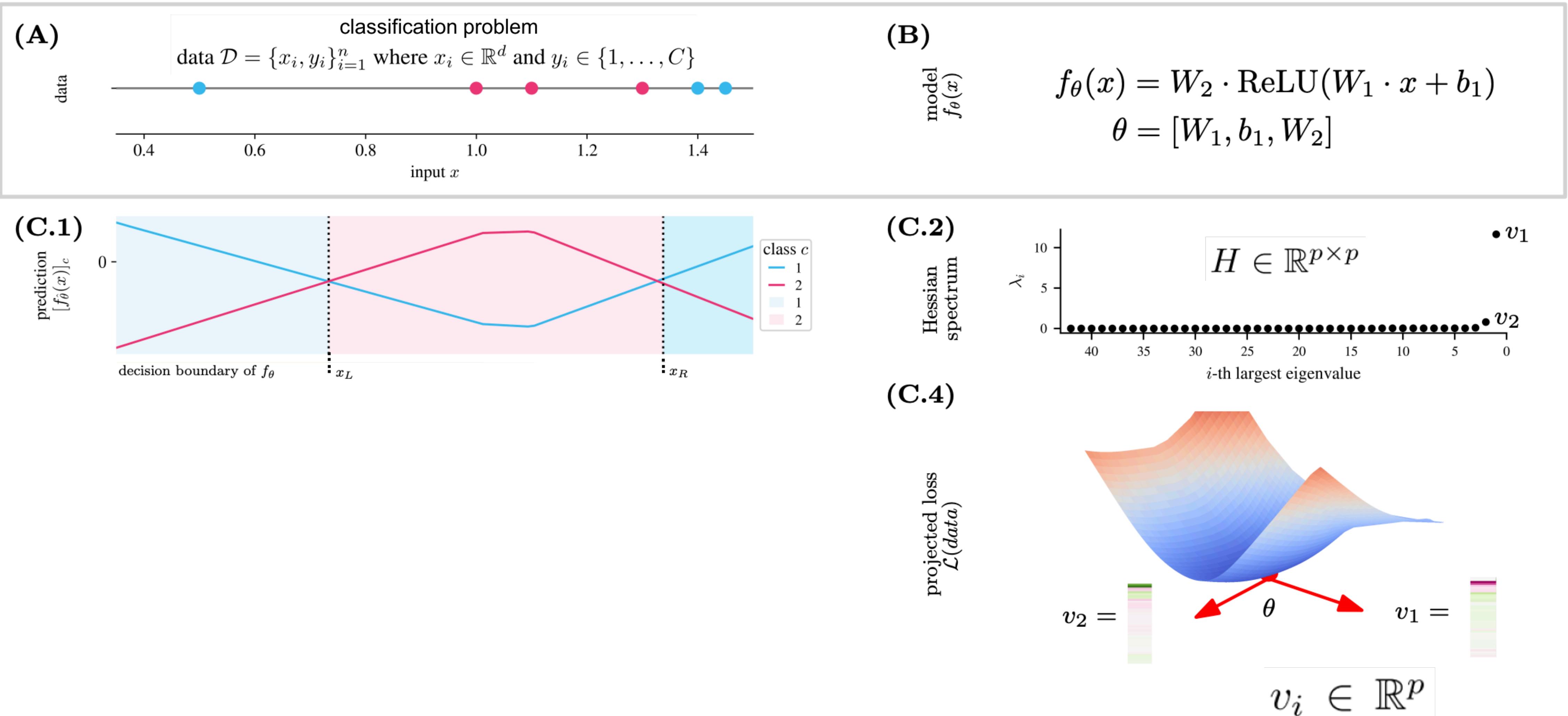
Toy example



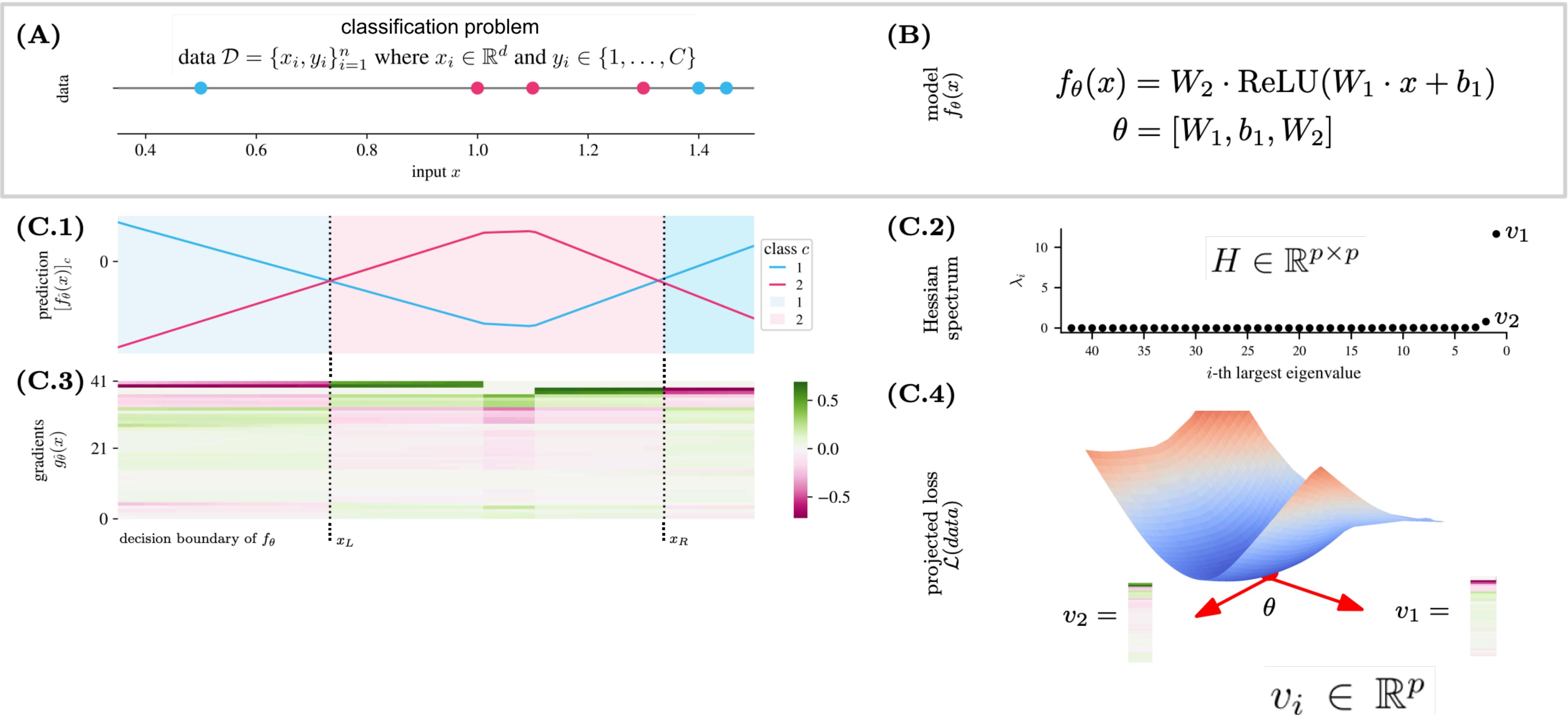
Toy example



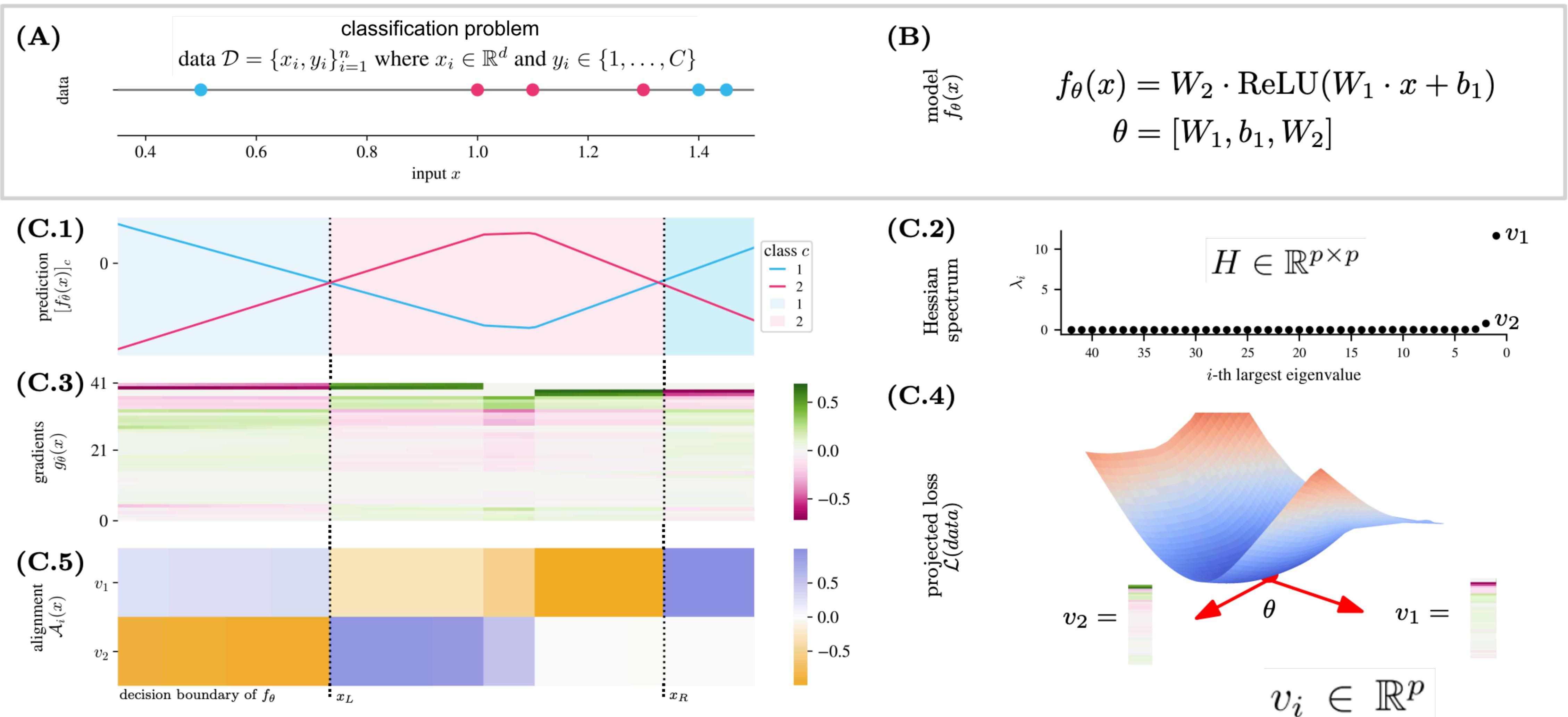
Toy example



Toy example



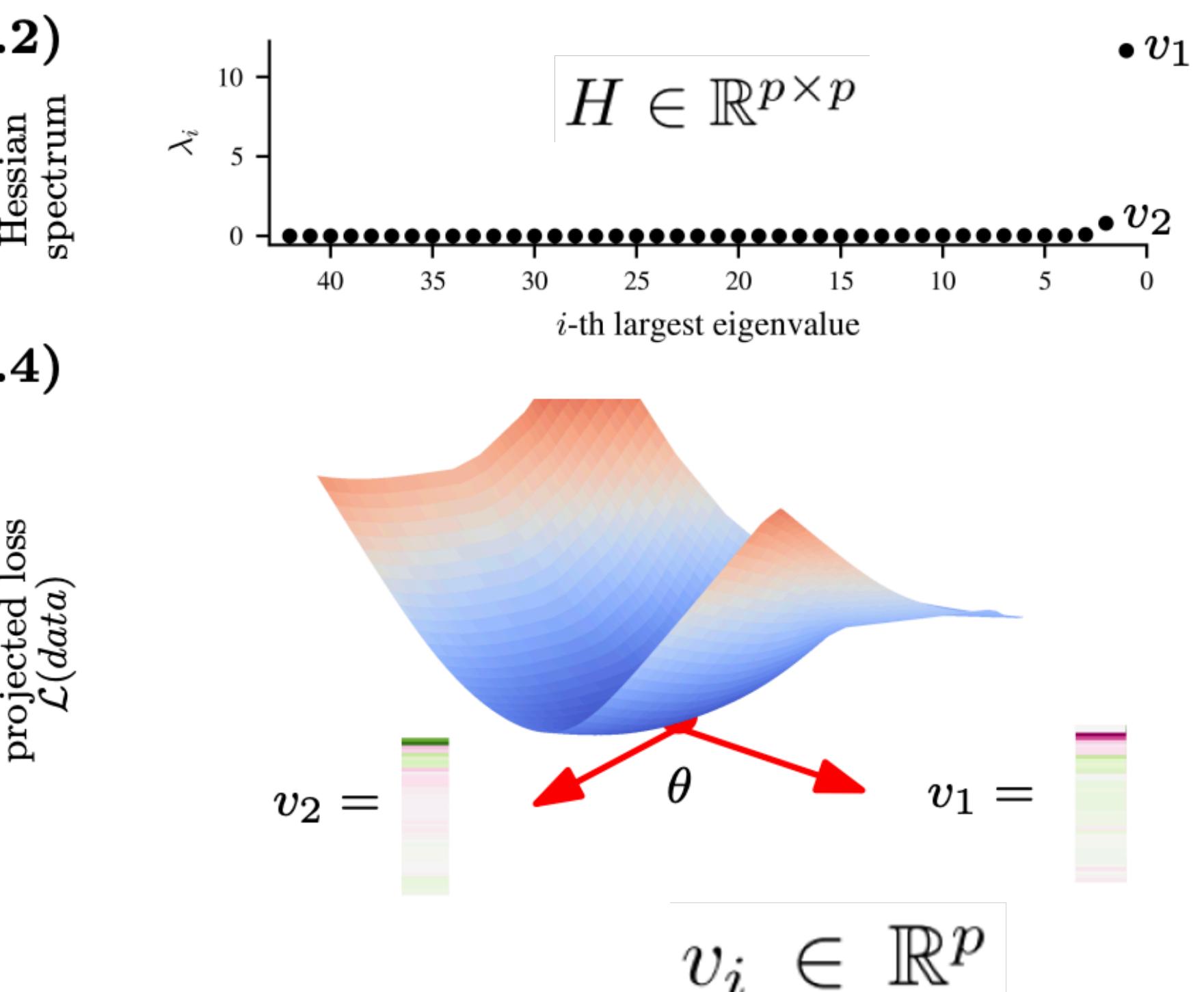
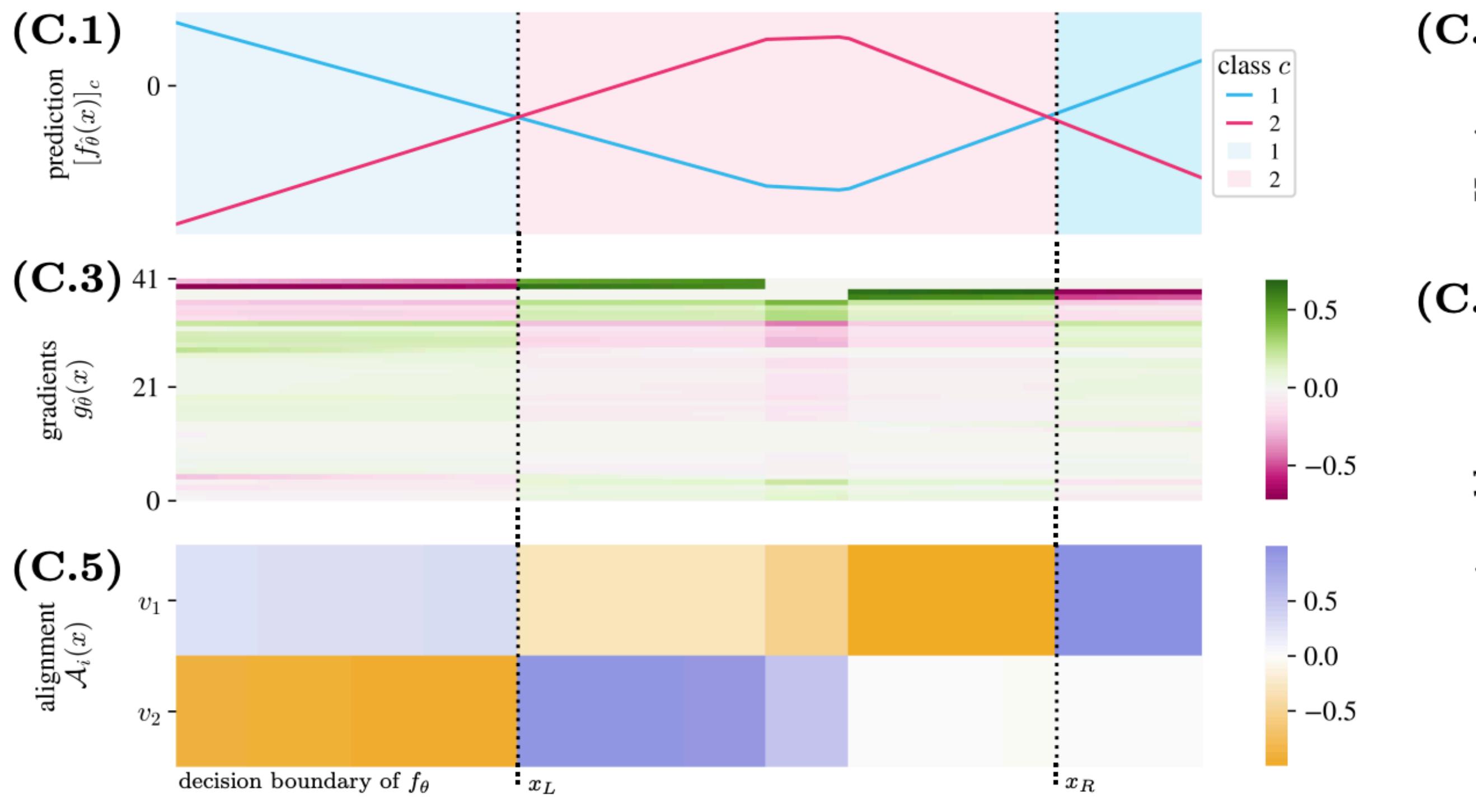
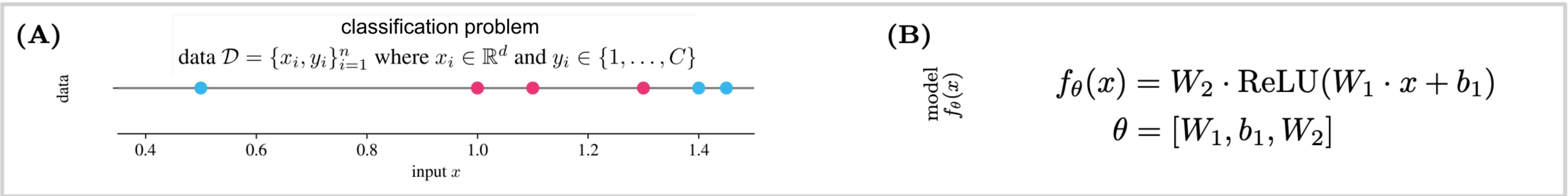
Toy example



Toy example

Alignment between the i -th Hessian eigenvector and the reinforcing gradient:

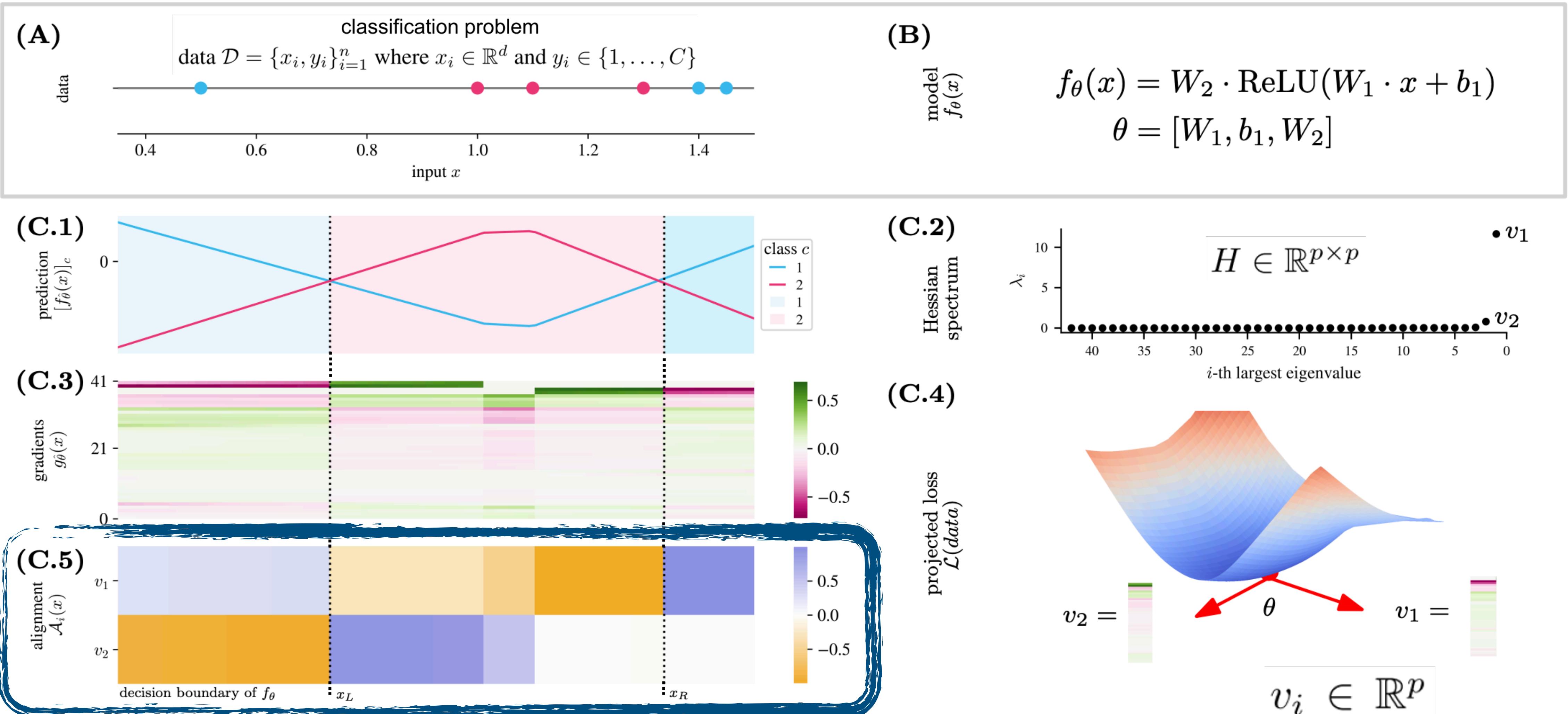
$$\mathcal{A}_i(x) = \frac{\langle g_\theta(x), v_i \rangle}{\|g_\theta(x)\| \|v_i\|}$$



Toy example

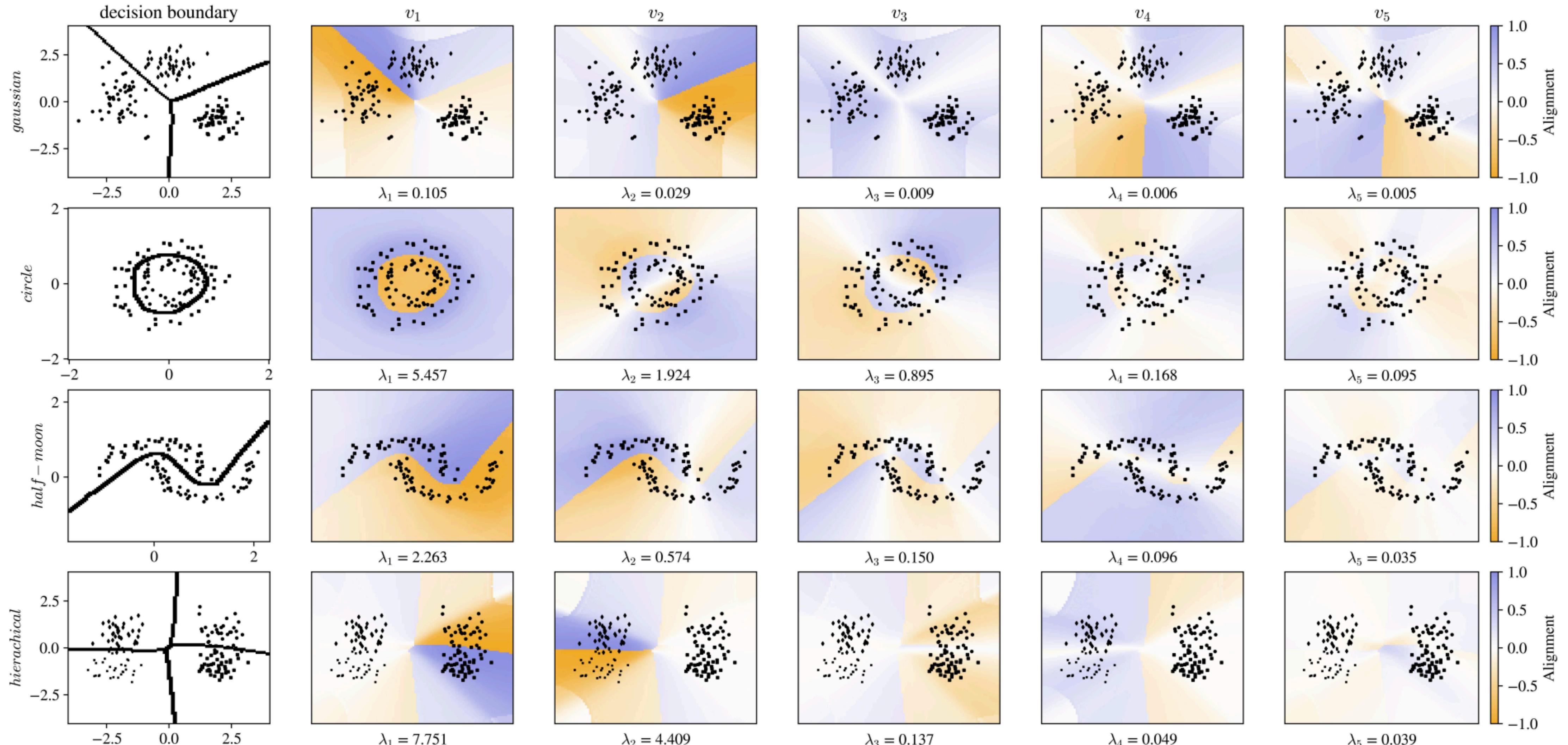
Alignment between the i -th Hessian eigenvector and the reinforcing gradient:

$$\mathcal{A}_i(x) = \frac{\langle g_\theta(x), v_i \rangle}{\|g_\theta(x)\| \|v_i\|}$$



2D datasets

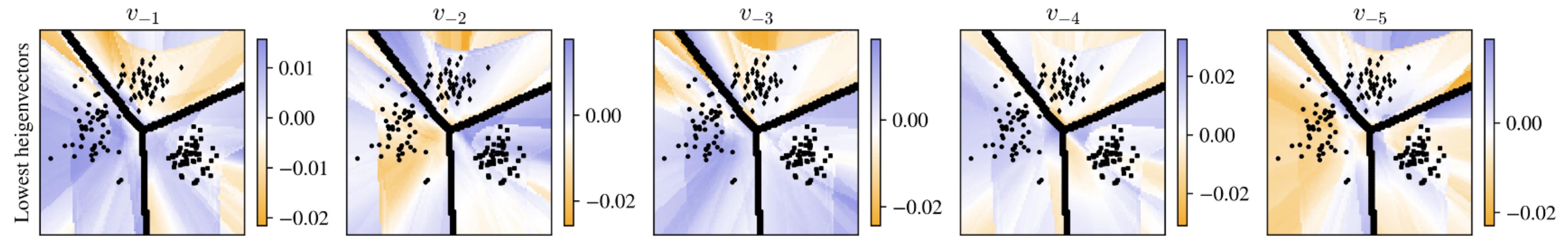
! The topmost eigenvectors have a close-to-one absolute alignment with gradients of loss of samples on the decision boundary of the network



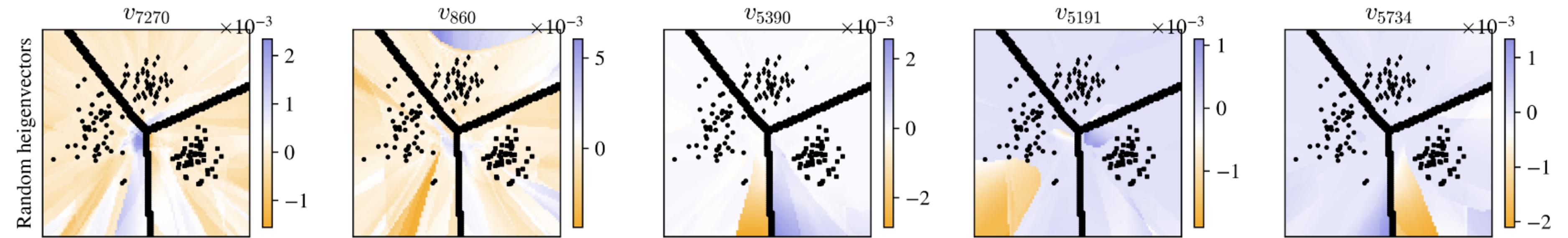
2D datasets

! Not accidental! Compare with the alignment with...

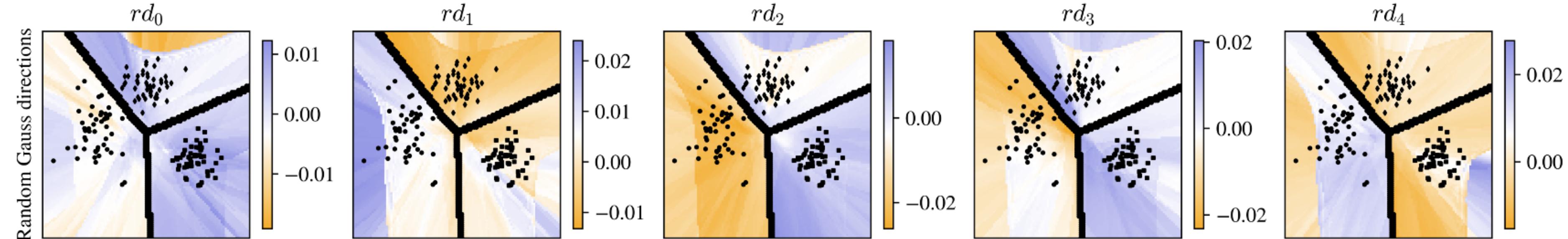
... Lowest eigenvectors



... Random eigenvectors



... Random directions



Outline



Empirical mess
in the Land of the Loss



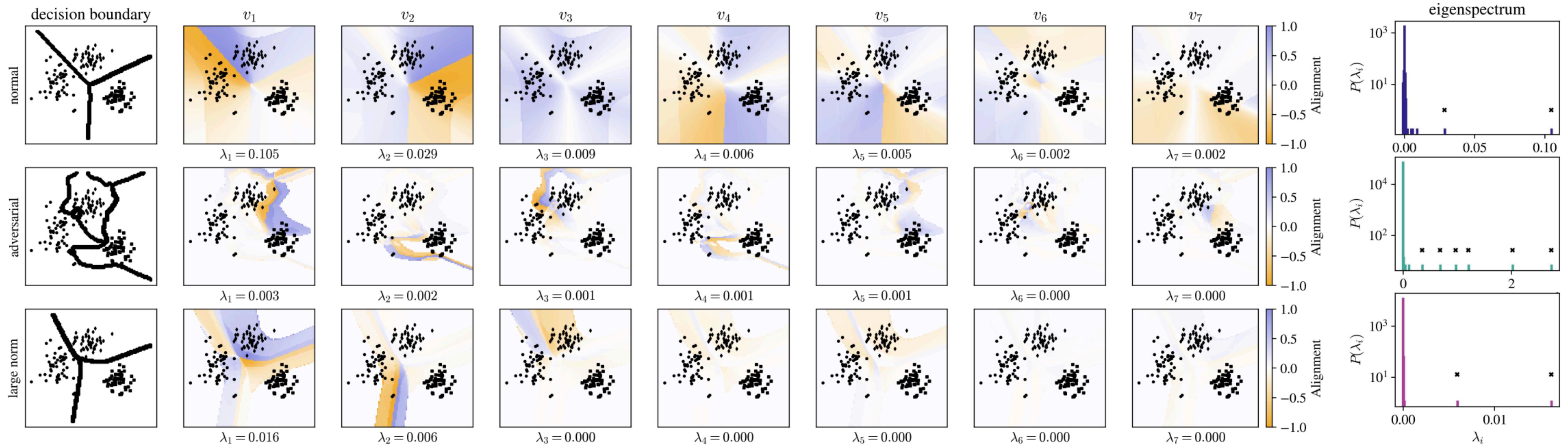
Hessian and the decision
boundary: toy example



**Selected interesting
consequences**

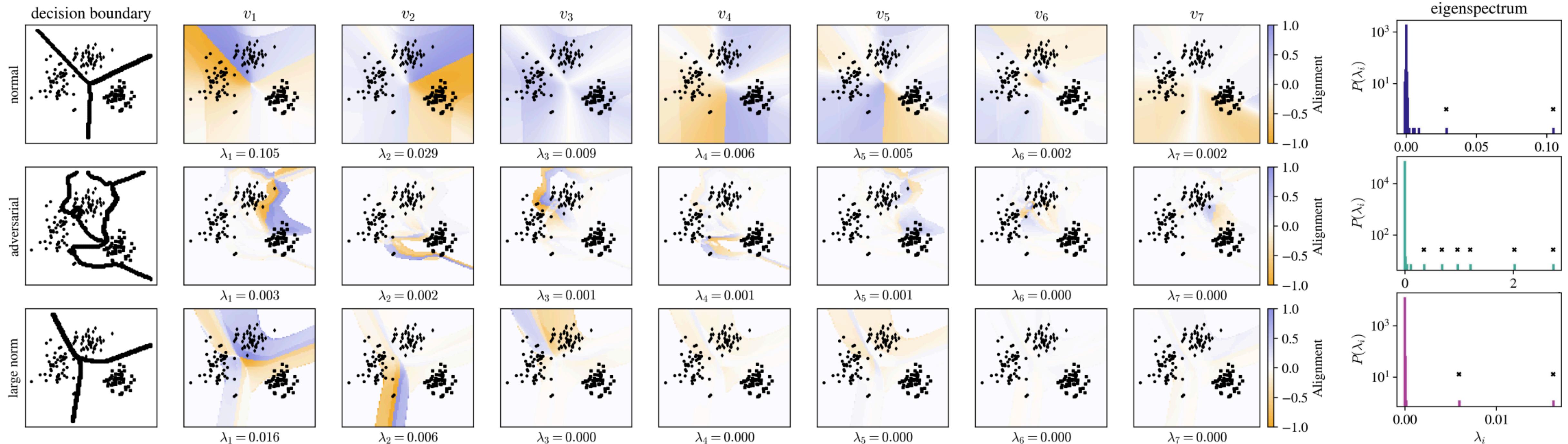
- measure of the decision boundary complexity
- perspectives for new generalization measure
- margin estimation

Decision boundary complexity



Decision boundary complexity

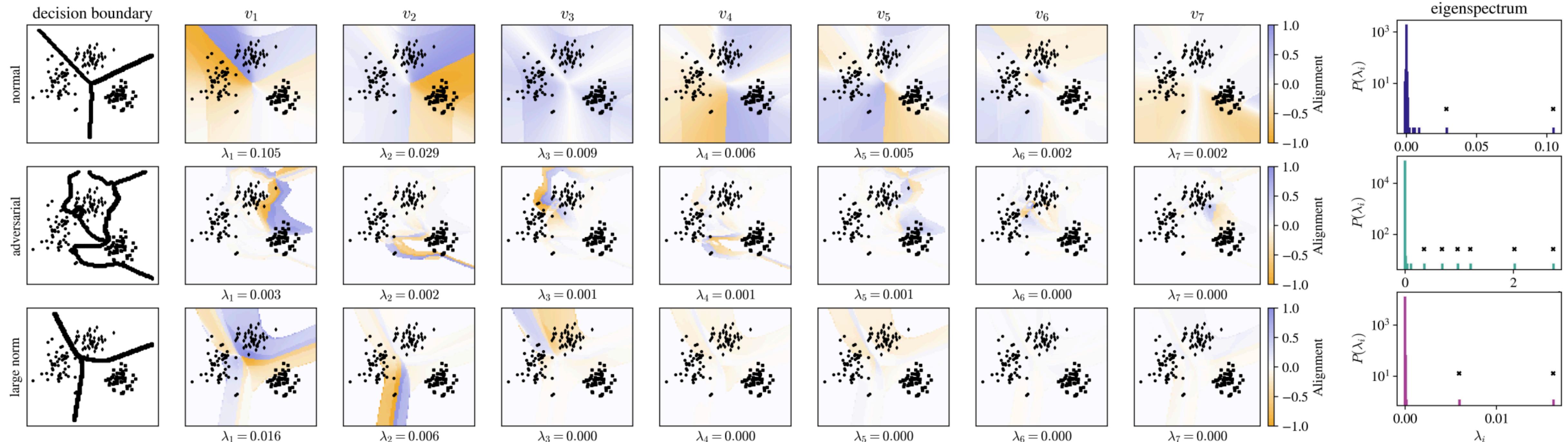
! More eigenvectors are needed to describe the decision boundary in the adversarial initialization and large norm case than in the normal training.



Decision boundary complexity

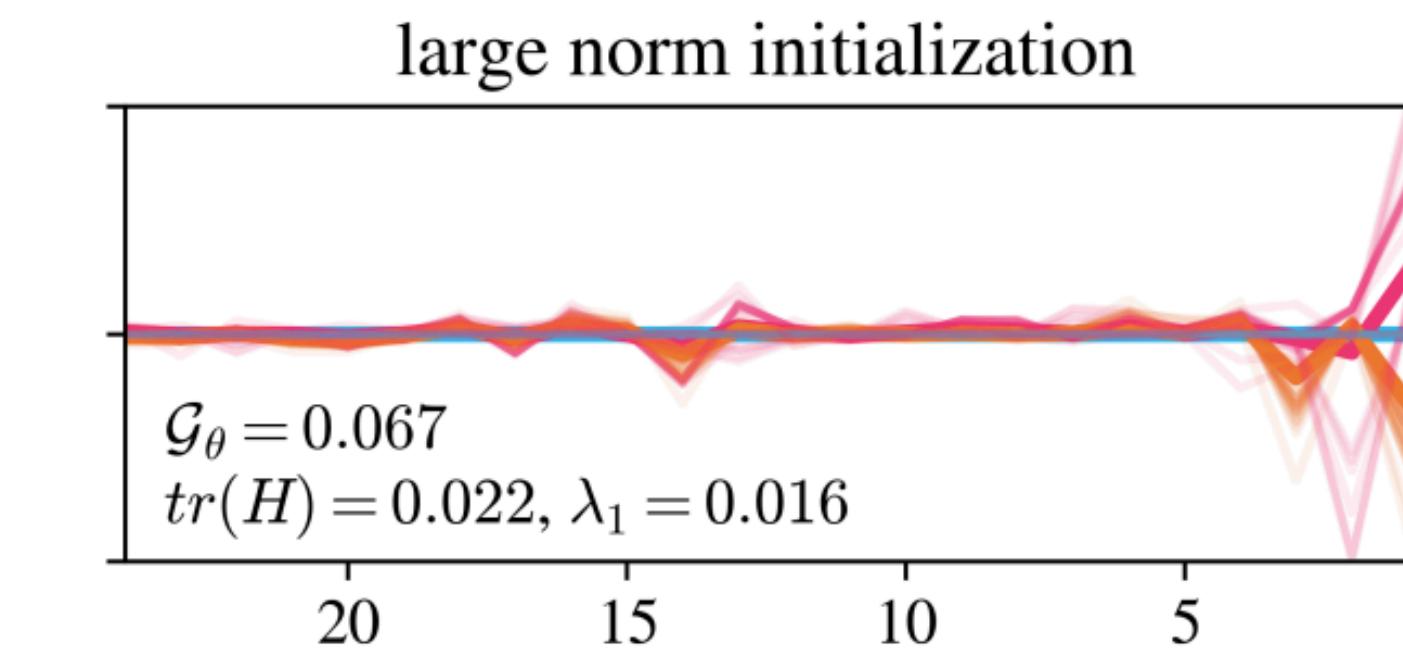
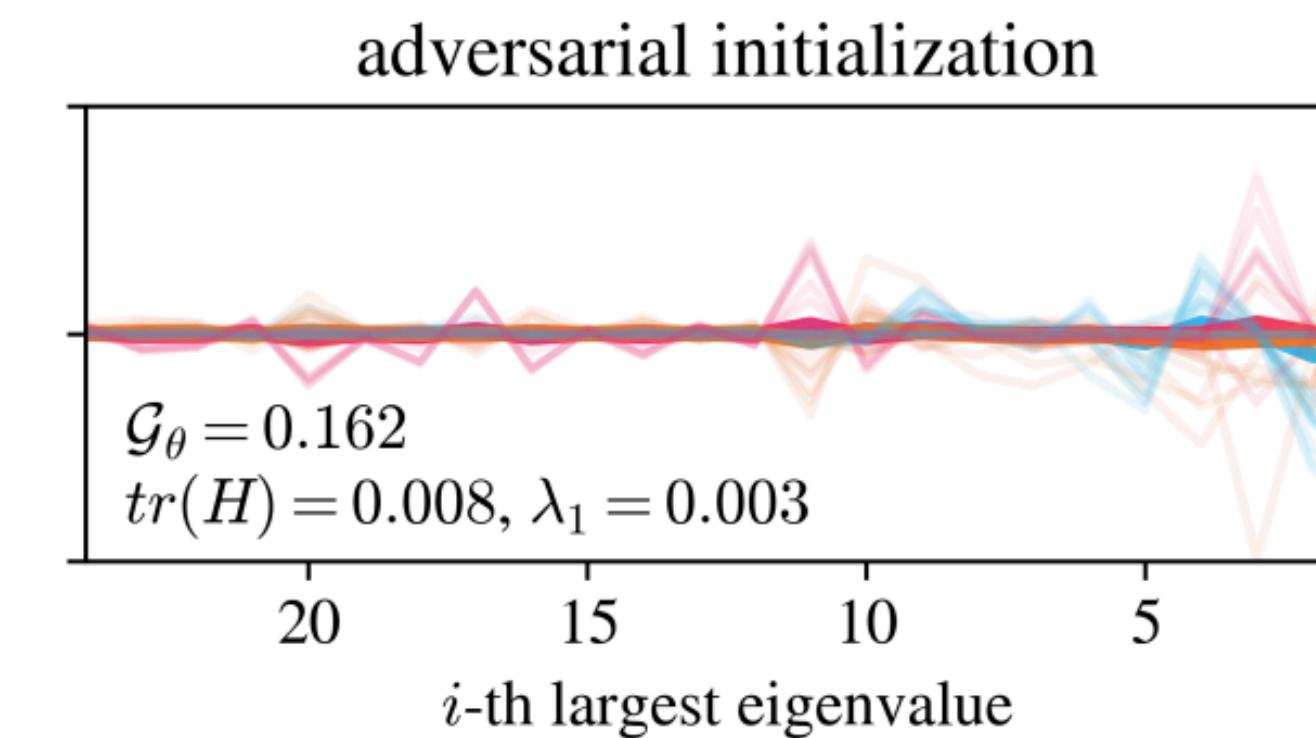
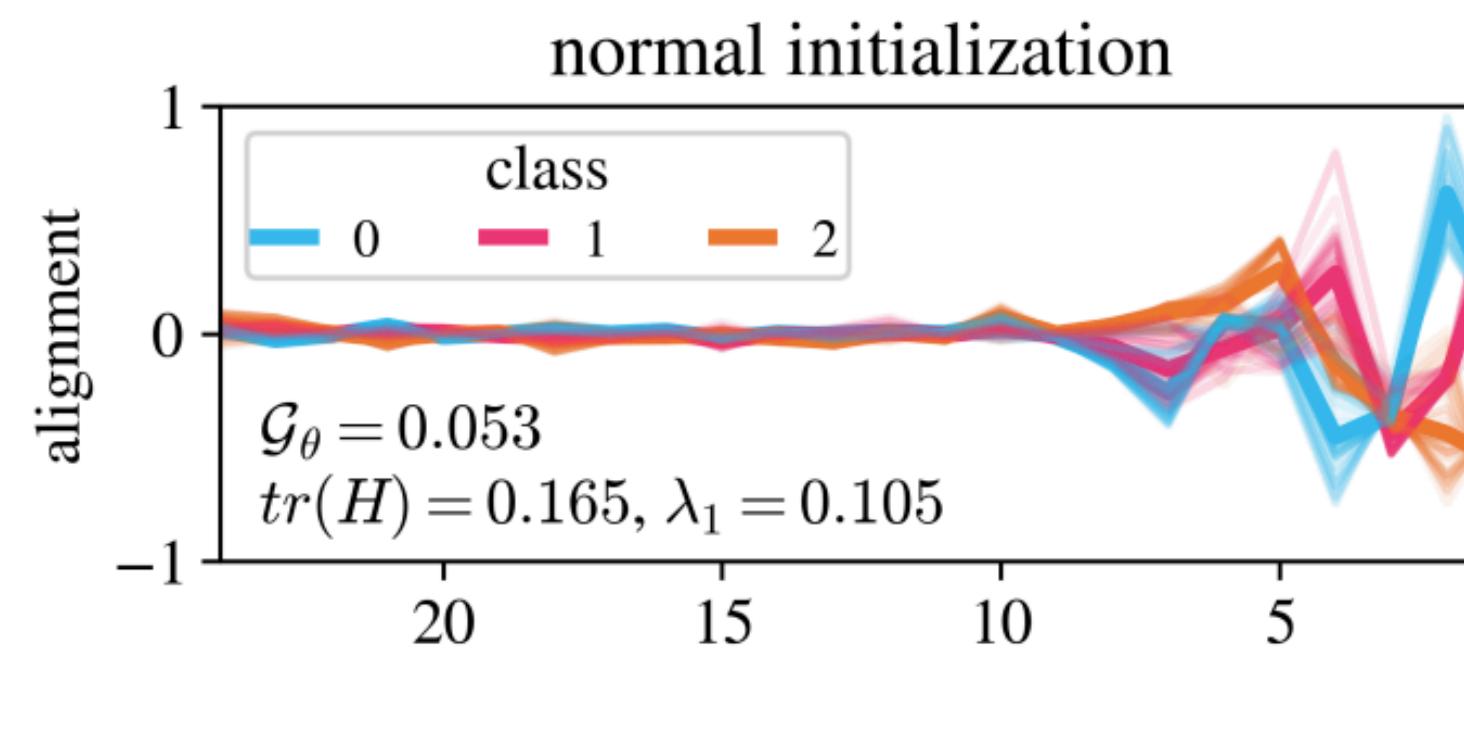
! More eigenvectors are needed to describe the decision boundary in the adversarial initialization and large norm case than in the normal training.

! Number of outliers in the spectrum depends on the decision boundary complexity (not only the number of classes)



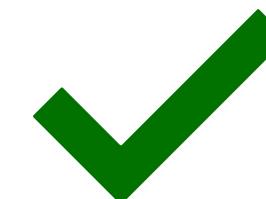
Decision boundary complexity measure!

Alignment of gradients of loss of individual training samples with top Hessian eigenvectors behaves differently for the simple and complex decision boundary

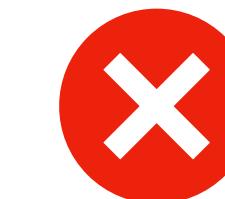


We quantified this different behavior in an attempt to create a generalization measure

$$m_i = \frac{1}{n} \sum_{s=1}^n |\mathcal{A}_i(x_s)| \quad ; \quad \mathcal{G}_\theta = \frac{1}{p} \sum_{i=1}^p \mathbb{1}[m_i > \epsilon]$$



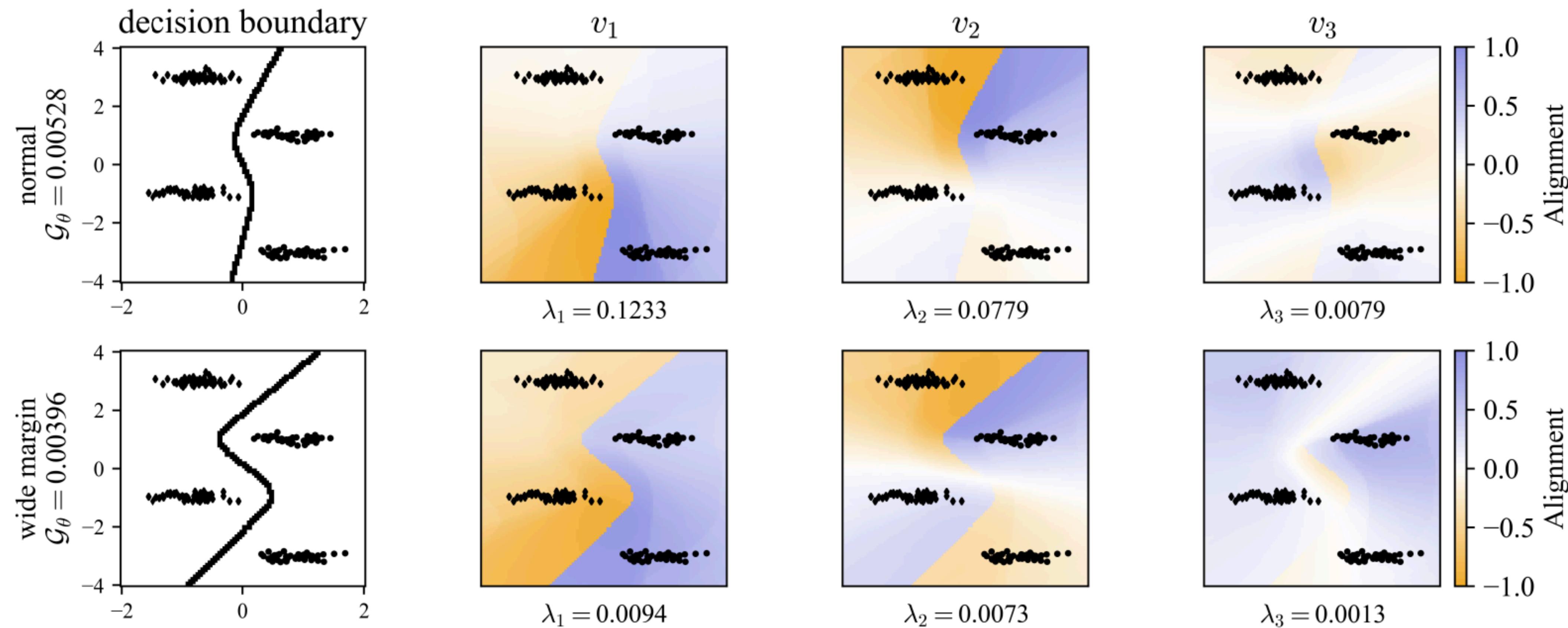
Invariant to reparametrization and promising results for 2D datasets, Iris, and few-digit MNIST



Expensive to compute and...

Simplicity bias though...

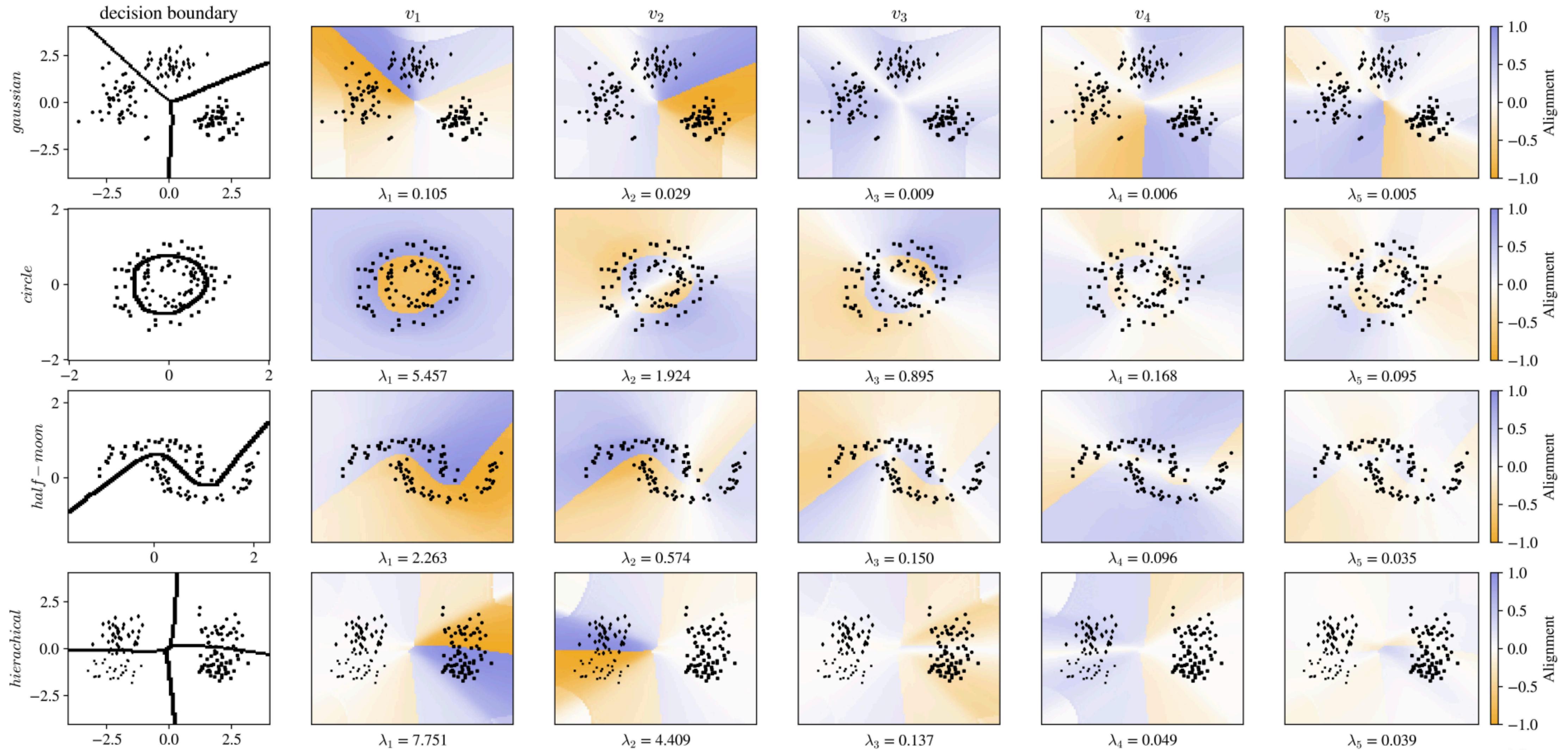
Simplicity of the decision boundary is not everything!
There is also a margin!



Complexity measure cannot detect the poorer generalization caused by the simplicity bias!



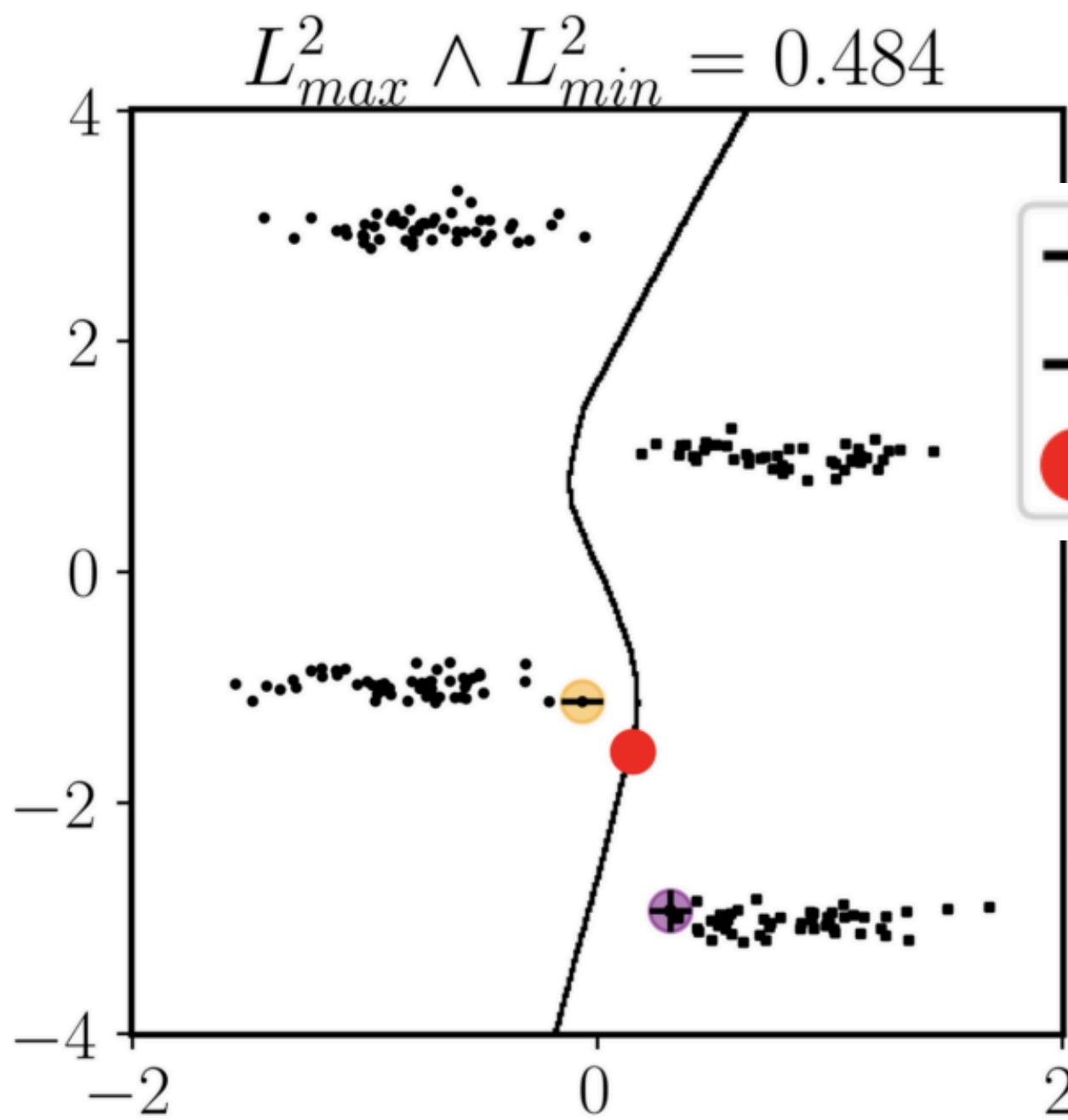
However, the order of the top eigenvectors follows the increasing margin of the sections of the boundary that they encode



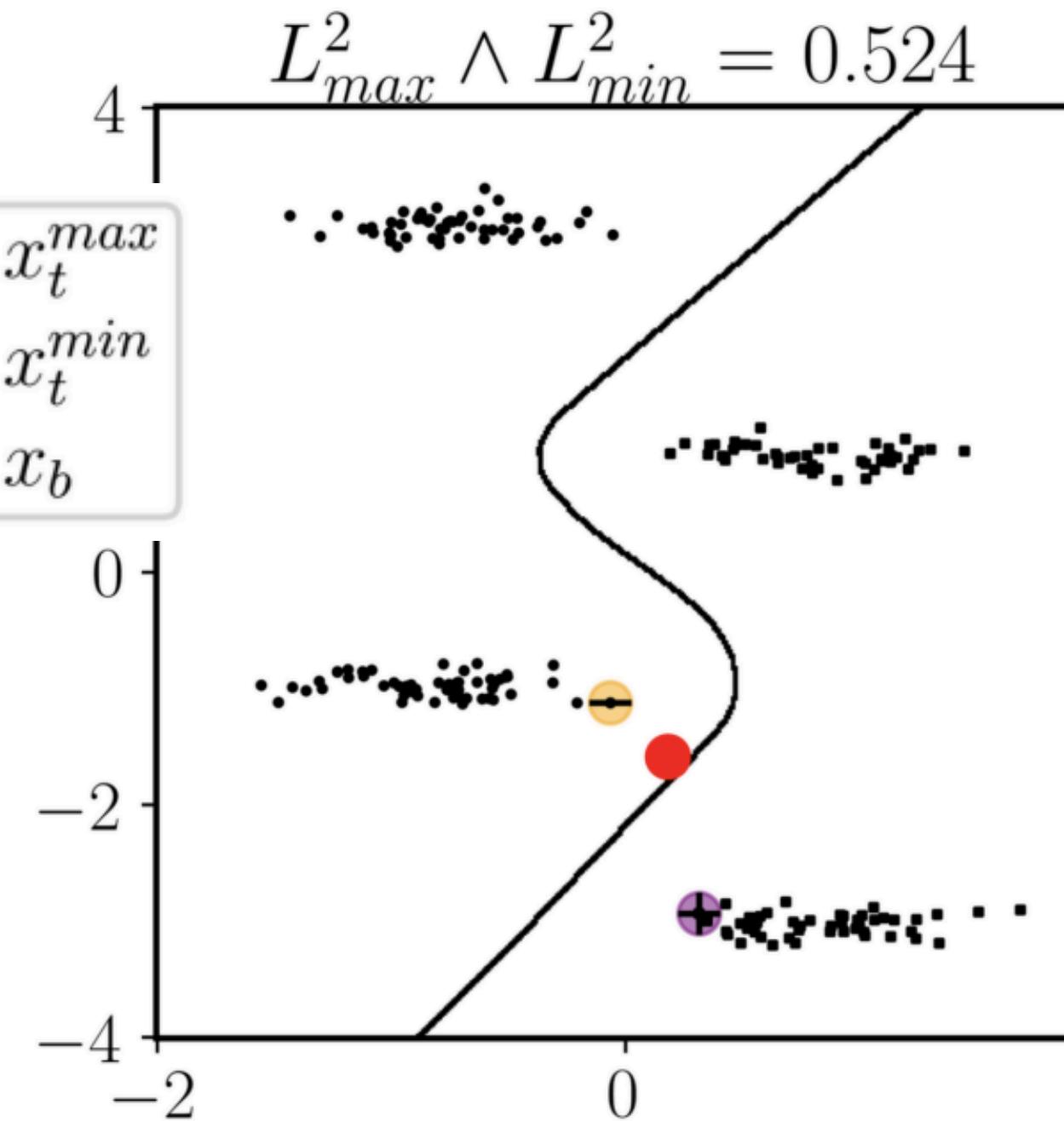
We can estimate the margin of various sections of the decision boundary!

Margin estimation

The top Hessian eigenvector allows to estimate the margin of the decision boundary!



(a) normal initialization



(b) wide-margin initialization



Training samples x_t that is the closest to the boundary in its narrowest part: they are **selected** as the one with the largest alignment with the top Hessian eigenvector v_1

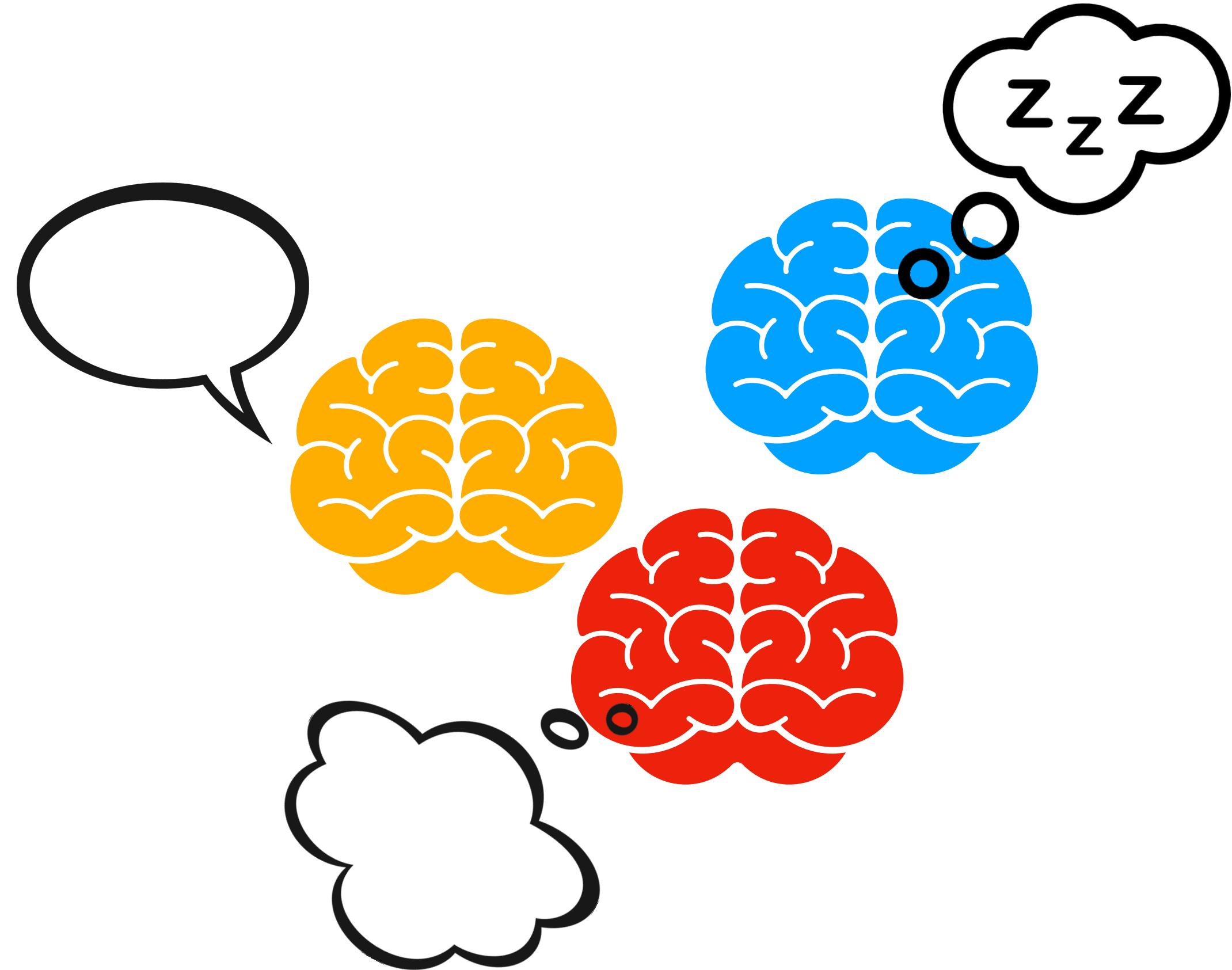


A sample on the boundary x_b in its narrowest part: it is **optimized** to have the largest alignment with the top Hessian eigenvector v_1

Together, the decision boundary complexity measure and the margin estimation technique indicate a better generalizing model!

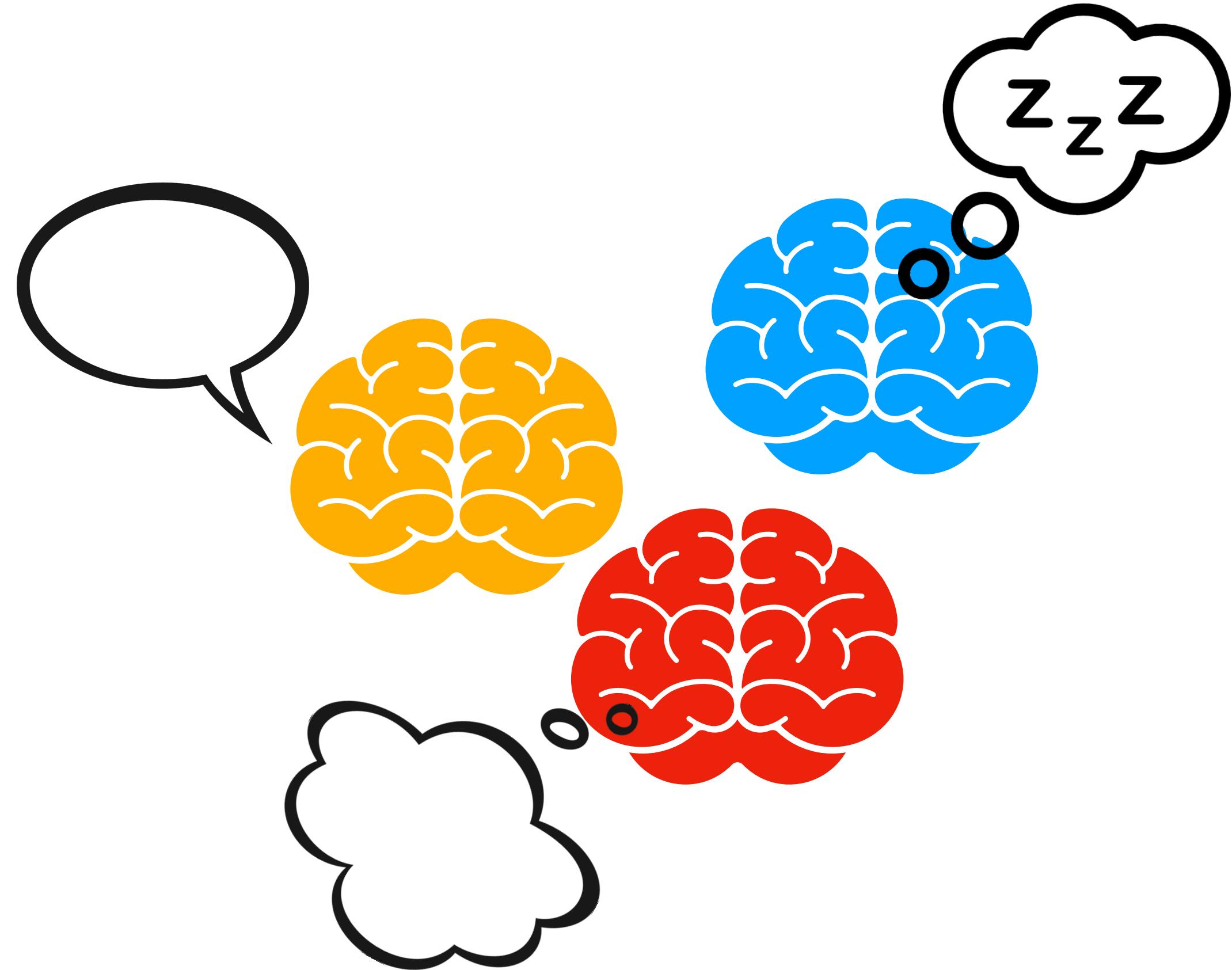
Open questions

- 1. Making the generalization measure useful**
- 2. Theoretical connection**
- 3. Higher dimensions**
 - Does separation between eigenvectors/boundaries hold?
 - Do we still need more eigenvectors for crooked high-dim boundaries?
 - How to even attack it?



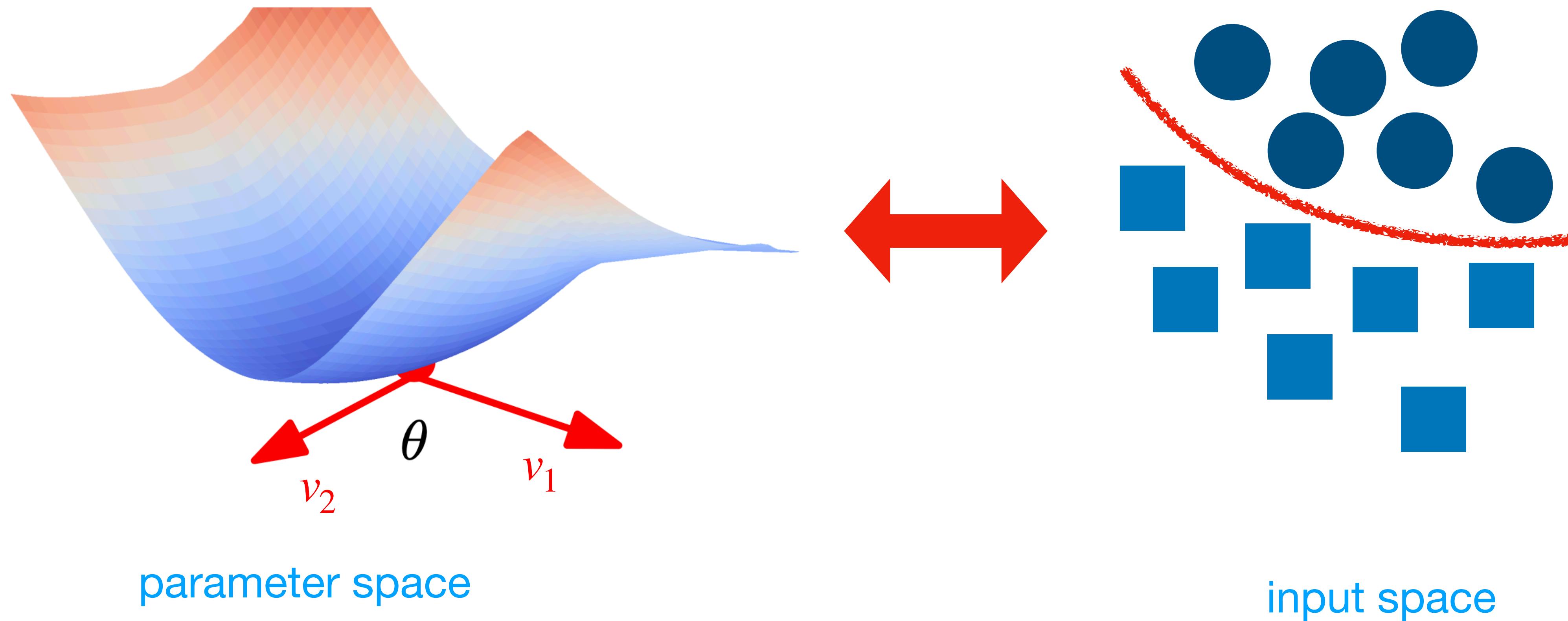
Cool new perspectives

1. Can it tell us anything about the model uncertainty?
2. Margin estimation and robustness
3. Do decision boundaries of the adversarially trained models look the same as the regularly trained models?



Take-home message

The top Hessian eigenvectors encode the decision boundary



Amazing collaborators



Mahalakshmi Sabanayagam
TU Munich



Freya Behrens
EPFL

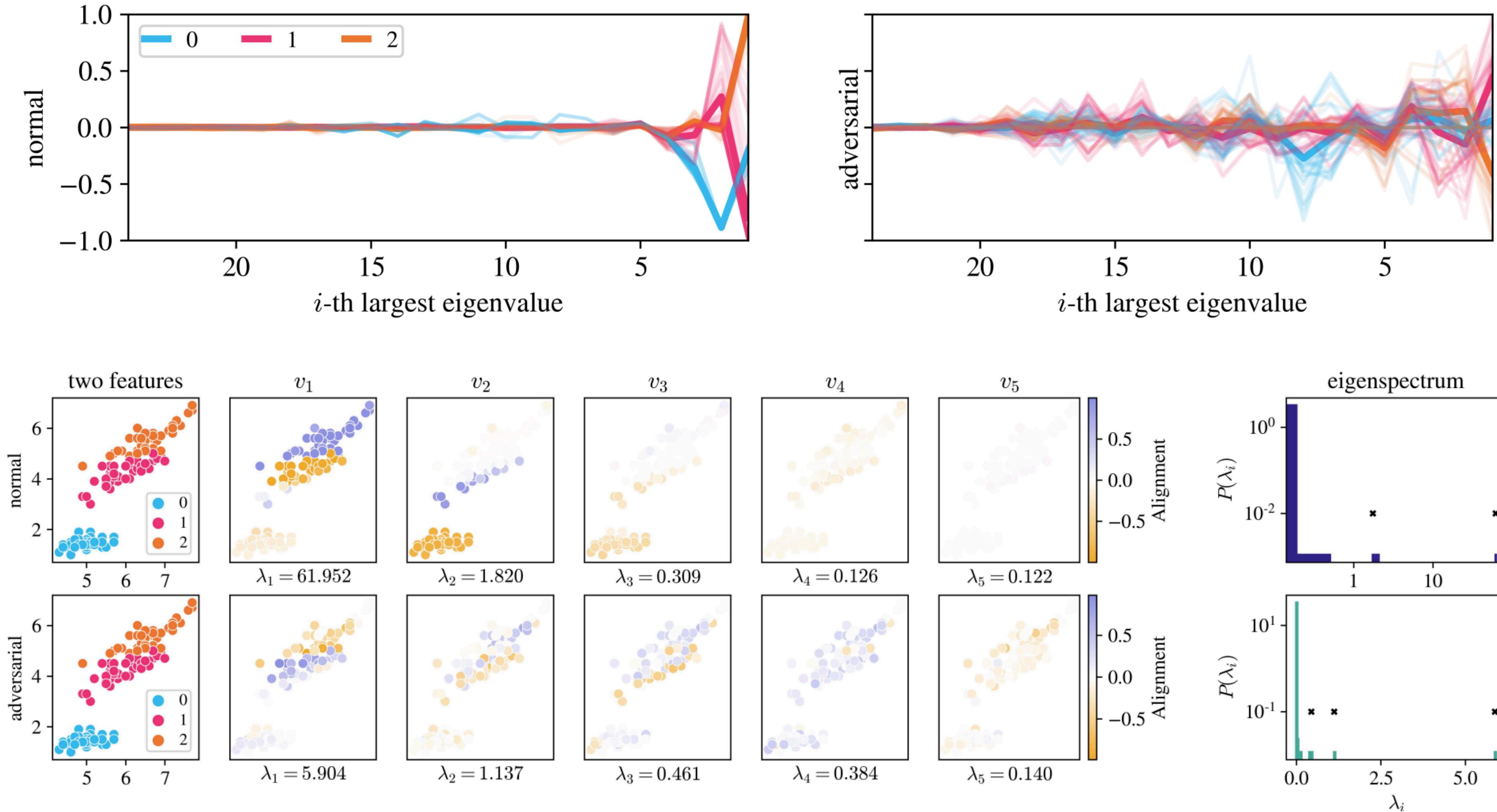


Urte Adomaityte
King's College London

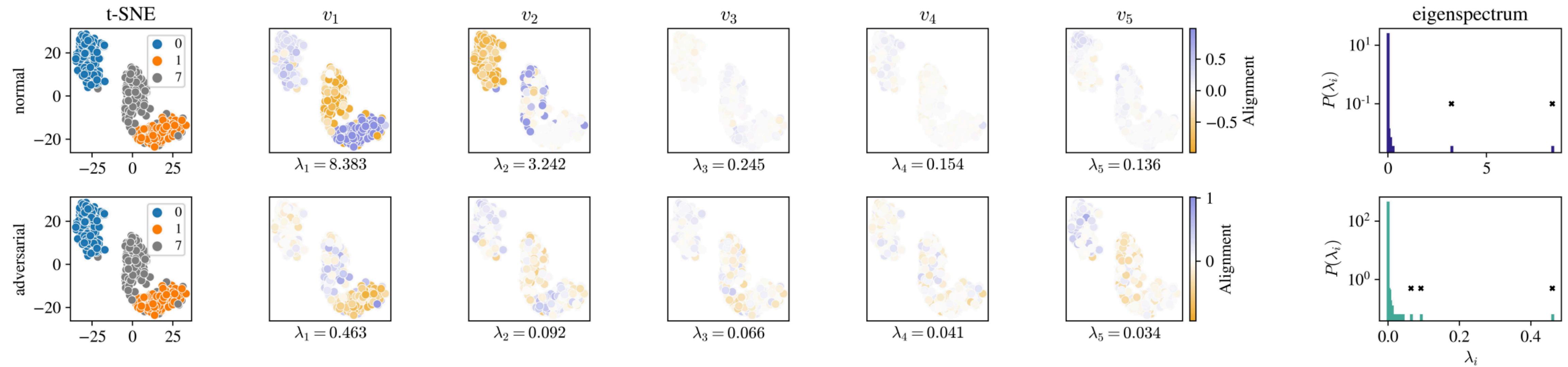
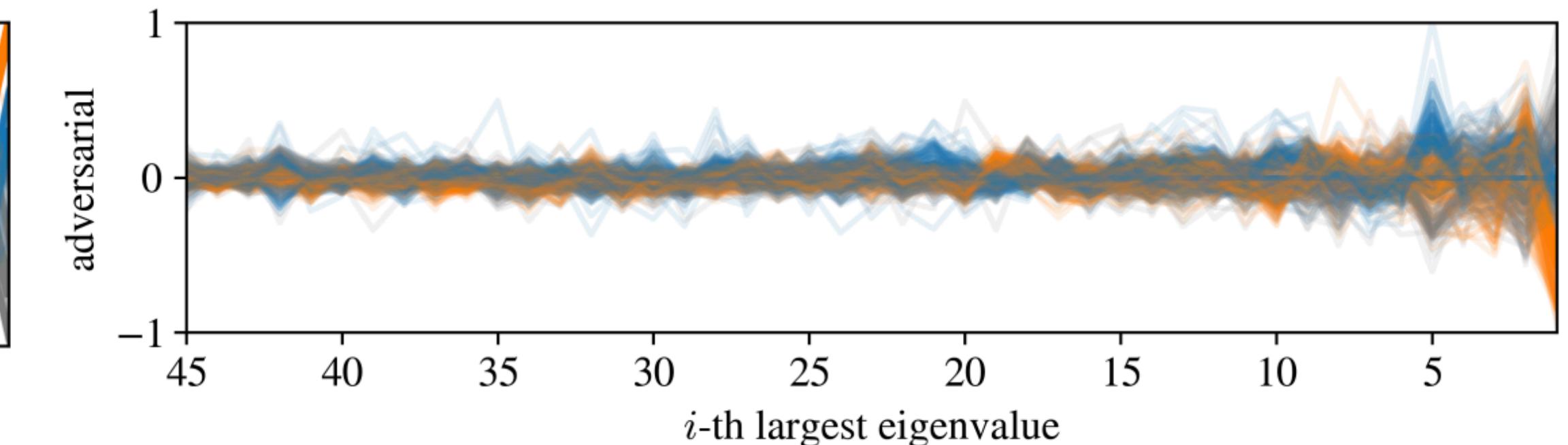
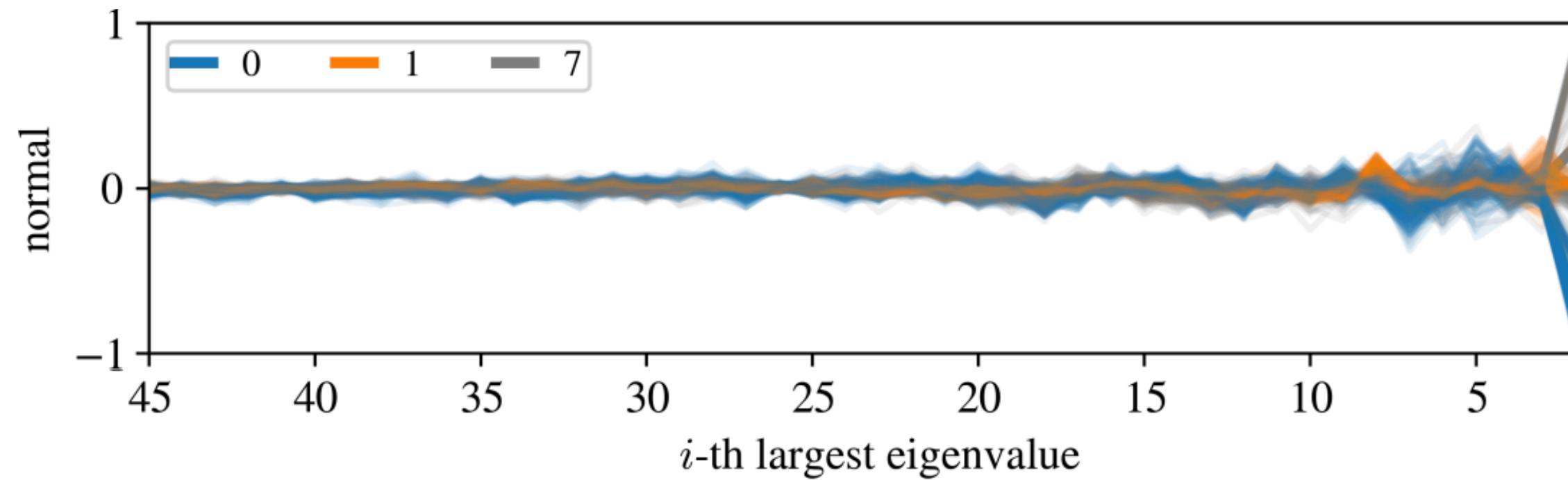
Thank you for your attention!



High(er)-dimensional results: Iris



High-dimensional results: MNIST-017



Mathematical analysis

We consider that the converged minimum $\theta := \theta^*$ is an exact minimum, meaning that the loss and its gradient at θ^* is zero, i.e., $\mathcal{L}(\theta^*; \mathcal{D}) = 0$ and $\nabla_{\theta}\mathcal{L}(\theta^*; \mathcal{D}) = \mathbf{0}$. Using this information, we expand the loss using the second-order Taylor's approximation at $\theta := \theta^*$.

$$\begin{aligned}\mathcal{L}(\theta^* + \Delta\theta; \mathcal{D}) &= \mathcal{L}(\theta^*; \mathcal{D}) + \nabla_{\theta}\mathcal{L}(\theta^*; \mathcal{D})^T \Delta\theta + \frac{1}{2} \Delta\theta^T \nabla_{\theta}^2 \mathcal{L}(\theta^*; \mathcal{D}) \Delta\theta \\ \mathcal{L}(\theta^* + \Delta\theta; \mathcal{D}) &= \frac{1}{2} \Delta\theta^T H_{\theta^*} \Delta\theta\end{aligned}$$

H_{θ^*} is the Hessian of the training loss function evaluated at the minimum. We denote its eigenvectors and corresponding eigenvalues as v_i and λ_i . Now, let's consider $\Delta\theta := \frac{g_{\theta}(x)}{\|g_{\theta}(x)\|}$ to be a reinforcing gradient of some input x in the dataset $\mathcal{D} := \{\mathcal{X}, \mathcal{Y}\}$, and an overparametrized classifier (e.g., neural network) with p parameters denoted by $f(\theta, \mathcal{X})$.

$$\mathcal{L} \left(f \left(\theta^* + \frac{g_\theta(x)}{\|g_\theta(x)\|}, \mathcal{X} \right), \mathcal{Y} \right) = \frac{1}{2} \frac{g_\theta(x)^T}{\|g_\theta(x)\|} \sum_{i=1}^p \lambda_i v_i v_i^T \frac{g_\theta(x)}{\|g_\theta(x)\|}$$

$$= \frac{1}{2} \sum_{i=1}^p \lambda_i \frac{\langle g_\theta(x), v_i \rangle}{\|g_\theta(x)\|} \frac{\langle g_\theta(x), v_i \rangle}{\|g_\theta(x)\|}$$

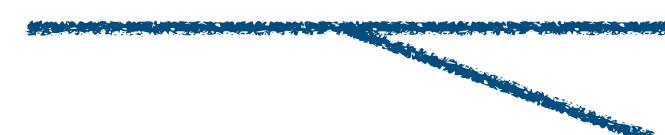
$$\mathcal{L} \left(f(\theta^*, \mathcal{X}) + \nabla_\theta f(\theta^*, \mathcal{X})^T \frac{g_\theta(x)}{\|g_\theta(x)\|}, \mathcal{Y} \right) = \frac{1}{2} \sum_{i=1}^p \lambda_i \mathcal{A}_i(x)^2$$

$$\mathcal{L} \left(\mathcal{Y} + \nabla_\theta f(\theta^*, \mathcal{X})^T \frac{g_\theta(x)}{\|g_\theta(x)\|}, \mathcal{Y} \right) = \frac{1}{2} \sum_{i=1}^p \lambda_i \mathcal{A}_i(x)^2$$



maximal for:

- $\nabla_\theta f(\theta^*, \mathcal{X})$ that changes predictions of maximal number of samples (shifts the decision boundary)
- $g_\theta(x)$ aligned with this $\nabla_\theta f(\theta^*, \mathcal{X})$



maximal when

$g_\theta(x)$ and v_1 are aligned