

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

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

Empirical Risk Minimization for Binary Classification

Teacher-Student Setting: Training data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \in \mathbb{R}^d \times \{-1, +1\}$ where $\mathbf{x}_i \sim \mathcal{N}(0, \Sigma_x)$ and $y_i \sim \mathbb{P}(y|\theta_0^\top \mathbf{x}_i)$ 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 $\hat{y}(\hat{\theta}(\mathcal{D}), \mathbf{x}) = \operatorname{sign}(\hat{\theta}(\mathcal{D}) \cdot \mathbf{x}/\sqrt{d})$ using adversarial training.

Generalisation Error:

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

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

Adversarial Generalisation Error:

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

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

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

$$\sum_{i=1}^n g\left(y_i \frac{\theta^\top \mathbf{x}_i}{\sqrt{d}} - \varepsilon_t \frac{\sqrt{\theta^\top \Sigma_\delta \theta}}{\sqrt{d}}\right) + r(\theta). \quad (3)$$

Block Features

$$\begin{aligned} \Sigma_x &= \text{blockdiag}(\psi_1 \mathbb{1}_{d_1}, \dots, \psi_k \mathbb{1}_{d_k}), \\ \Sigma_\delta &= \text{blockdiag}(\Delta_1 \mathbb{1}_{d_1}, \dots, \Delta_k \mathbb{1}_{d_k}), \\ \Sigma_v &= \text{blockdiag}(\Upsilon_1 \mathbb{1}_{d_1}, \dots, \Upsilon_k \mathbb{1}_{d_k}), \\ \Sigma_\theta &= \text{blockdiag}(t_1 \mathbb{1}_{d_1}, \dots, t_k \mathbb{1}_{d_k}), \end{aligned} \quad (4)$$

Usefulness and robustness

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

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

Main Theorem For the ERM estimator of the risk function with ℓ_2 regularisation $r(\theta) = \frac{\lambda}{2} \|\theta\|_2^2$ and $\lambda \geq 0$, under the data model defined in ?? and in the high dimensional proportional limit, the generalisation error E_{gen} and the boundary error E_{bnd} concentrate to

$$E_{\text{gen}} = \frac{1}{\pi} \arccos\left(m/\sqrt{(\rho + \tau^2)q}\right), \quad (7)$$

$$E_{\text{bnd}} = \int_0^{\varepsilon_g \frac{\sqrt{A}}{\sqrt{q}}} \operatorname{erfc}\left(\frac{-\frac{m}{\sqrt{q}}\nu}{\sqrt{2(\rho + \tau^2 - m^2/q)}}\right) \frac{e^{-\frac{\nu^2}{2}}}{\sqrt{2\pi}} d\nu, \quad (8)$$

and the adversarial generalisation error concentrates to $E_{\text{adv}} = E_{\text{gen}} + E_{\text{bnd}}$.

The values of m and q are the solutions of a system of eight self-consistent equations for the unknowns $(m, q, V, P, \hat{m}, \hat{q}, \hat{V}, \hat{P})$. The first four equations are dependant on the loss function g and the adversarial training strength ε_t and read

$$\begin{cases} \hat{m} = \alpha \mathbb{E}_\xi \left[\int_{\mathbb{R}} dy \partial_\omega \mathcal{Z}_0 f_g(y, \sqrt{q}\xi, P) \right] \\ \hat{q} = \alpha \mathbb{E}_\xi \left[\int_{\mathbb{R}} dy \mathcal{Z}_0 f_g^2(y, \sqrt{q}\xi, P) \right] \\ \hat{V} = -\alpha \mathbb{E}_\xi \left[\int_{\mathbb{R}} dy \mathcal{Z}_0 \partial_\omega f_g(y, \sqrt{q}\xi, P) \right] \\ \hat{P} = -\frac{\varepsilon_t}{2\sqrt{P}} \alpha \mathbb{E}_\xi \left[\int_{\mathbb{R}} dy y \mathcal{Z}_0 f_g(y, \sqrt{q}\xi, P) \right] \end{cases}, \quad (9)$$

where $\xi \sim \mathcal{N}(0, 1)$ and $\mathcal{Z}_0 = 1/2 \operatorname{erfc}(-y\omega/\sqrt{2(V+\tau^2)})$ and $f_g(y, \omega, V, P) = (\mathcal{P}(\omega) - \omega)/V$, where \mathcal{P} is the following proximal operator

$$\mathcal{P}(\omega) = \min_x \left[\frac{(x - \omega)^2}{2V} + g(yx - \varepsilon_t \sqrt{P}) \right]. \quad (10)$$

The second set of equation depend on the spectral distribution of the matrices Σ_x, Σ_δ and on

the limiting distribution of the elements of $\bar{\theta}$. The equations read

$$\begin{cases} m = \mathbb{E}_\mu \left[\frac{\hat{m} \bar{\theta}^2}{\lambda + \hat{V} \omega + \hat{P} \delta} \right] \\ q = \mathbb{E}_\mu \left[\frac{\hat{m}^2 \bar{\theta}^2 \omega + \hat{q} \omega^2}{(\lambda + \hat{V} \omega + \hat{P} \delta)^2} \right] \\ V = \mathbb{E}_\mu \left[\frac{\omega}{\lambda + \hat{V} \omega + \hat{P} \delta} \right] \\ P = \mathbb{E}_\mu \left[\frac{\zeta \hat{m}^2 \bar{\theta}^2 + \hat{q} \omega^2}{(\lambda + \hat{V} \omega + \hat{P} \delta)^2} \right] \end{cases}. \quad (11)$$

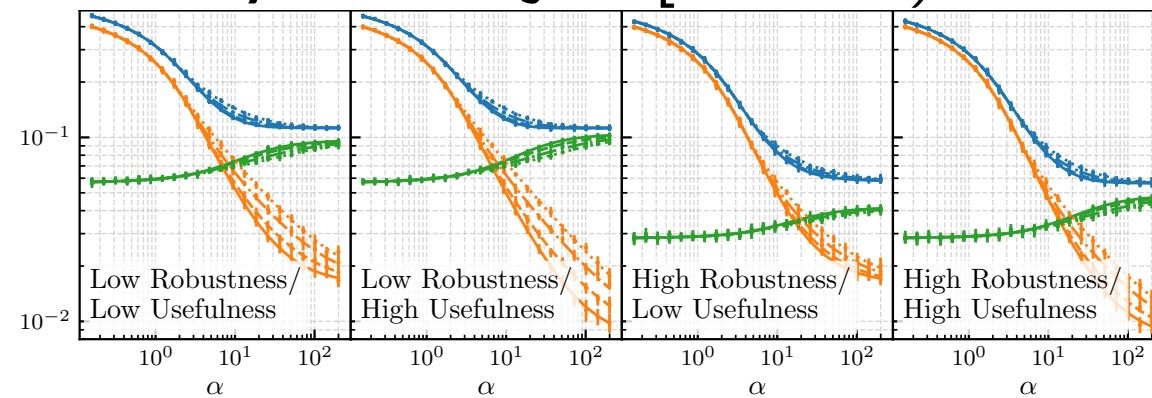
The value of A can be obtained from the solution of the same system of self consistent equations as

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

Trade-Offs

$$E_{\text{adv}} = E_{\text{gen}}(\vartheta, \mathcal{U}_{\theta_0}) + \int_0^{\varepsilon_g \kappa} f(\xi; \vartheta, \mathcal{U}_{\theta_0}) d\xi, \quad (13)$$

where we introduce the variable $\vartheta = m/\sqrt{\rho q}$ and $\kappa = \sqrt{A}/\sqrt{q}$. ϑ is the cosine of the angle between the teacher weights θ_0 and the student estimate $\hat{\theta}$ in the geometry of Σ_x and κ 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 ϑ for any fixed $\xi \in [0, +\infty)$.



We notice that the values for E_{gen} and E_{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

$$E_{\text{gen}} \geq \frac{1}{\pi} \arccos\left(\sqrt{\frac{\pi}{2\rho}} \mathcal{U}_{\theta_0}\right). \quad (14)$$

We note that ρ and \mathcal{U}_{θ_0} only depend on Σ_x and θ_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

$$\begin{aligned} E_{\text{bnd}} &\leq 2\mathsf{T}(\varepsilon_g \mathcal{A} \mathcal{B}, \mathcal{A}^{-1}) - \frac{1}{\pi} \arctan(\mathcal{A}^{-1}) \\ &\quad - \frac{1}{\pi} \operatorname{erf}\left(\frac{\varepsilon_g \mathcal{B}}{\sqrt{2}}\right) \operatorname{erfc}\left(\frac{\varepsilon_g \mathcal{A} \mathcal{B}}{\sqrt{2}}\right), \end{aligned} \quad (15)$$

where $\mathcal{B} = \max_i \sqrt{(\Sigma_v)_{ii}/(\Sigma_x)_{ii}}$, $\mathcal{A} = \sqrt{\pi} \mathcal{U}_{\theta_0}/\sqrt{2\rho}$ 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 α big enough exist two positive numbers M_1, M_2 such that

$$\begin{aligned} |E_{\text{adv}}(\varepsilon_g, \varepsilon_t) - E_{\text{adv}}(\varepsilon_g, \varepsilon_t = 0)| &< M_1/\alpha, \\ |E_{\text{gen}}(\varepsilon_t) - E_{\text{gen}}(\varepsilon_t = 0)| &< M_2/\alpha, \end{aligned} \quad (16)$$

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 ε_t and evaluated for ε_g .

Directional Defences and structured data

Consider the SWFM defined in ?? where the defence matrix is $\Sigma_\delta = \text{blockdiag}((\Delta_1 + \delta_1 \varrho) \mathbb{1}_{d_1}, (\Delta_2 + \delta_2 \varrho) \mathbb{1}_{d_2}, \dots, (\Delta_k + \delta_k \varrho) \mathbb{1}_{d_k})$ with ϱ the parameter that makes the defence

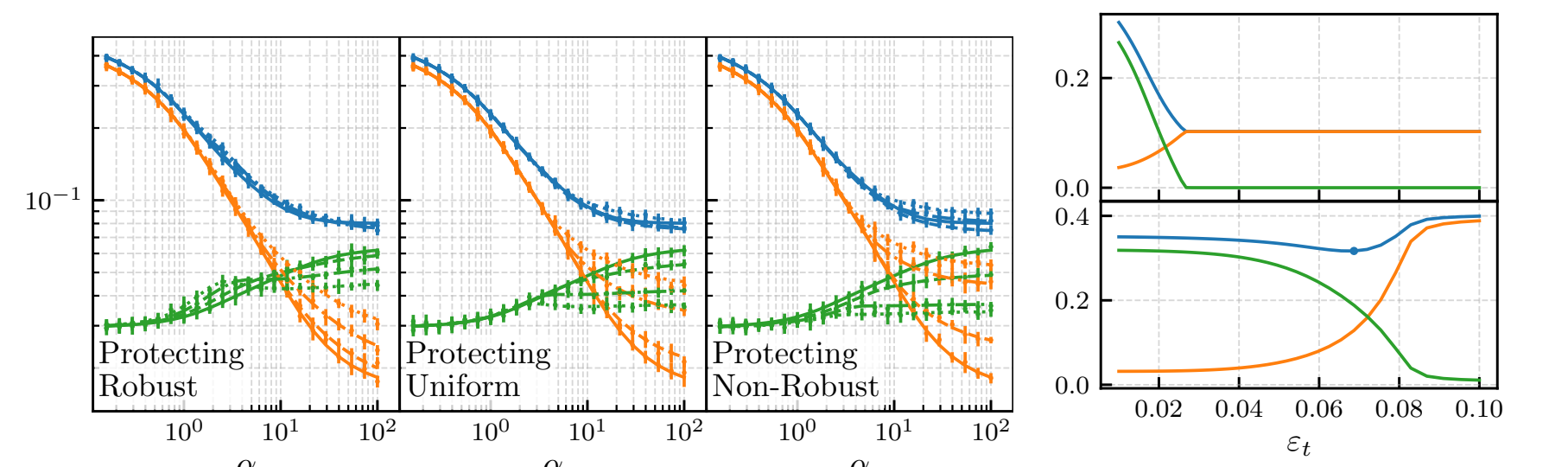
matrix change. Assume also that $\psi_1 > \psi_2, \Delta_2 \psi_1 \geq \Delta_1 \psi_2$ and $\Upsilon_i = 1$. In the $\alpha \rightarrow \infty$ (taken after the $n, d \rightarrow \infty$) there exists $\kappa > 0$ such that $\forall \delta_1 > \kappa, \delta_2 = -\delta_1$ one has that

$$\begin{aligned} E_{\text{bnd}}(\varrho) &= E_{\text{bnd}}^0 + E_{\text{bnd}}^1 \varrho + \mathcal{O}(\varrho^2), \\ E_{\text{gen}}(\varrho) &= E_{\text{gen}}^0 + E_{\text{gen}}^1 \varrho + \mathcal{O}(\varrho^2), \end{aligned} \quad (17)$$

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_{adv} at order ϱ iff the following condition is satisfied

$$\frac{\varepsilon_g}{\sqrt{2}} \operatorname{erfc}\left(-\frac{\vartheta_0 u_0 \varepsilon_g}{\sqrt{2-2\vartheta_0^2}}\right) < \frac{e^{-\frac{\vartheta_0^2 u_0^2 \varepsilon_g^2}{2(1-\vartheta_0^2)}}}{\sqrt{\pi} \sqrt{1-\vartheta_0^2}}, \quad (18)$$

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 ε_g small enough this condition is always verified.



TODO add caption

Tradeoff directions and innocuous directions

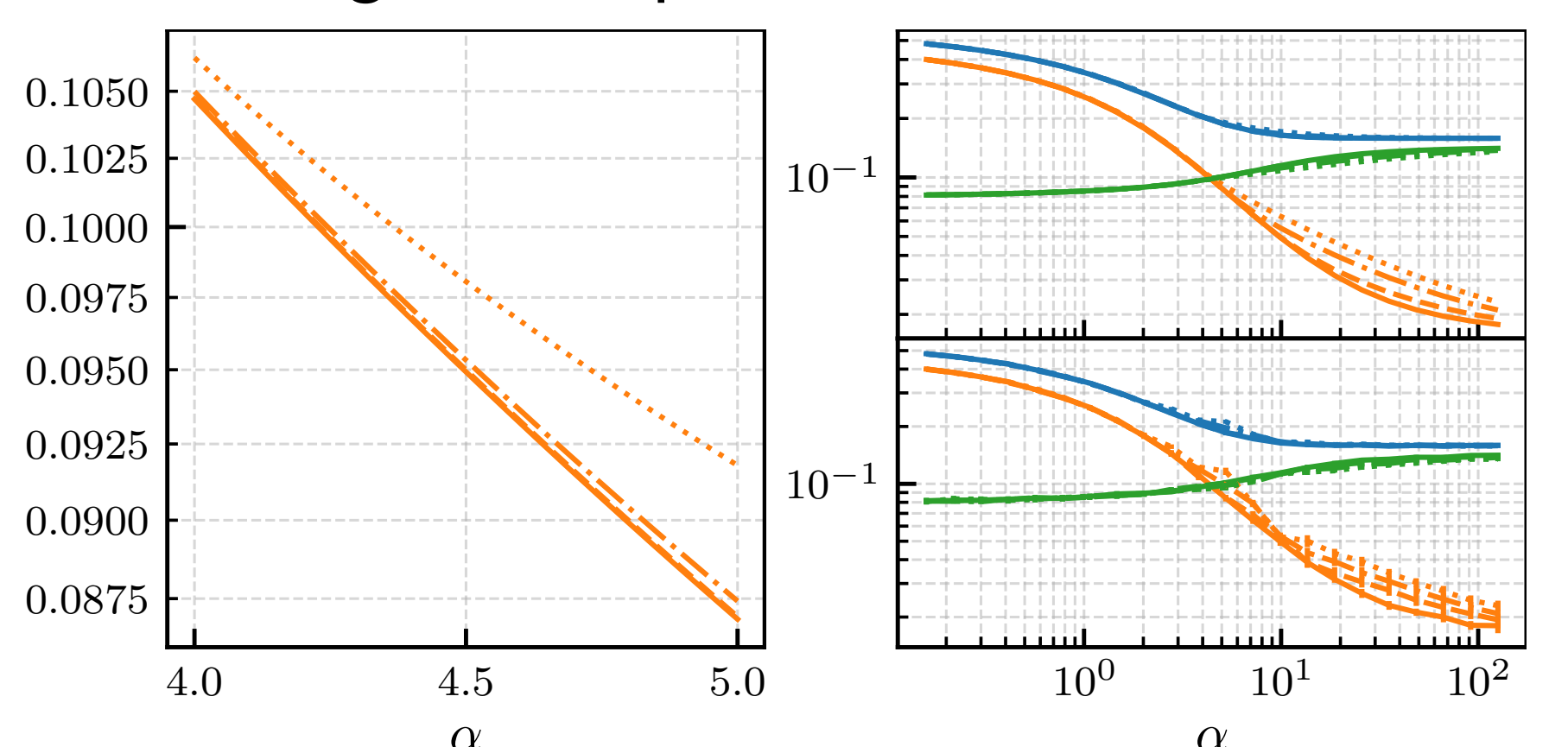
We now investigate the effect that different types of geometries have on the trade-off between E_{gen} and E_{bnd} . Depending on the attack geometry Σ_v one can choose different defence geometries Σ_δ and ask if and for which, protection without trade-off is possible. Any attack matrix Σ_v eigenvalues can be split into directions orthogonal to the teacher and directions aligned with the teacher. ?? (Left) considers the effect of the adversarial training strength ε_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 ε_t is

$$\sum_{i=1}^n g\left(y_i \frac{\theta^\top \mathbf{x}_i}{\sqrt{d}}\right) + \tilde{\lambda}_1 \sqrt{\theta^\top \Sigma_\delta \theta} + \tilde{\lambda}_2 \theta^\top \Sigma_\delta \theta \quad (19)$$

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