

# A rigorous, closed-form characterisation of adversarial generalisation errors.

## A High Dimensional Statistical Model for Adversarial Training: Geometry and Trade-Offs

### Binary Classification

**Teacher-Student Setting:** Training data  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \in \mathbb{R}^d \times \{-1, +1\}$  where  $x_i \sim \mathcal{N}(0, \Sigma_x)$  and  $y_i \sim \mathbb{P}(y|\theta_0^\top x_i)$  for fixed  $\theta_0 \in \mathbb{R}^d$ .

Focus on probit model  $\mathbb{P}(y|z) = 1/2 \operatorname{erfc}(-z/\sqrt{2}\tau)$ , with noise parameter  $\tau > 0$ .

**Problem Statement:** Learn  $\hat{y}(\hat{\theta}(\mathcal{D}), x) = \operatorname{sign}(\hat{\theta}(\mathcal{D}) \cdot x/\sqrt{d})$  using adversarial training.

**Generalisation Error:**

$$E_{\text{gen}} = \mathbb{E}_{y,x} [\mathbb{1}(y \neq \hat{y}(\hat{\theta}, x))], \quad (1)$$

**Adversarial Attack:** Chose  $\|v\|_{\Sigma_v^{-1}} \leq \varepsilon_g$  to misclassify the sample  $x_i$ :  $y(x_i) \neq \hat{y}(\hat{\theta}, x_i + v)$ .

**Adversarial Generalisation Error:**

$$E_{\text{adv}} = \mathbb{E}_{y,x} \left[ \max_{\|\delta\|_{\Sigma_v^{-1}} \leq \varepsilon_g} \mathbb{1}(y \neq \hat{y}(\hat{\theta}, x + \delta)) \right], \quad (2)$$

Note,  $E_{\text{gen}} = E_{\text{adv}}(\varepsilon_g = 0)$ .

**Boundary Error:**  $E_{\text{adv}} = E_{\text{gen}} + E_{\text{bnd}}$  where  $E_{\text{bnd}}$  are the attackable samples.

### Adversarial ERM

$$\sum_{i=1}^n \max_{\|\delta\|_{\Sigma_v^{-1}} \leq \varepsilon_t} g\left(y_i \frac{\theta^\top (x_i + \delta_i)}{\sqrt{d}}\right) + r(\theta), \quad (3)$$

where  $g$  is a convex loss,  $r(\theta)$  is a convex regularisation and  $\Sigma_\delta$  is positive definite. We set  $r(\theta) = \frac{\lambda}{2} \|\theta\|_2^2$ .

The inner maximisation has a closed form solution:

$$\sum_{i=1}^n g\left(y_i \frac{\theta^\top x_i}{\sqrt{d}} - \varepsilon_t \frac{\sqrt{\theta^\top \Sigma_\delta \theta}}{\sqrt{d}}\right) + r(\theta). \quad (4)$$

**The High-Dimensional Proportional Limit (hdl).** The dimension  $d$  and the number of training samples  $n$  are large  $d, n \rightarrow \infty$ , but the sample complexity is fixed  $\alpha := n/d$ .

### Usefulness and robustness

Usefulness relates to generalisation error and robustness relates to boundary error.

$$\mathcal{U}_{\theta_0} = \frac{1}{\sqrt{d}} \mathbb{E}_{x,y} [y \theta_0^\top x], \quad (5)$$

$$\mathcal{R}_{\theta_0} = \frac{1}{\sqrt{d}} \mathbb{E}_{x,y} \left[ \inf_{\|\delta\|_{\Sigma_v^{-1}} \leq \varepsilon_g} y \theta_0^\top (x + \delta) \right]. \quad (6)$$

### Block Features

Craft artificial usefulness and robustness:

$$\begin{aligned} \Sigma_x &= \text{blockdiag}(\psi_1 \mathbb{1}_{d_1}, \dots, \psi_k \mathbb{1}_{d_k}), \\ \Sigma_\delta &= \text{blockdiag}(\Delta_1 \mathbb{1}_{d_1}, \dots, \Delta_k \mathbb{1}_{d_k}), \\ \Sigma_v &= \text{blockdiag}(\Upsilon_1 \mathbb{1}_{d_1}, \dots, \Upsilon_k \mathbb{1}_{d_k}), \\ \Sigma_\theta &= \text{blockdiag}(t_1 \mathbb{1}_{d_1}, \dots, t_k \mathbb{1}_{d_k}), \end{aligned} \quad (7)$$

### Main Theorem Assumptions:

- Well defined spectral distribution in the hdl  $\Sigma_x = S^\top \text{diag}(\omega_i) S$ ,  $\zeta_i = \text{diag}(S \Sigma_\delta S^\top)_i$  and  $v_i = \text{diag}(S \Sigma_v S^\top)_i$ .
- Assume that  $\theta_0^\top \Sigma_x \theta_0 / d$  converges to  $\rho$  in the limit and that the entries of  $\bar{\theta} = S \Sigma_x^\top \theta_0 / \sqrt{\rho}$  converge as well to a limiting distribution.
- Assume that in the hdl the spectral distributions for the matrices and the distributions of the elements of the vectors just defined converge jointly to a p.d.f., i.e.  $1/d \sum_{i=1}^d \delta(\omega - \omega_i) \delta(\bar{\theta} - \bar{\theta}_i) \delta(\zeta - \zeta_i) \delta(v - v_i) \rightarrow \mu(\omega, \bar{\theta}, \zeta, v)$ .

For  $\lambda \geq 0$

$$E_{\text{gen}} = \frac{1}{\pi} \arccos\left(m/\sqrt{(\rho + \tau^2)q}\right), \quad (8)$$

$$E_{\text{bnd}} = \int_0^{\varepsilon_g \frac{\sqrt{A}}{\sqrt{q}}} \operatorname{erfc}\left(\frac{-\frac{m}{\sqrt{q}} \nu}{\sqrt{2(\rho + \tau^2 - m^2/q)}}\right) \frac{e^{-\frac{\nu^2}{2}}}{\sqrt{2\pi}} d\nu, \quad (9)$$

and  $E_{\text{adv}} = E_{\text{gen}} + E_{\text{bnd}}$ .

$(m, q, V, P, \hat{m}, \hat{q}, \hat{V}, \hat{P})$  are the solutions of the following system of equations. Four are dependent on the loss

function  $g$  and the adversarial training strength  $\varepsilon_t$ :

$$\begin{cases} \hat{m} = \alpha \mathbb{E}_\xi \left[ \int_{\mathbb{R}} dy \partial_\omega Z_0 f_g(y, \sqrt{q} \xi, P) \right] \\ \hat{q} = \alpha \mathbb{E}_\xi \left[ \int_{\mathbb{R}} dy Z_0 f_g^2(y, \sqrt{q} \xi, P) \right] \\ \hat{V} = -\alpha \mathbb{E}_\xi \left[ \int_{\mathbb{R}} dy Z_0 \partial_\omega f_g(y, \sqrt{q} \xi, P) \right] \\ \hat{P} = -\frac{\varepsilon_t}{2\sqrt{P}} \alpha \mathbb{E}_\xi \left[ \int_{\mathbb{R}} dy y Z_0 f_g(y, \sqrt{q} \xi, P) \right] \end{cases}, \quad (10)$$

where  $\xi \sim \mathcal{N}(0, 1)$  and  $Z_0 = 1/2 \operatorname{erfc}(-y\omega/\sqrt{2(V+\tau^2)})$  and  $f_g(y, \omega, V, P) = (P(\omega) - \omega)/V$ , where  $P$  is the following proximal operator

$$P(\omega) = \min_x \left[ \frac{(x - \omega)^2}{2V} + g(yx - \varepsilon_t \sqrt{P}) \right]. \quad (11)$$

The other four are dependent on the spectral distribution of the matrices  $\Sigma_x, \Sigma_\delta$  and on the limiting distribution of the elements of  $\bar{\theta}$ :

$$\begin{cases} m = \mathbb{E}_\mu \left[ \frac{\hat{m} \bar{\theta}^2}{\lambda + \hat{V} \omega + \hat{P} \delta} \right] \\ q = \mathbb{E}_\mu \left[ \frac{\hat{m}^2 \bar{\theta}^2 \omega + \hat{q} \omega^2}{(\lambda + \hat{V} \omega + \hat{P} \delta)^2} \right] \\ V = \mathbb{E}_\mu \left[ \frac{\omega}{\lambda + \hat{V} \omega + \hat{P} \delta} \right] \\ P = \mathbb{E}_\mu \left[ \zeta \frac{\hat{m}^2 \bar{\theta}^2 + \hat{q} \omega^2}{(\lambda + \hat{V} \omega + \hat{P} \delta)^2} \right] \end{cases}. \quad (12)$$

$A$  is given by

$$A = \mathbb{E}_\mu \left[ v \frac{\hat{m}^2 \bar{\theta}^2 \omega + \hat{q} \omega^2}{(\lambda + \hat{V} \omega + \hat{P} \delta)^2} \right]. \quad (13)$$

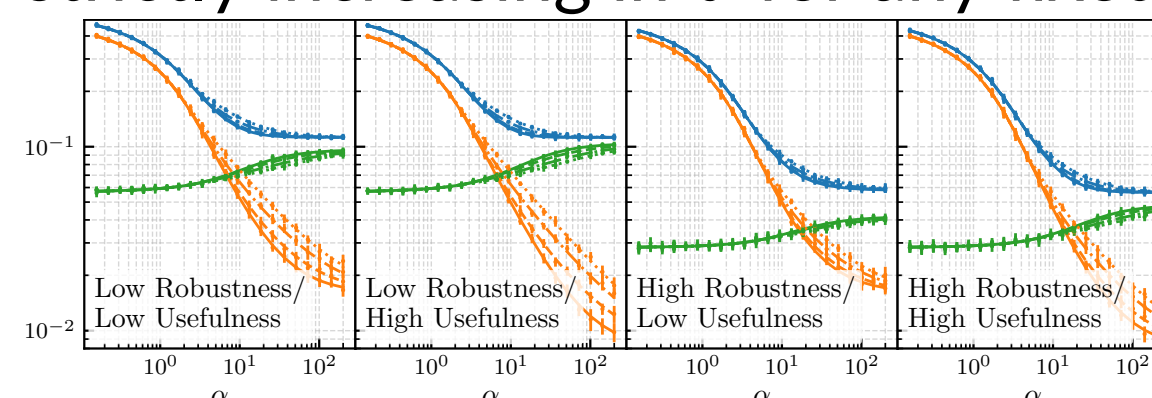
$m, q, P$  and  $A$  are interpretable:

$$\begin{aligned} m &= \mathbb{E}_D \left[ \frac{1}{q} \theta_0^\top \Sigma_x \hat{\theta} \right], & q &= \mathbb{E}_D \left[ \frac{1}{q} \hat{\theta}^\top \Sigma_x \hat{\theta} \right], \\ P &= \mathbb{E}_D \left[ \frac{1}{q} \hat{\theta}^\top \Sigma_\delta \hat{\theta} \right], & A &= \mathbb{E}_D \left[ \frac{1}{q} \hat{\theta}^\top \Sigma_v \hat{\theta} \right]. \end{aligned} \quad (14)$$

### Trade-Offs

$$E_{\text{adv}} = E_{\text{gen}}(\vartheta, \mathcal{U}_{\theta_0}) + \int_0^{\varepsilon_g \kappa} f(\xi; \vartheta, \mathcal{U}_{\theta_0}) d\xi, \quad (15)$$

where we introduce the variable  $\vartheta = m/\sqrt{\rho q}$  and  $\kappa = \sqrt{A}/\sqrt{q}$ .  $\vartheta$  is the cosine of the angle between the teacher weights  $\theta_0$  and the student estimate  $\hat{\theta}$  in the geometry of  $\Sigma_x$  and  $\kappa$  is the norm of  $\hat{\theta}$  under the attack matrix. The function  $f(\xi; \vartheta)$  is positive  $\forall \vartheta, \forall \xi \in [0, +\infty)$  and it is strictly increasing in  $\vartheta$  for any fixed  $\xi \in [0, +\infty)$ .



We notice that the values for  $E_{\text{gen}}$  and  $E_{\text{bnd}}$  change by varying the usefulness and robustness of the features for fixed types of attacks. Intuitively, we have that the more usefulness one has the less generalisation error one makes, indeed we can write a lower bound for the generalisation error

$$E_{\text{gen}} \geq \frac{1}{\pi} \arccos\left(\sqrt{\frac{\pi}{2\rho}} \mathcal{U}_{\theta_0}\right). \quad (16)$$

We note that  $\rho$  and  $\mathcal{U}_{\theta_0}$  only depend on  $\Sigma_x$  and  $\theta_0$ .

Robustness only affects the boundary error. High robustness implies less sensibility to adversarial attacks: robust features have less samples within an attack range of the student decision boundary. The highest value that the boundary error can achieve is limited by both the robustness and the usefulness as

$$\begin{aligned} E_{\text{bnd}} &\leq 2T(\varepsilon_g A \mathcal{B}, A^{-1}) - \frac{1}{\pi} \arctan(A^{-1}) \\ &\quad - \frac{1}{\pi} \operatorname{erf}\left(\frac{\varepsilon_g \mathcal{B}}{\sqrt{2}}\right) \operatorname{erfc}\left(\frac{\varepsilon_g A \mathcal{B}}{\sqrt{2}}\right), \end{aligned} \quad (17)$$

where  $\mathcal{B} = \max_i \sqrt{(\Sigma_v)_{ii}/(\Sigma_x)_{ii}}$ ,  $A = \sqrt{\pi} \mathcal{U}_{\theta_0}/\sqrt{2\rho}$  and  $T$  is the Owen function. This previous bound is a decreasing function of the robustness.

Under the same setting as ?? and considering a BFM with a single type of feature, i.e.  $k = 1$  one has that  $\forall \varepsilon_g, \varepsilon_t \geq 0$  for  $\alpha$  big enough exist two positive numbers  $M_1, M_2$  such that

$$\begin{aligned} |E_{\text{adv}}(\varepsilon_g, \varepsilon_t) - E_{\text{adv}}(\varepsilon_g, \varepsilon_t = 0)| &< M_1/\alpha, \\ |E_{\text{gen}}(\varepsilon_t) - E_{\text{gen}}(\varepsilon_t = 0)| &< M_2/\alpha, \end{aligned} \quad (18)$$

where  $E_{\text{adv}}(\varepsilon_g, \varepsilon_t)$  and  $E_{\text{gen}}(\varepsilon_t)$  define the adversarial and generalisation error of  $\hat{\theta}$  trained with  $\varepsilon_t$  and evaluated for  $\varepsilon_g$ .

### Directional Defences and structured data

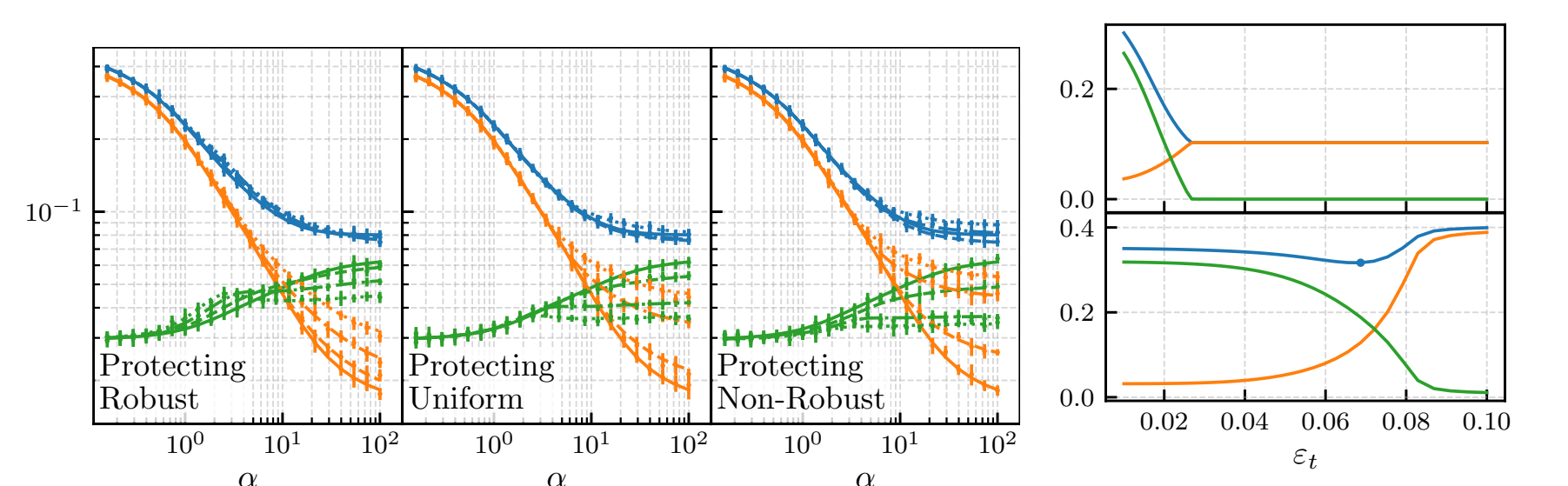
Consider the SWFM defined in ?? where the defence matrix is  $\Sigma_\delta = \text{blockdiag}((\Delta_1 + \delta_1 \varrho) \mathbb{1}_{d_1}, (\Delta_2 + \delta_2 \varrho) \mathbb{1}_{d_2})$ , with  $\varrho$  the parameter that makes the defence matrix change. Assume also that  $\psi_1 > \psi_2$ ,  $\Delta_2 \psi_1 \geq \Delta_1 \psi_2$  and  $\Upsilon_i = 1$ . In the  $\alpha \rightarrow \infty$  (taken after the  $n, d \rightarrow \infty$ ) there exists  $\kappa > 0$  such that  $\forall \delta_1 > \kappa, \delta_2 = -\delta_1$  one has that

$$\begin{aligned} E_{\text{bnd}}(\varrho) &= E_{\text{bnd}}^0 + E_{\text{bnd}}^1 \varrho + \mathcal{O}(\varrho^2), \\ E_{\text{gen}}(\varrho) &= E_{\text{gen}}^0 + E_{\text{gen}}^1 \varrho + \mathcal{O}(\varrho^2), \end{aligned} \quad (19)$$

where  $E_{\text{gen}}^1 > 0$ ,  $E_{\text{bnd}}^1 < 0$  and  $E_{\text{bnd}}^0, E_{\text{gen}}^0$  are the errors when  $\varrho = 0$ . Additionally, this leads to an improved value of  $E_{\text{adv}}$  at order  $\varrho$  iff the following condition is satisfied

$$\frac{\varepsilon_g}{\sqrt{2}} \operatorname{erfc}\left(-\frac{\vartheta_0 u_0 \varepsilon_g}{\sqrt{2-2\vartheta_0^2}}\right) < \frac{e^{-\frac{\vartheta_0^2 u_0^2 \varepsilon_g^2}{2(1-\vartheta_0^2)}}}{\sqrt{\pi} \sqrt{1-\vartheta_0^2}}, \quad (20)$$

where  $\vartheta_0 = m_0/\sqrt{\rho q_0}$  and  $u_0 = \sqrt{A_0}/\sqrt{q_0}$  the solution of the problem with  $\varrho = 0$ . Notice that for  $\varepsilon_g$  small enough this condition is always verified.



TODO add caption

### Tradeoff directions and innocuous directions

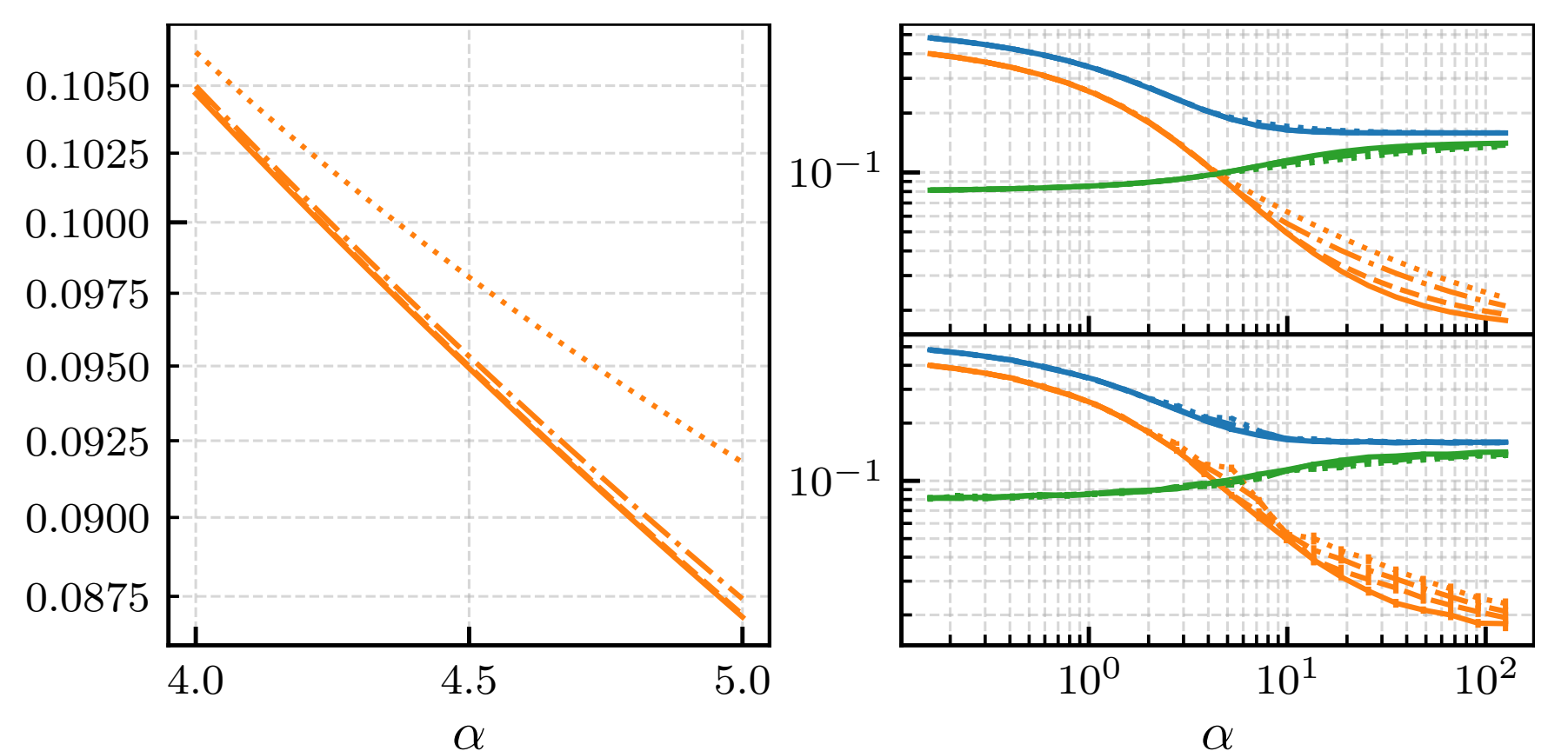
We now investigate the effect that different types of geometries have on the trade-off between  $E_{\text{gen}}$  and  $E_{\text{bnd}}$ . Depending on the attack geometry  $\Sigma_v$  one can choose different defence geometries  $\Sigma_\delta$  and ask if and for which, protection without trade-off is possible. Any attack matrix  $\Sigma_v$  eigenvalues can be split into directions orthogonal to the teacher and directions aligned with the teacher. ?? (Left) considers the effect of the adversarial training strength  $\varepsilon_t$  on the errors for different choices of matrices  $\Sigma_\delta = \Sigma_v$ . In the the top we consider matrices whose biggest eigenvalues are orthogonal to the teacher vector and in the bottom one matrices where there is a leading eigenvector in the direction of the teacher.

### Data Dependent Regularisation

The form for the approximate loss in the case of small  $\varepsilon_t$  is

$$\sum_{i=1}^n g\left(y_i \frac{\theta^\top x_i}{\sqrt{d}}\right) + \tilde{\lambda}_1 \sqrt{\theta^\top \Sigma_\delta \theta} + \tilde{\lambda}_2 \theta^\top \Sigma_\delta \theta \quad (21)$$

where  $\tilde{\lambda}_1$  and  $\tilde{\lambda}_2$  depend on the model's parameters and perturbed margins of the points that shift sign under perturbation.



TODO better caption

### Acknowledgements

Bruno Loureiro acknowledges support from the *Choose France - CNRS AI Rising Talents* program, and Florent Krzakala from the Swiss National Science Foundation grant SNFS OperaGOST (grant number 200390).





Kasimir Tanner  
Matteo Vilucchio  
Bruno Loureiro  
Florent Krzakala

EPFL