

Lec 21: Adversarial Robustness

Yaoliang Yu

July 21, 2020



UNIVERSITY OF
WATERLOO

Supervised Learning



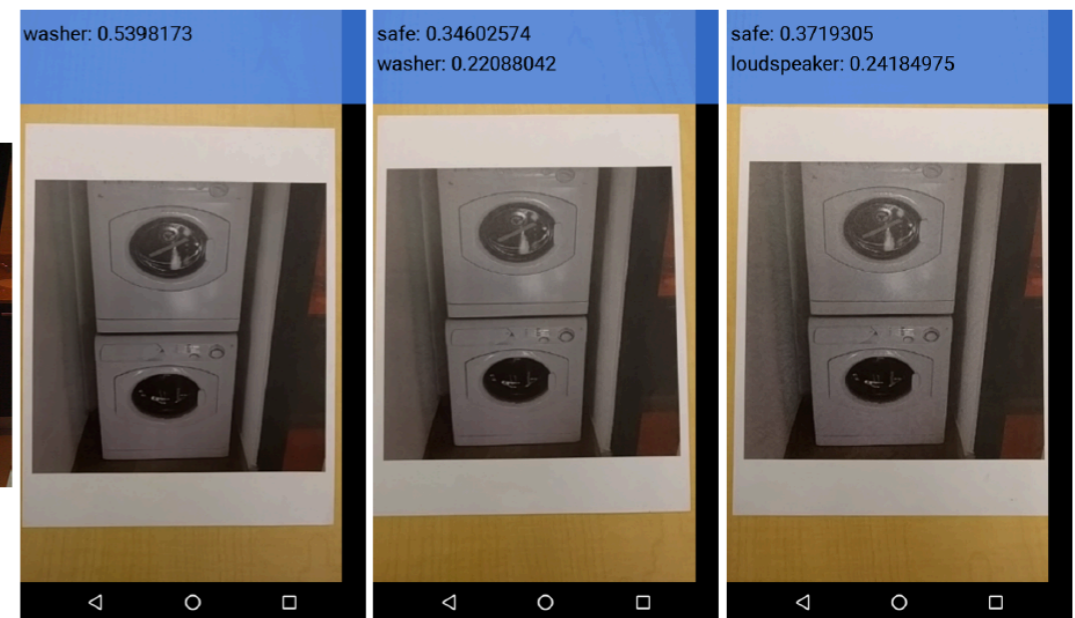
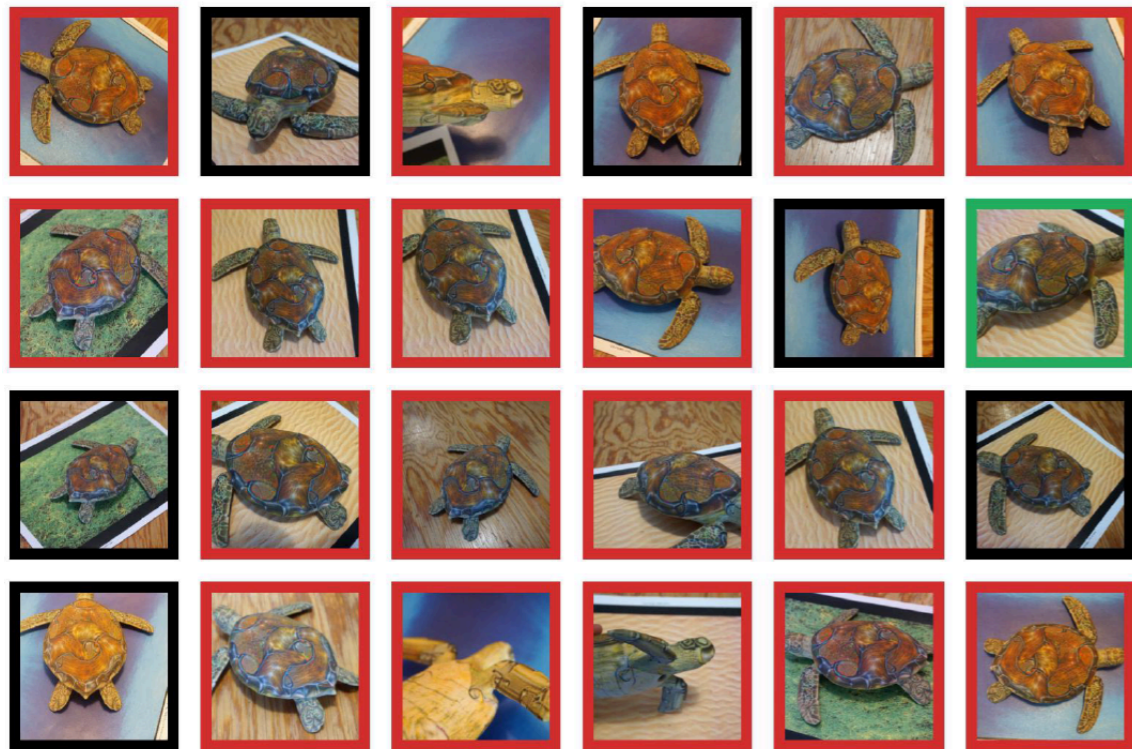
Formally

- Given a **training** set of **pairs** of examples $(\mathbf{x}_i, y_i) \in X \times Y$
- Return a function (classifier) $f : X \rightarrow Y$
- On an **unseen test** example x , output $f(x)$
- The goal is to do well on unseen test data
 - usually do not care about performance on training set

Performance Metric

- Accuracy (top-1, top-10 error, precision, recall, etc.)
- Training time
- Memory
- Test time
- Robustness
- Privacy
- Fairness
- Interpretability

And then the surprise



Szegedy et al. *Intriguing properties of neural networks*, ICLR 2014
Athalye et al. *Synthesizing Robust Adversarial Examples*, ICML 2018
Kurakin, Goodfellow, Bengio. *Adversarial Examples in the Physical World*, ICLR workshop 2017

Why should we care?

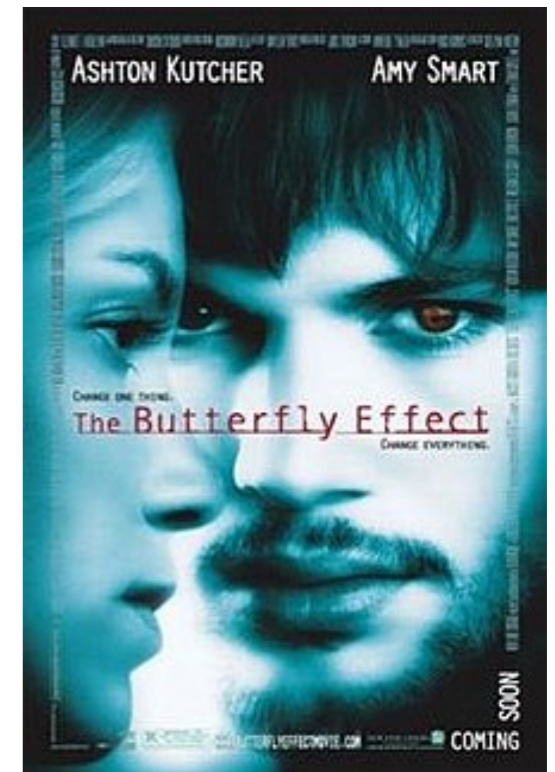


Goodfellow, P. McDaniel and N. Papernot. Making machine learning robust against adversarial inputs, CACM (2018)
Gilmer et al. Motivating the Rules of the Game for Adversarial Example Research, arXiv:1807.06732 (2018)

Formally

- **Exist** small $\Delta \mathbf{x}$ such that $f(\mathbf{x} + \Delta x) \neq y(\mathbf{x})$
- Practically, **exist** small $\Delta \mathbf{x}$ such that $f(\mathbf{x} + \Delta x) \neq f(\mathbf{x})$
 - similar if f is **very accurate**
 - such examples \mathbf{x} are called “adversarial”
- Intuitive explanation:

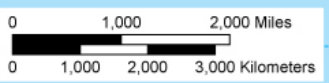
f is not sufficiently smooth (continuous)
- Or in fancier words, f is not **robust**



WORLD MAP



- | | | | |
|---------------------------|---------------------|-----------------------|----------------------------------|
| 1. Netherlands | 10. Austria | 20. Ghana | 29. Liechtenstein |
| 2. Belgium | 11. Hungary | 21. Togo | 30. Montenegro |
| 3. Luxembourg | 12. Serbia | 22. Benin | 31. Kosovo |
| 4. Switzerland | 13. Moldova | 23. Cameroon | 32. Palestinian Territories |
| 5. Slovenia | 14. North Macedonia | 24. Equatorial Guinea | 33. St. Vincent & the Grenadines |
| 6. Croatia | 15. Albania | 25. Rwanda | |
| 7. Bosnia and Herzegovina | 16. Cyprus | 26. Cambodia | |
| 8. Czechia | 17. Lebanon | 27. Panama | |
| 9. Slovakia | 18. Guinea-Bissau | 28. Malawi | |



**“Adversarial” examples exist for
any (non-constant) classifier**

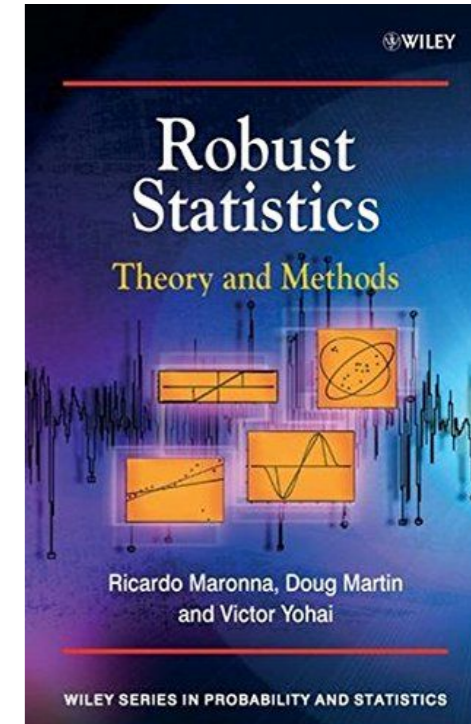
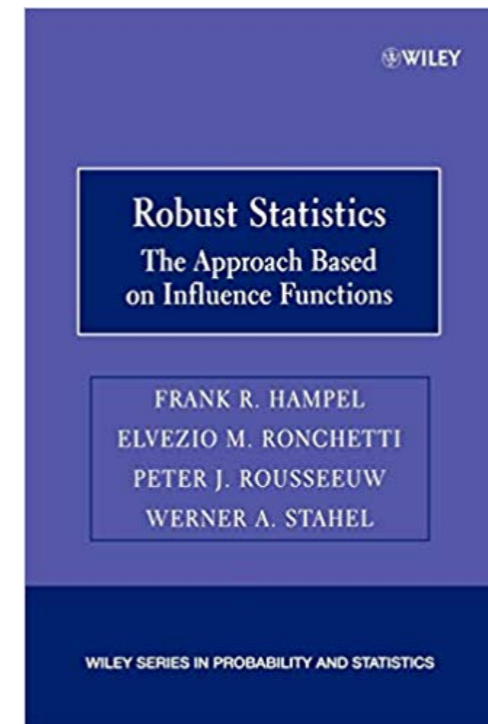
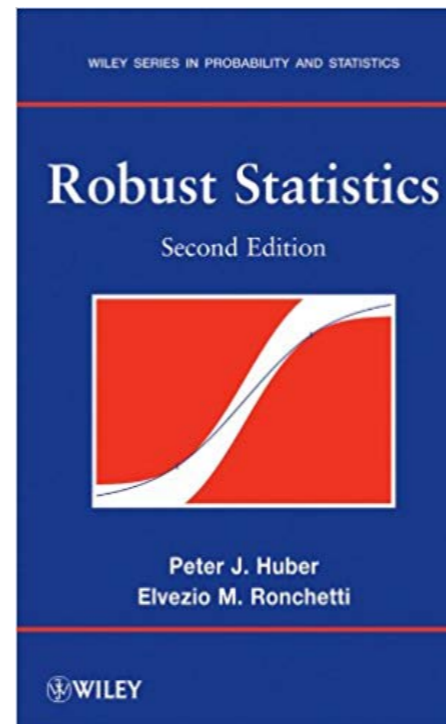
**Existence is not surprising;
universality is**

Robustness is not new

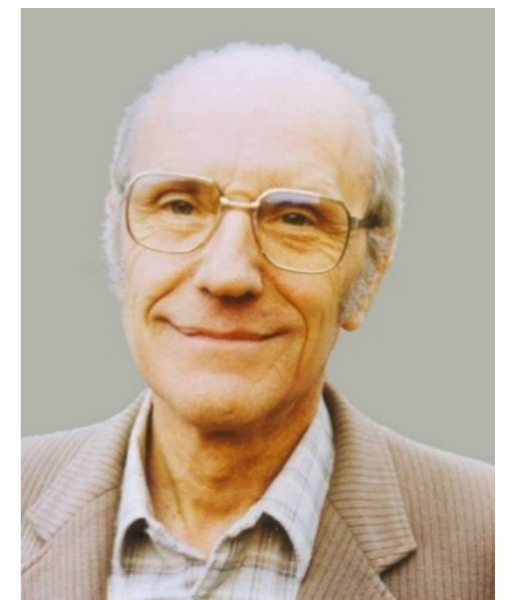
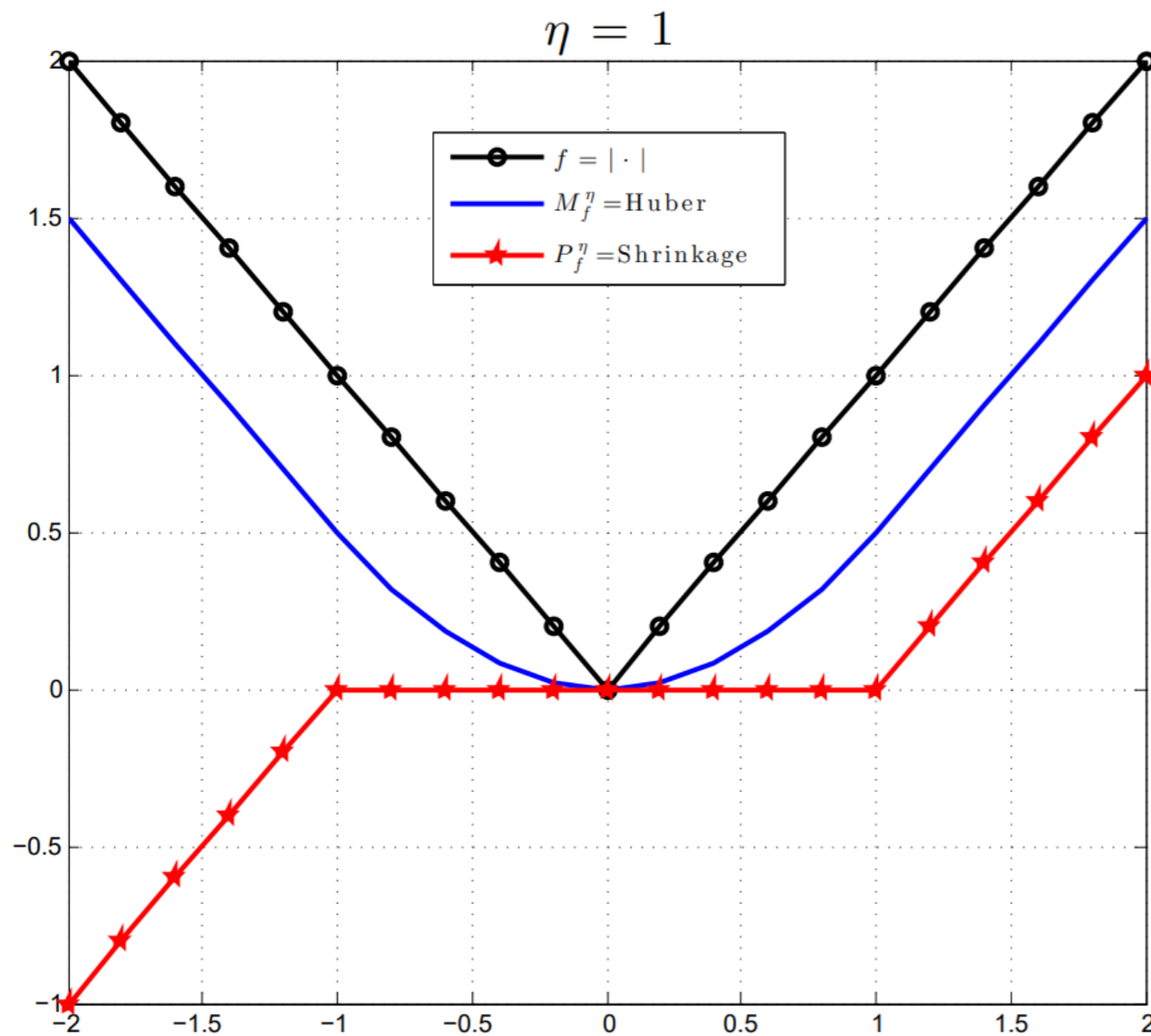
least-squares vs least absolute deviation

$$\|Xw - y\|_2 \quad \text{vs} \quad \|Xw - y\|_1$$

S. PORTNOY AND R. KOENKER



Huber's loss is Moreau's envelope

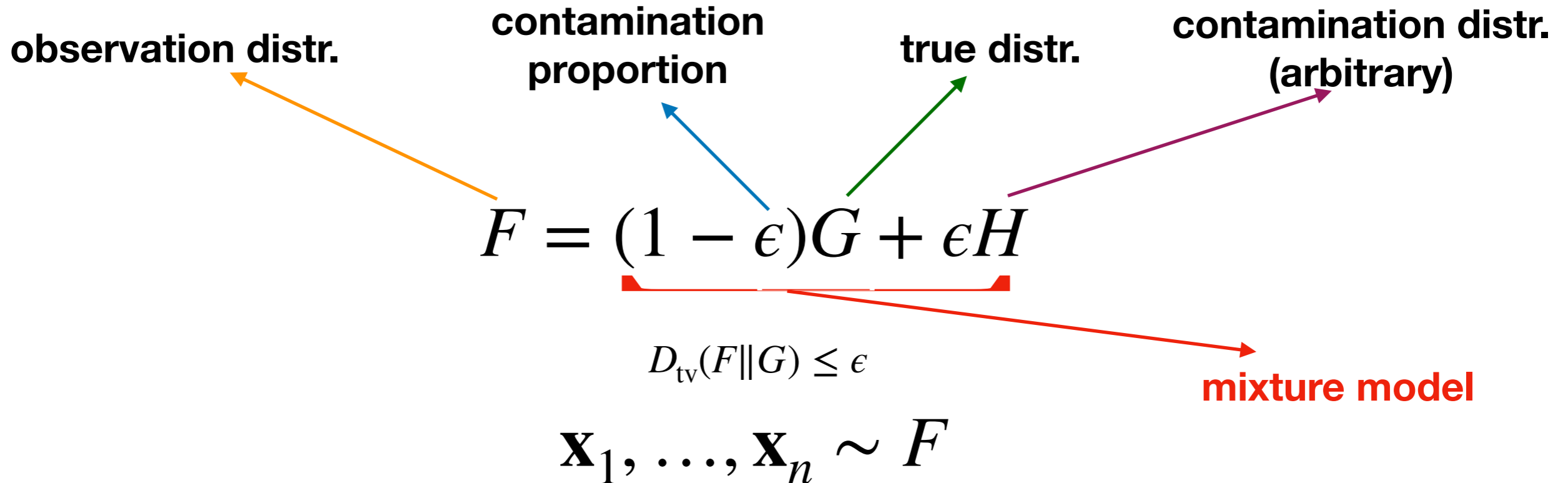


$$\begin{aligned} M_{|\cdot|}^\eta(t) &= \min_s \frac{1}{2}(s - t)^2 + \eta|s| \\ &= \begin{cases} \frac{1}{2}t^2, & |t| \leq \eta \\ \eta|t| - \frac{1}{2}\eta^2, & |t| \geq \eta \end{cases} \end{aligned}$$

Peter J. Huber. Robust estimation of a location parameter, *Annals of Statistics* (1964)

Jean J. Moreau. Fonctions convexes duales et points proximaux dans un espace hilbertien. *C.R.A.S.* (1962)

Huber's Contamination



- with probability (w.p.) ϵ , \mathbf{x}_i is from contamination H
- and w.p. $1-\epsilon$, \mathbf{x}_i is from true distribution
- so roughly ϵ proportion of training set is (arbitrarily) contaminated
- difficulty lies in don't know which data example is authentic or not

Score classifier

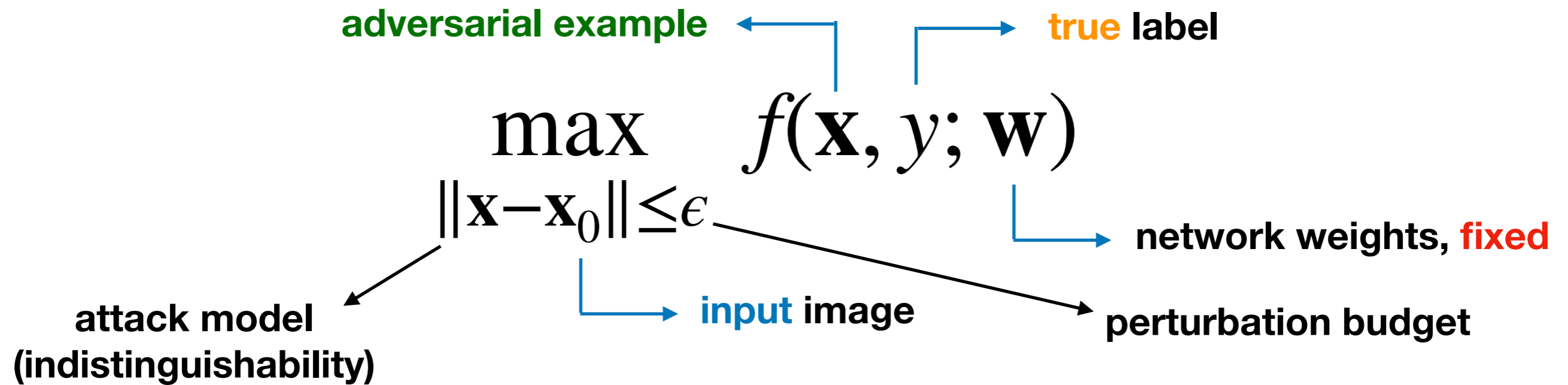
- A score function (classifier) $f : X \rightarrow \mathbb{R}$
- Input universe X is usually subset of \mathbb{R}^d
- Intuitively, small perturbation on x should result in small perturbation of $f(x)$
- Lipschitz continuity:

$$|f(\mathbf{x}) - f(\mathbf{z})| \leq L \cdot \|\mathbf{x} - \mathbf{z}\|$$

Lipschitz constant

$$\|\nabla f\| \leq L$$

Attack Algorithms



Fast Gradient Sign Method (FGSM)

$$\mathbf{x} \leftarrow \mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} f(\mathbf{x}, y; \mathbf{w}))$$

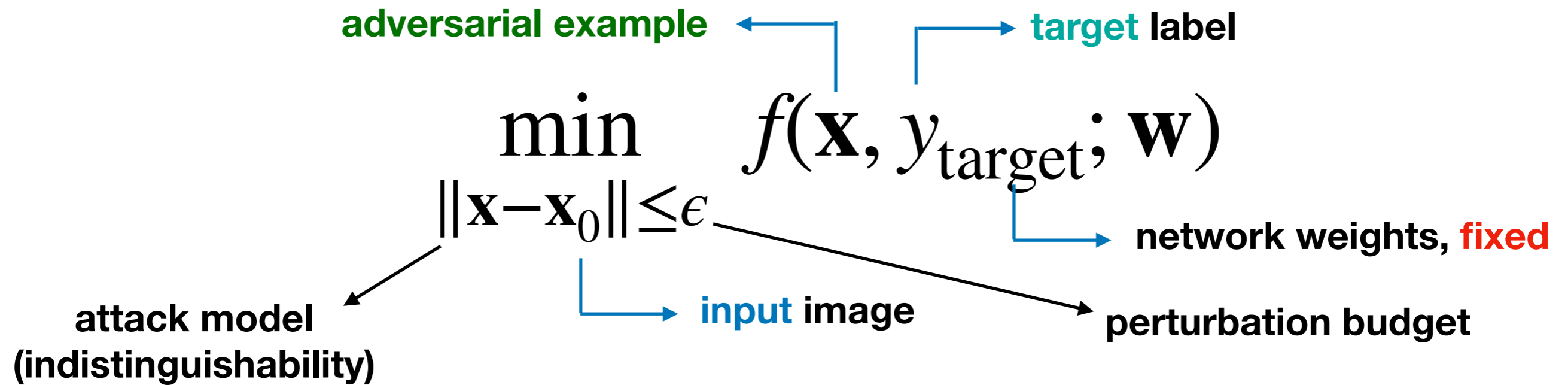
$$\mathbf{x} \leftarrow \text{Proj}(\mathbf{x})$$

Projected Gradient Method (PGM)

$$\mathbf{x} \leftarrow \mathbf{x} + \eta \cdot \nabla_{\mathbf{x}} f(\mathbf{x}, y; \mathbf{w})$$

$$\mathbf{x} \leftarrow \text{Proj}(\mathbf{x})$$

Targeted Attack



Fast Gradient Sign Method (FGSM)

$$\mathbf{x} \leftarrow \mathbf{x} - \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} f(\mathbf{x}, y_{\text{target}}; \mathbf{w}))$$

$$\mathbf{x} \leftarrow \text{Proj}(\mathbf{x})$$

Projected Gradient Method (PGM)

$$\mathbf{x} \leftarrow \mathbf{x} - \eta \cdot \nabla_{\mathbf{x}} f(\mathbf{x}, y_{\text{target}}; \mathbf{w})$$

$$\mathbf{x} \leftarrow \text{Proj}(\mathbf{x})$$

Lipschitz regularization

$$\mathcal{F}_\epsilon = \{F : W(F||G) \leq \epsilon\}$$

$$\max_{f: \text{Lip}(f) \leq 1} \mathbb{E}_{X \sim F} f(X) - \mathbb{E}_{Z \sim G} f(Z)$$

$$\min_{\mathbf{w}} \max_{F \in \mathcal{F}_\epsilon} \mathbb{E} \ell(\mathbf{w}^\top X), \quad X \sim F$$

||

$$\min_{\gamma \geq 0} \gamma \epsilon^p - \mathbb{E} M_{-\rho}^\gamma(X; \mathbf{w}), \quad \rho(\mathbf{x}) = \ell(\mathbf{w}^\top \mathbf{x}), \quad X \sim G$$

|∧

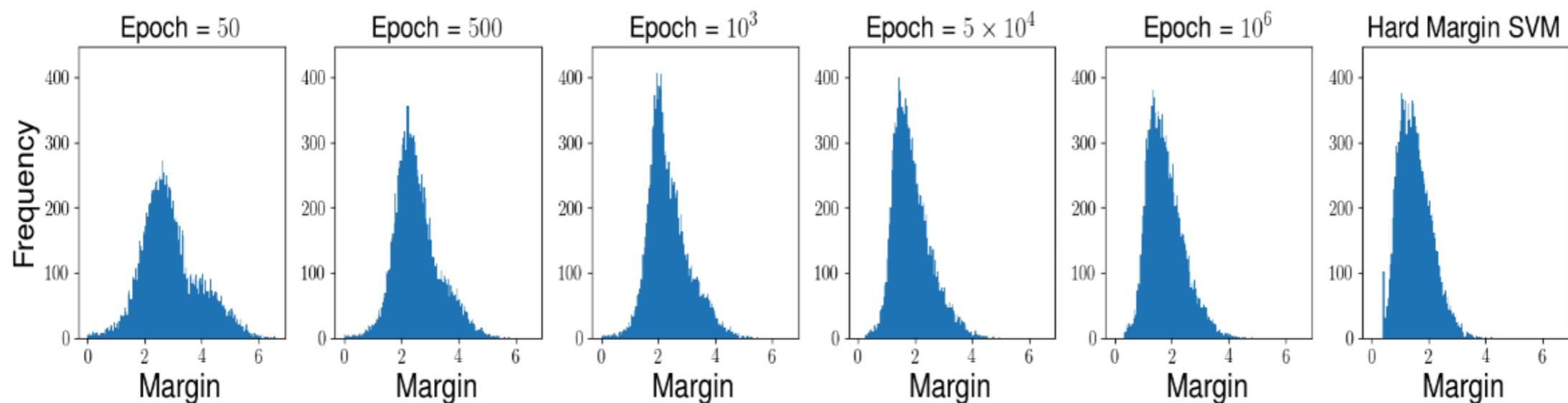
$$\mathbb{E} \ell(\mathbf{w}^\top \mathbf{x}) + \epsilon \cdot \text{Lip}(\ell_{\mathbf{w}}), \quad X \sim G$$

Margin story

$$m(\mathbf{x}, y; \{F_k\}) := \text{sign}(\hat{y}(\mathbf{x}), y) \cdot d(\mathbf{x}, \partial F_{\hat{y}})$$

Theorem 1 (Soudry et al. 2018 [3]) For almost all linearly separable binary datasets and any smooth decreasing loss with an exponential tail, **gradient descent** with small constant step size and any starting point \mathbf{w}_0 converges to the (unique) solution $\hat{\mathbf{w}}$ of hard-margin SVM, *i.e.*

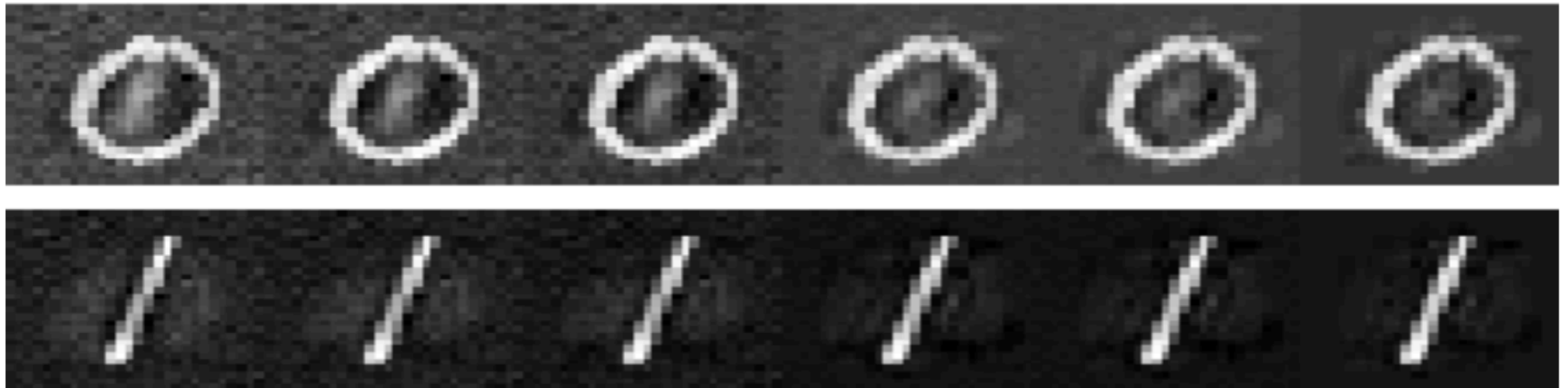
$$\lim_{t \rightarrow \infty} \frac{\mathbf{w}_t}{\|\mathbf{w}_t\|} = \frac{\hat{\mathbf{w}}}{\|\hat{\mathbf{w}}\|}. \quad \blacksquare$$



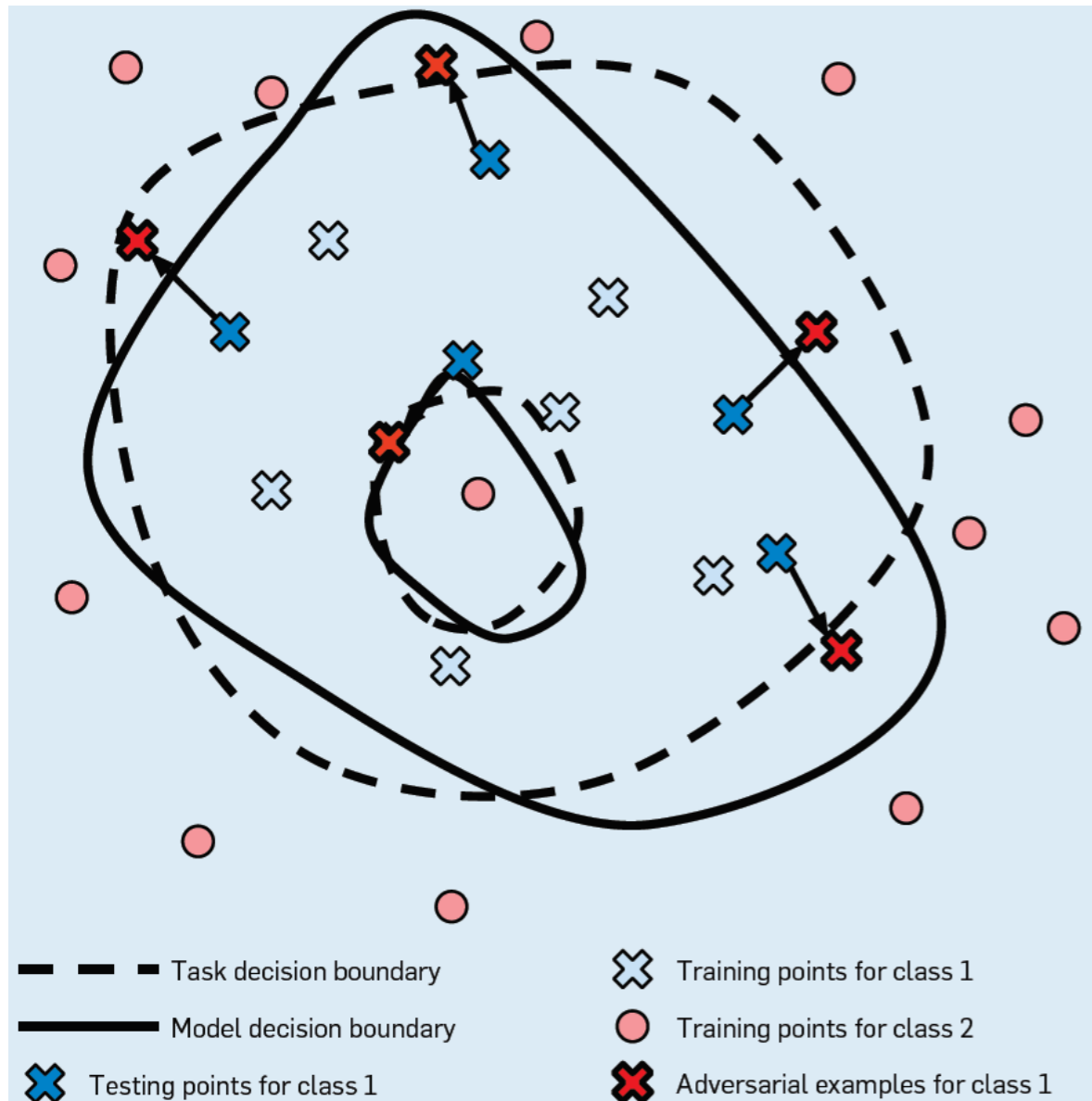
Binary linear classifiers

For a binary linear classifier which predicts positive if $\mathbf{w}^\top \mathbf{x} > 0$, can construct

$$\mathbf{x}^{\text{adv}} = \mathbf{x} - \frac{\mathbf{w}^\top \mathbf{x}}{\|\mathbf{w}\|^2} \mathbf{w}$$



The deep challenge



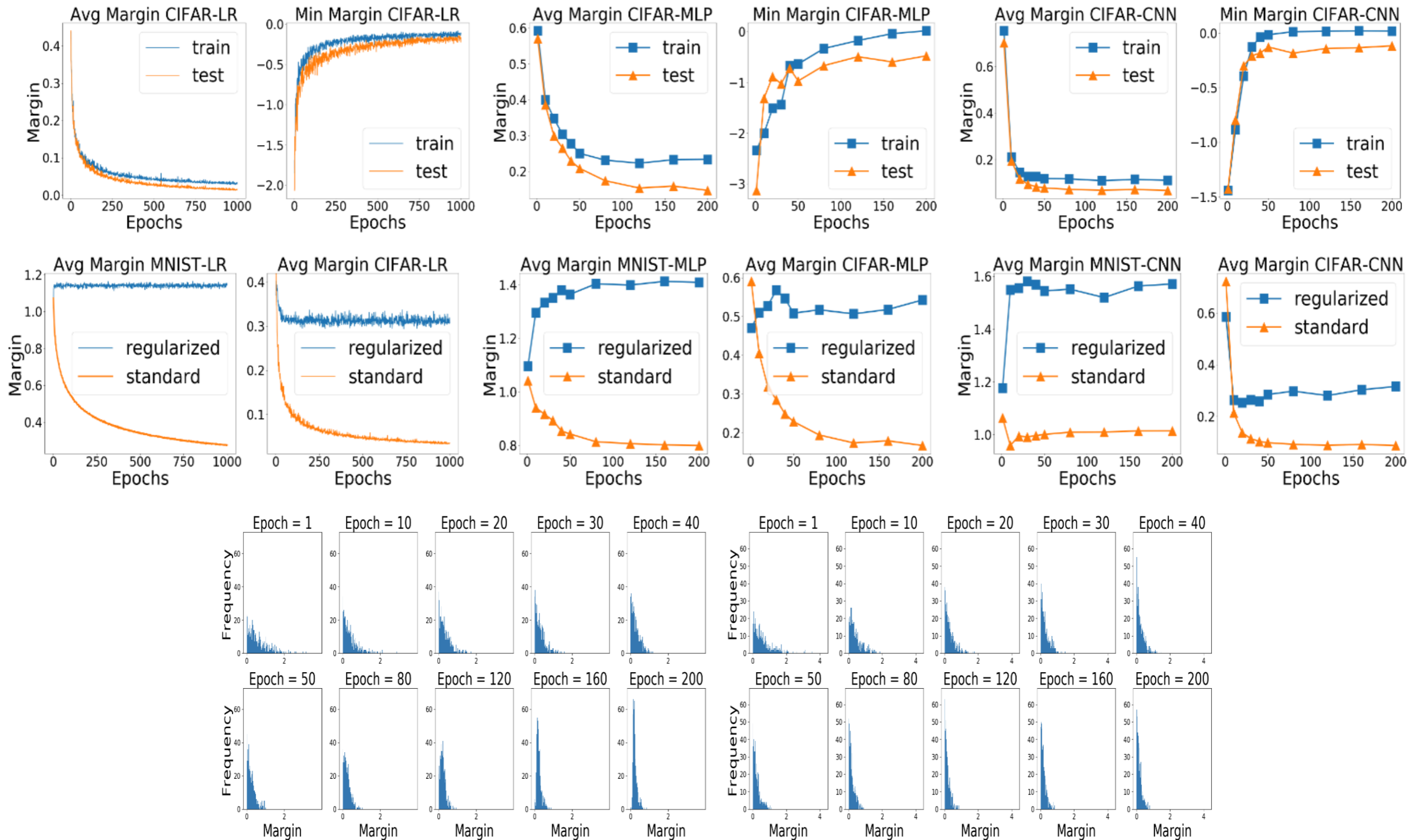
Deep Learning

$$\mathbf{x} \xrightarrow{\varphi} \varphi(\mathbf{x}) \xrightarrow{\mathbf{w}} \mathbf{w}^T \varphi(\mathbf{x}) =: \hat{y}$$

$$\min_{\mathbf{w}, \varphi} \ell(y, \hat{y})$$

- nonlinear transformation φ
- linear classifier w
- trained **jointly** by SGD

Observation & Solution?



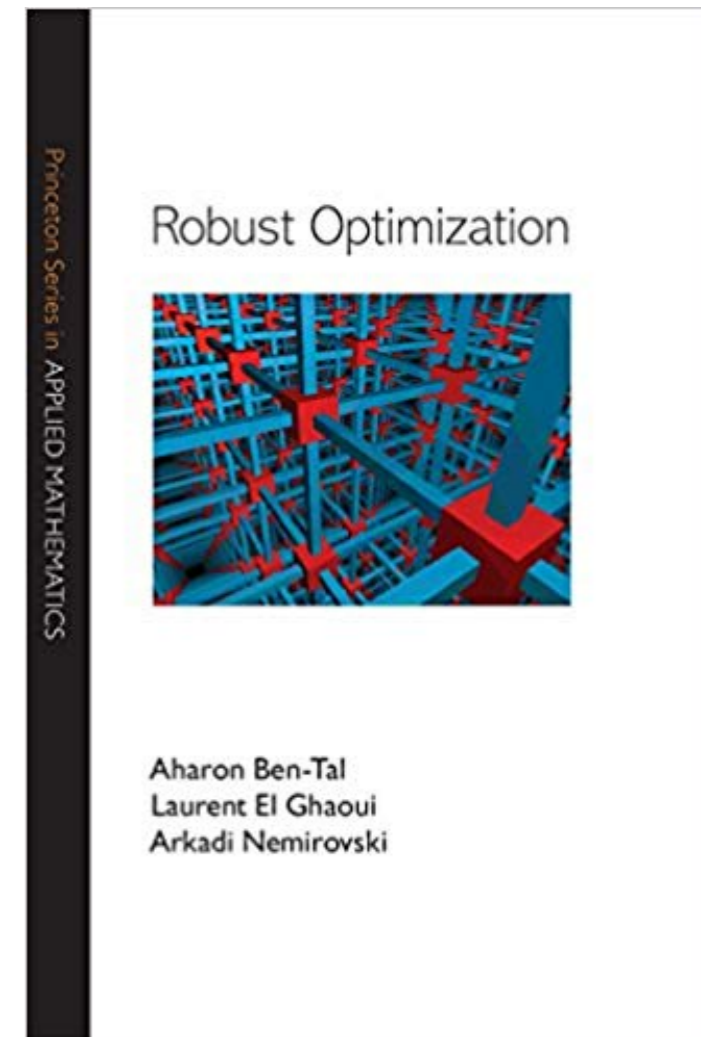
[WY]. Understanding Adversarial Robustness: The Trade-off between Minimum and Average Margin, (2019)

Adversarial training

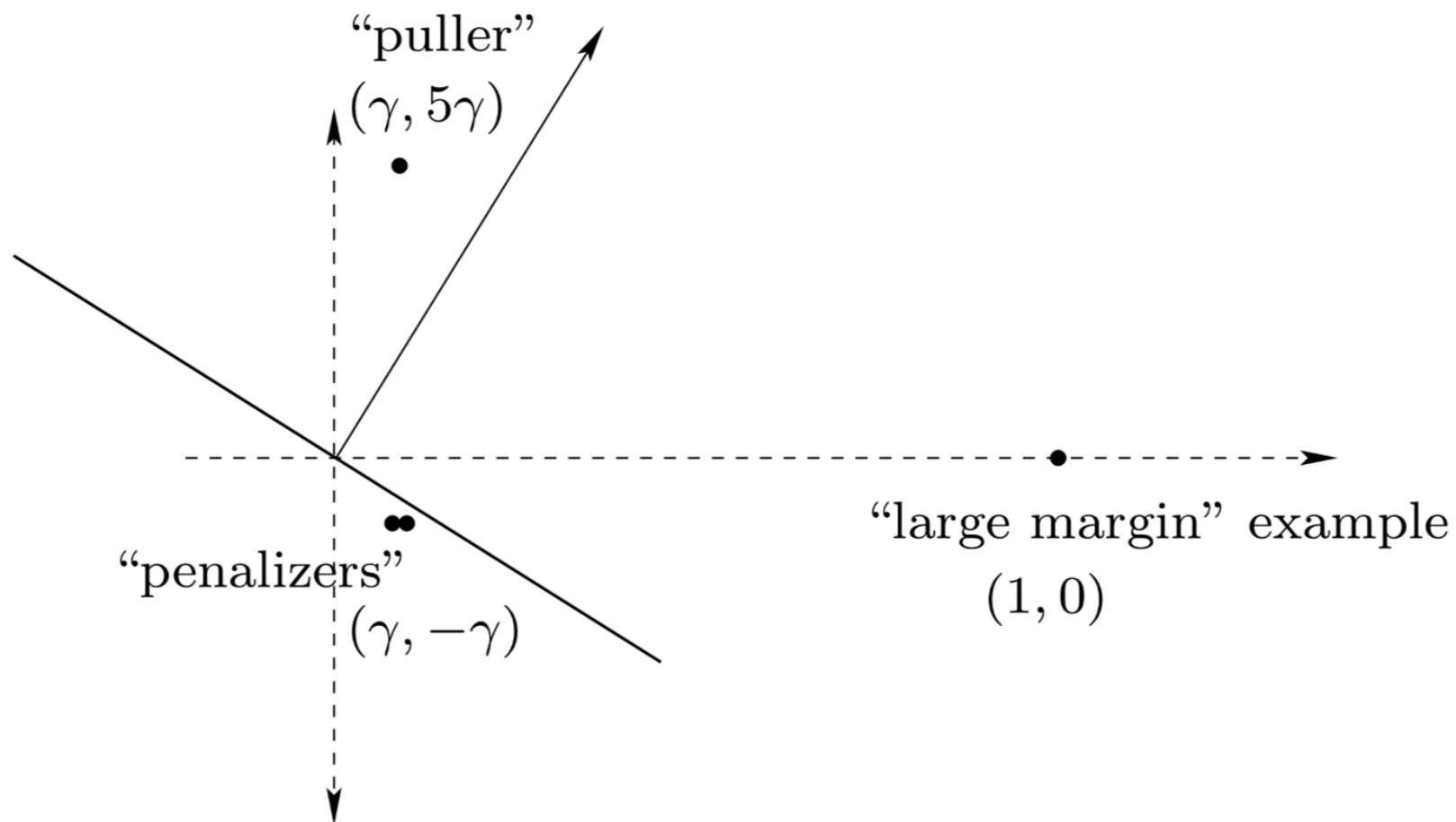
$$\min_{\mathbf{w}} \mathbf{E} \left[\max_{\|\Delta \mathbf{x}\| \leq \epsilon} \ell(\mathbf{x} + \Delta \mathbf{x}; \mathbf{w}) \right]$$

$\bar{\ell}_{\epsilon}(\mathbf{x}; \mathbf{w})$

- Amounts to changing the loss...
- One of the best defensive mechanisms
- Min-max formulation
- Inner max solved by an attack algorithm



A very negative result

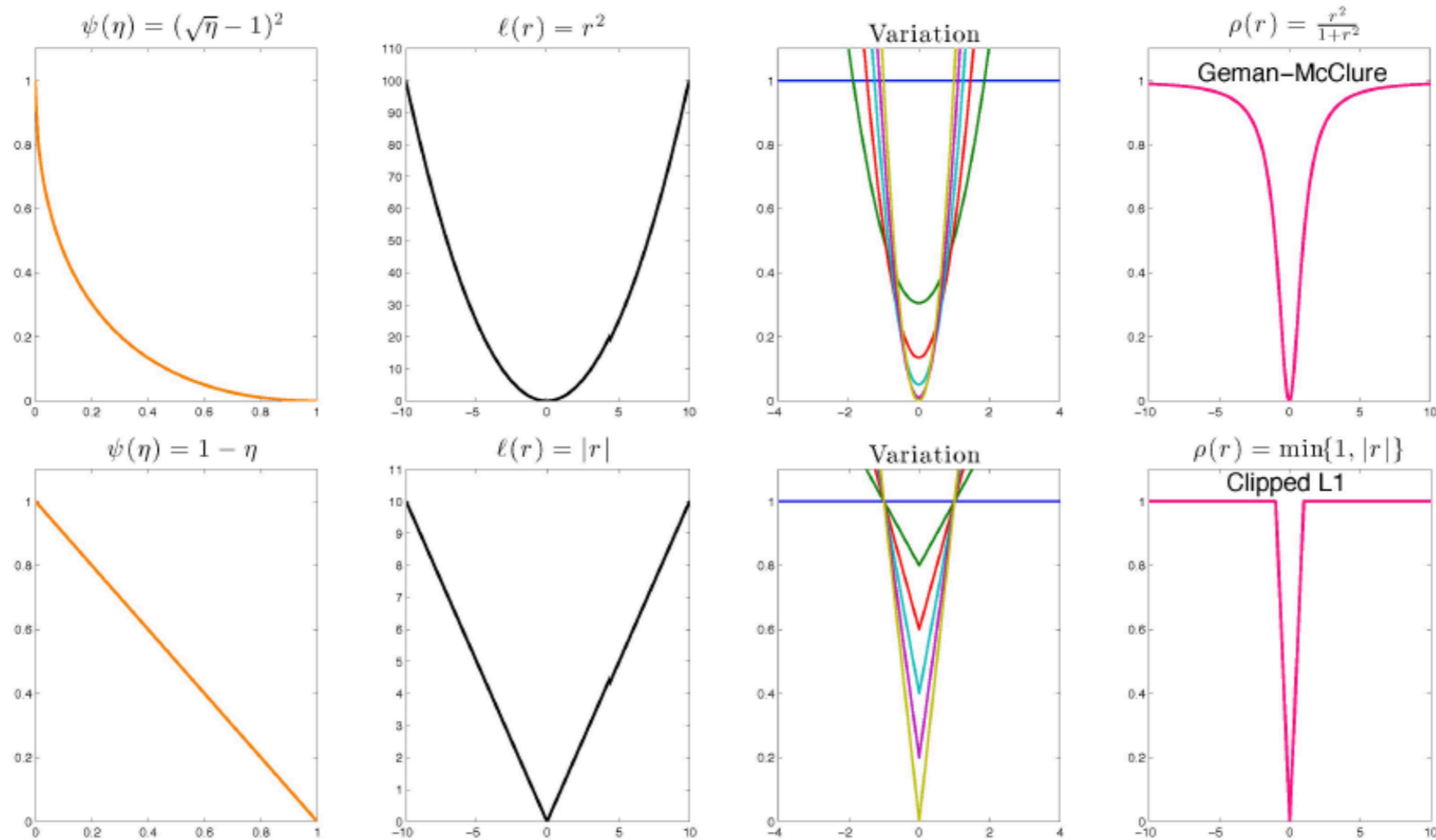


Theorem

*All convex potential function based boosters can **not** tolerate random classification noise at rate $\eta \in (0, 1/2)$.*

Variational loss

$$r(t) = \min_{0 \leq \eta \leq 1} \eta \ell(t) + \psi(\eta)$$



Certification

