

Explaining Deep Neural Networks Through Inverse Classification

PhD Defense

Shpresim Sadiku

October 13, 2025





Explaining DNNs Through Inverse Classification

| 1. | Introduction | 6 slides |
|------------|--|----------|
| 2. | GSE: Group-wise Sparse and Explainable Adversarial Attacks | 6 slides |
| 3. | S-CFE: Simple Counterfactual Explanations | 6 slides |
| 4. | Learning from Counterfactual Explanations | 6 slides |
| 5 . | Future Directions | 1 slide |

Goal. Enhance transparency/interpretability of ML models by providing intelligible justifications for decisions in high-stakes domains.

Interpretability. Degree to which a human can consistently predict the model's output. **Explainability.** Degree to which a human can understand the cause of a decision.

Why it matters.

- DNNs are accurate yet opaque ("black-box"). Trust, accountability, and governance require explanations.
- Distinguish prediction (model output) from prescription (human action).
- Aim: alignment of model reasoning with domain knowledge and real-world expectations.

1. Introduction

Inverse Classification and Adversarial Perturbations

Given a trained classifier f_{θ} and input $\mathbf{x} \in \mathbb{R}^d$ with $f_{\theta}(\mathbf{x}) = y$, **inverse classification** seeks a minimally modified $\tilde{\mathbf{x}} = \mathbf{x} + \mathbf{r}$ with desired label $\tilde{\mathbf{y}} \neq \mathbf{y}$:

$$\min_{\boldsymbol{r} \in \mathbb{R}^d} \mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{x} + \boldsymbol{r}), \tilde{\boldsymbol{y}}) \quad \text{s.t.} \quad \|\boldsymbol{r}\|_{p} \le \epsilon.$$
 (1)

Two major classes of such perturbations:

- Counterfactual Explanations (CFEs): human-oriented; prioritize plausibility, feasibility, and recourse.
- Adversarial Attacks: robustness evaluation; imperceptible yet effective (esp. in images).

Shared maths: both solve constrained optimization that minimally alters the decision; they differ in downstream *constraints* (plausibility vs. worst-case failure).

Explanations: Mapping and Desiderata

What is an explanation? A mapping $E(x, f_{\theta})$ to a human-interpretable object (textual, visual, symbolic).

Causal query. Why output y for input x? (e.g., loan rejection due to low credit score.)

Desirable properties

- Comprehensibility: understandable to non-experts.
- Stability: small input changes ⇒ similar explanations.
- Consistency: same input ⇒ similar explanations across runs/models.
- Realism: counterfactuals should be feasible/in-manifold.

Open question: what constitutes a "good" explanation remains unsettled.

Existence of Adversarial Perturbations (Theory)

For two-layer ReLU networks $f_{\theta}(\mathbf{x}) = \sum_{i=1}^{m} u_{i} \, \sigma(\mathbf{w}_{i}^{\top} \mathbf{x})$ with inputs near a subspace $P \subset \mathbb{R}^{d}$ (dim $P^{\perp} = \ell$):

• After training, the input gradient has a large P^{\perp} component with high probability:

$$\|\Pi_{P^{\perp}}(\partial f_{\theta}/\partial x)\| \geq \sqrt{\frac{k\ell}{2md}},$$
 (2)

where $k = |\{active neurons\}|$.

• There exists a universal $r \in P^{\perp}$ such that $\mathrm{sign}(f_{\theta}(\mathbf{x}+\mathbf{r})) \neq \mathrm{sign}(f_{\theta}(\mathbf{x}))$ and

$$\|\mathbf{r}\| \leq \mathcal{O}\left(f_{\theta}(\mathbf{x}) \cdot \sqrt{\frac{m}{k_{y}}} \cdot \sqrt{\frac{d}{\ell}}\right),$$
 (3)

with k_v neurons aligned with label y.

Takeaway. Even small off-manifold perturbations can flip decisions; universal directions may exist.

Adversarial and Counterfactual Methods

FGSM (targeted/untargeted).

$$\tilde{\mathbf{x}} = \mathbf{x} - \epsilon \operatorname{sign}(\nabla_{\mathbf{x}} \mathcal{L}(f_{\theta}(\mathbf{x}), \tilde{\mathbf{y}})). \tag{4}$$

PGD (ℓ_{∞} -bounded).

$$\mathbf{x}_{t+1} = \Pi_{\mathcal{B}_{\epsilon}^{\infty}(\mathbf{x})} \Big(\mathbf{x}_{t} - \alpha \operatorname{sign}(\nabla_{\mathbf{x}} \mathcal{L}(f_{\theta}(\mathbf{x}_{t}), \tilde{\mathbf{y}})) \Big).$$
 (5)

CFEs. Solve variants of (1) with *plausibility* constraints/regularisation:

- Distance-based (e.g., Manhattan/Mahalanobis), actionability constraints.
- Density-aware: hard constraints via GMM components; or soft regularisers (e.g., LOF-based).

Structured attacks (images). Group-sparse (e.g., ADMM-based), nuclear-group norms, homotopy sparse attacks.

1. Introduction

Optimization Toolbox for Perturbations

First-order methods.

■ GD/SGD: $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla_{\mathbf{x}} s_{\theta}(\mathbf{x}_t)$; momentum and Nesterov variants improve stability/speed.

Proximal gradient (PG).

$$\mathbf{x}_{t+1} = \operatorname{prox}_{\lambda g} \left(\mathbf{x}_t - \lambda \nabla h(\mathbf{x}_t) \right), \quad \lambda \approx 1/L.$$
 (6)

Acceleration (FISTA). Nesterov-type extrapolation yields $\mathcal{O}(1/t^2)$ in convex settings. Thresholding operators.

- *Soft* (ℓ_1) and *hard* (ℓ_0) thresholding (closed forms).
- Nonconvex $\ell_{1/2}$: explicit proximal update using $\phi_{2\lambda}(x_i)$ and threshold $g(2\lambda)$.

Message. These tools enable principled trade-offs between imperceptibility, sparsity/structure, and target success.



Explaining DNNs Through Inverse Classification

| 1. | Introduction | 6 slides |
|----|--|----------|
| 2. | GSE: Group-wise Sparse and Explainable Adversarial Attacks | 6 slides |
| 3. | S-CFE: Simple Counterfactual Explanations | 6 slides |
| 4. | Learning from Counterfactual Explanations | 6 slides |
| 5. | Future Directions | 1 slide |



Introduction and Motivation

- DNNs are vulnerable to adversarial perturbations across tasks: classification, captioning, retrieval, QA, autonomous driving, face recognition/detection, etc. (e.g., Carlini et al. 2017, Athalye et al. 2018, Zhang et al. 2020).
- Beyond ℓ_p with $p \ge 1$, the p = 0 (sparse) regime is compelling: few pixels changed without constraints on where and by how much \Rightarrow perceptible artifacts (Su et al. 2019).
- Need: impose structure ⇒ group-wise sparse perturbations targeted at the object of interest (Xu et al. 2018, Zhu et al. 2021, Imtiaz et al. 2022, Kazemi et al. 2023).

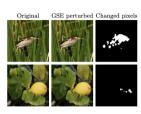


Figure: Adversarial attacks generated by GSE algorithm.

 Bridges human perception vs. machine features (Ilyas et al. 2019); perturbations become explainable.



Contributions (GSE)

Two-phase algorithm for group-wise sparse, low-magnitude, explainable attacks.

- 1. Phase I (Selection). Non-convex regularisation with proximal splitting + a proximity-based update of per-pixel tradeoff parameters λ to select salient pixel groups.
- 2. **Phase II (Refinement).** Nesterov's accelerated gradient (projected onto selected coordinates) with ℓ_2 -regularisation to minimise perturbation magnitude.
- 3. **Empirics.** CIFAR-10 and ImageNet: up to **50.9%** (CIFAR-10) and **38.4%** (ImageNet) higher group-wise sparsity (targeted, average case) at **100%** ASR.

Evaluation: ASR; sparsity (ACP), grouping (ANC, $d_{2,0}$), magnitude (ℓ_2), explainability (ASM-based IS), runtime.

Explainability. Quantitatively aligns perturbations with salient regions (ASM/CAM), outperforming SOTA sparse and group-wise sparse attacks.



Related Work (Sparse & Group-wise Sparse Attacks)

- Sparse (p=0): one-pixel (Su et al. 2019); local search (Narodytska et al. 2016); evolutionary methods (Croce et al. 2019); ℓ₁ relaxations, e.g., SparseFool (Modas et al. 2019). Often perceptible; location/magnitude unconstrained.
- Group-wise sparse:
 - StrAttack (Xu et al. 2018): ADMM with sliding masks.
 - SAPF (Fan et al. 2020): ℓ_p -Box ADMM with binary selections.
 - Homotopy-Attack (Zhu et al. 2021): nmAPG; SLIC-based 2, 0—'norm' regularisation.
 - FWnucl (Kazemi et al. 2023): nuclear group norm.
- Contrary to benchmarks, GSE method does not depend on pixel partitionings.
- Links to explanations: hitting-set duality on MNIST (Ignatiev et al. 2019); perturbations trace discriminative features (Xu et al. 2018).

Adversarial Attack Formulation

Feasible images: $\mathcal{X} = [I_{\min}, I_{\max}]^{M \times N \times C}$. Benign image $\mathbf{x} \in \mathcal{X}$ with label $y \in \mathbb{N}$, target $\tilde{y} \in \mathbb{N}$ $(\tilde{y} \neq y)$. Classifier f_{θ} and loss \mathcal{L} .

$$\min_{\mathbf{r} \in \mathbb{R}^{M \times N \times C}} \mathcal{L}(f_{\theta}(\mathbf{x} + \mathbf{r}), \tilde{\mathbf{y}}) + \lambda \mathcal{D}(\mathbf{r}). \tag{7}$$

$$\max_{\boldsymbol{r} \in \mathbb{R}^{M \times N \times C}} \mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{x} + \boldsymbol{r}), y) - \lambda \mathcal{D}(\boldsymbol{r}). \tag{8}$$

Sparse regularisation: $\mathcal{D}(\cdot) = \|\cdot\|_p^p$, 0 .

1/2-Quasinorm Regularisation and FBS

Quasinorm-regularised objective (sparse attacks):

$$\min_{\mathbf{r}} \mathcal{L}(f_{\theta}(\mathbf{x} + \mathbf{r}), \mathbf{y}) + \lambda \|\mathbf{r}\|_{p}^{p}, \quad 0$$

For $p = \frac{1}{2}$, the proximal operator admits a **closed form** (component-wise).

Algorithm Forward–Backward Splitting Attack (sketch)

- 1: Initialise $r_0 \leftarrow \mathbf{0}$
- 2: **for** t = 0, ..., T 1 **do**
- 3: $r_{t+1} \leftarrow \operatorname{prox}_{\alpha_t \lambda \|\cdot\|_p^p} \left(r_t \alpha_t \nabla_r \mathcal{L}(f_{\theta}(\mathbf{x} + r_t), y) \right)$
- 4: end for
- 5: Return $\tilde{\mathbf{r}} = \mathbf{r}_T$

Limitation: yields very sparse but often **large-magnitude** and poorly localised perturbations (Fan et al. 2020).

GSE: Group-wise Sparse, Low-Magnitude Attack

Phase I (Select coordinates):

- Use a per-pixel vector $\lambda \in \mathbb{R}^{M \times N \times C}_{\geq 0}$ in the $\frac{1}{2}$ -quasinorm proximal step.
- Build $m = \text{sign}\left(\sum_{c=1}^{C} |r_t|_{:,:,c}\right)$; blur with a Gaussian kernel K to obtain M = m * *K.
- Form $\overline{\pmb{M}}$ via $\overline{M}_{ij}=M_{ij}+1$ if $M_{ij}\neq 0$, else $q\in (0,1];$ update $\lambda_{t+1}^{i,j,:}=\frac{1}{\overline{M}_{ii}}\,\lambda_t^{i,j,:}.$
- After \hat{t} iters, define selected subspace $V = \text{span}\{e_{i,j,c} \mid \lambda_{\hat{t}}^{i,j,c} < \lambda_{0}^{i,j,c}\}$.

Phase II (Refine on V):

$$\min_{\mathbf{r} \in V} \mathcal{L}(f_{\theta}(\mathbf{x} + \mathbf{r}), y) + \mu \|\mathbf{r}\|_{2}, \text{ solve by projected NAG.}$$
 (10)

Lemma (Sadiku, Wagner, and Pokutta 2025)

The projected NAG solving Eq. (10) converges as NAG solving an unconstrained problem.



Explaining DNNs Through Inverse Classification

| 1. | Introduction | 6 slides |
|------------|--|----------|
| 2. | GSE: Group-wise Sparse and Explainable Adversarial Attacks | 6 slides |
| 3. | S-CFE: Simple Counterfactual Explanations | 6 slides |
| 4. | Learning from Counterfactual Explanations | 6 slides |
| 5 . | Future Directions | 1 slide |



Counterfactual Explanations (CFEs): Motivation

- ML systems operate in high-stakes domains (finance, healthcare, justice, hiring). Opacity ⇒ transparency, fairness, accountability concerns.
- CFEs answer what-if: minimal (feasible) changes to flip the decision to a target label (Wachter et al. 2017).
- Contrast with LRP/LIME: attribution of present features (Bach et al. 2015, Ribeiro et al. 2016) vs. CFEs identify absent features whose presence would change the outcome.

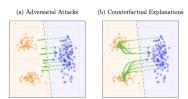


Figure: (a) Without the plausibility term, points cluster near the blue factual data but far from the orange distribution. (b) With the plausibility term, points lie in high-density regions. The dashed line shows the linear decision boundary.



Principles: Proximity, Validity, Actionability, Plausibility, Sparsity

Basic principles.

- Proximity (small ℓ_2 distance to factual) and Validity $(f_{\theta}(\tilde{\mathbf{x}}) = \tilde{\mathbf{y}})$.
- Actionability: respect feature ranges; avoid impossible edits.
- Plausibility: move toward target class manifold (not merely across boundary).
- Sparsity: change as few features as possible (short explanations are preferred (Mothilal et al. 2020, Naumann et al. 2021)).

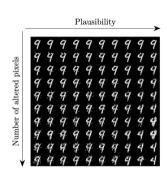


Figure: CFEs for changing $9 \rightarrow 4$: sparsity alone gives adversarial results, plausibility gives unrealistic ones, combining both yields sparse and realistic CFEs.

Canonical CFE Optimisation and Challenges

Canonical form for a factual **x** (conceptual):

$$\min_{x' \in \text{actionable set}} \left[\underbrace{\text{CFE loss}}_{\text{validity}} + \underbrace{\text{dist}(x', x)}_{\text{proximity}} + \underbrace{\text{dist to manifold}}_{\text{plausibility}} + \underbrace{\#\text{changes}}_{\text{sparsity}} \right]. \tag{11}$$

Difficulties.

- Nonconvex classifier losses; non-smooth sparsity terms (e.g., ℓ_0); complex manifold penalties; box constraints.
- Prior work tackles subsets: linear/trees with GMM constraints (Artelt et al. 2020); ReLU MIP with LOF (Tsiourvas et al. 2024); density-regularised relaxations (Zhang et al. 2023).

- Early CFEs: weighted ℓ_1 /Mahalanobis for sparsity and proximity (Wachter et al. 2017, Verma et al. 2024, Karimi et al. 2020).
- DNNs with VAEs for plausibility (CEM) (Dhurandhar et al. 2018); density-based plausibility with elastic-net (DCFE) (Zhang et al. 2023).
- Convex/GMM approach for simple classifiers (PCFE) (Artelt et al. 2020).
- MIP over ReLU polytopes with LOF constraint (limited to ReLU) (Tsiourvas et al. 2024).

S-CFE: A Simple APG Framework (FISTA-style)

Relaxed objective (penalty form):

$$\min_{\mathbf{x}'} h(\mathbf{x}', \tilde{\mathbf{y}}) + g_{\rho}(\mathbf{x}'), \quad h = \|\mathbf{x}' - \mathbf{x}\|_{2}^{2} + \gamma \mathcal{L}(f_{\theta}(\mathbf{x}'), \tilde{\mathbf{y}}) - \tau \, \hat{q}(\mathbf{x}', \tilde{\mathbf{y}}), \tag{12}$$

$$g_{\rho} = I_{\mathcal{A}}(\mathbf{x}') + \beta \|\mathbf{x}' - \mathbf{x}\|_{\rho}^{\rho}, \quad \rho \in \{\frac{1}{2}, \frac{2}{3}, 1\}$$
 (13)

APG step (cf. FISTA):

$$\mathbf{x}'_{t+1} = \operatorname{prox}_{\sigma_t \mathbf{g}_p} \left(\mathbf{z}_t - \sigma_t \nabla h(\mathbf{z}_t, \tilde{\mathbf{y}}) \right), \quad \mathbf{z}_{t+1} = \mathbf{x}'_{t+1} + \alpha_t (\mathbf{x}'_{t+1} - \mathbf{x}'_t). \tag{14}$$

Plausibility choices: differentiable $\hat{q} \in \{\hat{q}_{KDE}, \hat{q}_{GMM}, \hat{q}_{kNN}\}.$

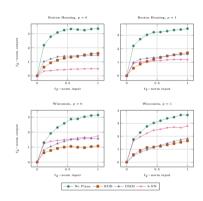
Sparsity control (constrained form):

$$g_0 = I_{\mathcal{A}} + \beta I_{\|\mathbf{x}' - \mathbf{x}\|_0 \le m} \Rightarrow \text{prox} = \text{projection onto } \{\|\mathbf{x}' - \mathbf{x}\|_0 \le m\} \cap \mathcal{A}. \tag{15}$$

Empirical Highlights and Robustness

Setup. Boston Housing, Wine, MNIST; logistic/DNN/CNN classifiers; metrics: Validity (%), proximity (ℓ_2), sparsity (ℓ_0), plausibility (LOF), runtime.

- S-CFE variants
 (KDE/GMM/kNN) produce
 sparse (ℓ₀ controlled), plausible
 (low LOF) CFEs with strong
 validity—while keeping proximity
 and runtime competitive.
- Projection onto $\{\|x'-x\|_0 \le m\} \cap \mathcal{A} \text{ offers}$ explicit sparsity control; density terms steer toward target manifolds.



 Robustness: plausibility constraints improve stability to small input shifts; promotes individual fairness.



Explaining DNNs Through Inverse Classification

| 1. | Introduction | 6 slides |
|----|--|----------|
| 2. | GSE: Group-wise Sparse and Explainable Adversarial Attacks | 6 slides |
| 3. | S-CFE: Simple Counterfactual Explanations | 6 slides |
| 4. | Learning from Counterfactual Explanations | 6 slides |
| 5. | Future Directions | 1 slide |



Learning from *plausible* counterfactuals (p-CFEs): Why?

- **Goal:** Flip a model's prediction via *minimal* input changes.
- Two worlds: <u>Adversarial attacks</u> vs. <u>p-CFEs</u> (plausible, manifold-aligned, interpretable).
- Recent theory: adversarial perturbations contain generalizable, class-specific features (Ilyas et al. 2019, Kumano et al. 2024).

Question: Do p-CFEs share this representational richness? And can they be *better* for learning—especially under spurious correlations?

• Claim: Training on p-CFEs attains competitive accuracy and mitigates spurious correlations (strong WGA gains).



Contributions

- 1. **Learning from p-CFEs:** Extend the *learning from perturbations* paradigm from adversarial examples to *plausible* counterfactuals.
- 2. **Accuracy:** Models trained on p-CFEs reach test accuracy comparable to models trained on adversarial examples (PGD ℓ_2 , ℓ_∞) and CFE- ℓ_2 .
- 3. **Spurious correlations:** p-CFE training substantially improves worst-group accuracy (WGA); on WaterBirds it **surpasses** standard training by $\approx 12\%$.

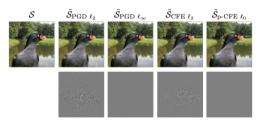


Figure: Random WaterBirds samples with perturbations (\times 40) targeting landbird labels from true waterbirds.

Learning from perturbations: setup & objectives

p-CFE (targeted)

Definition (Learning from perturbations)

Given a dataset $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, create a perturbed set $\tilde{S} = \{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^n$ by targeting labels $\tilde{y}_i \neq y_i$; then train a new model on \tilde{S} and evaluate on the clean test set.

PGD (targeted)

$$\min_{\tilde{\mathbf{x}}} \mathcal{L}(f_{\theta}(\tilde{\mathbf{x}}), \tilde{\mathbf{y}})$$
s.t. $\|\tilde{\mathbf{x}} - \mathbf{x}\|_{p} \leq \epsilon$.
$$\tilde{\mathbf{x}} = \arg\min_{\mathbf{x}' \in \mathcal{A}} \left\{ \|\mathbf{x}' - \mathbf{x}\|_{2}^{2} + \gamma \mathcal{L}(f_{\theta}(\mathbf{x}'), \tilde{\mathbf{y}}) - \tau \, \hat{q}(\mathbf{x}', \tilde{\mathbf{y}}) + \beta \, \|\mathbf{x}' - \mathbf{x}\|_{0} \right\}.$$

• Key distinction: the **plausibility** term $(-\tau \hat{q})$ pulls counterfactuals toward the target-class manifold; ℓ_0 promotes **sparsity**.



Experimental setup: data, training, metrics

- Datasets with spurious correlations:
 - WaterBirds (Sagawa et al. 2019): label (land vs. water) spuriously correlates with background.
 - SpuCoAnimals (Joshi et al. 2023): big vs. small dogs spuriously correlate with indoor/outdoor.

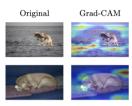


Figure: Grad-CAM visualizations show misclassifications: a landbird on water labeled as a waterbird and a big dog indoors as a small dog.

- **Training:** Fine-tune ResNet50 on perturbed sets (PGD- ℓ_2 , PGD- ℓ_∞ , CFE- ℓ_2 , p-CFE- ℓ_0); target labels \tilde{y} chosen uniformly at random.
- Metrics: Train/Test accuracy.
- Worst-Group Accuracy (WGA) to quantify spurious reliance.



Results: Accuracy and Worst-Group Accuracy (WGA)

Test accuracy (%)

| | PGD- ℓ_2 | PGD- ℓ_∞ | $CFE	ext{-}\ell_2$ | p-CFE | Orig. |
|----------------------------|----------------|--------------------|--------------------|-----------------------|-------|
| WaterBirds SpuCoAnimals | 86.08 78.10 | 86.02 79.43 | 88.58 79.00 | 86.54 81.78 | |

Worst-Group Acc. (%)

| | $PGD\text{-}\ell_2$ | $PGD\text{-}\ell_\infty$ | $CFE\text{-}\ell_2$ | p-CFE | Orig. |
|--------------|---------------------|--------------------------|---------------------|-------|-------|
| WaterBirds | 56.58 | 61.72 | 63.04 | 76.05 | |
| SpuCoAnimals | 56.06 | 57.53 | 56.60 | 63.53 | 65.60 |

■ Takeaways. p-CFE training: (i) matches adversarial/CFE- ℓ_2 on accuracy; (ii) strongly mitigates spurious correlations— $\pm 11-12\%$ WGA vs. standard training on WaterBirds.



Qualitative evidence (Grad-CAM) & conclusions

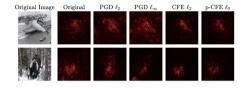


Figure: Saliency maps for a landbird and dog: original, standard, PGD (ℓ_2, ℓ_∞) , CFE (ℓ_2) , and p-CFE (ℓ_0) models.

Observed focus (Grad-CAM):

- Standard/PGD/CFE-ℓ₂ tend to over-weight background.
- p-CFE shifts attention to semantic object (bird/dog).
- Simple, model-agnostic recipe—no group labels needed

Conclusions

- p-CFEs are effective training signals: accurate & robust to spurious cues.
- Manifold alignment (plausibility) steers learning toward <u>semantic</u> features.



Explaining DNNs Through Inverse Classification

| 5 . | Future Directions | 1 slide |
|------------|--|----------|
| 4. | Learning from Counterfactual Explanations | 6 slides |
| 3. | S-CFE: Simple Counterfactual Explanations | 6 slides |
| 2. | GSE: Group-wise Sparse and Explainable Adversarial Attacks | 6 slides |
| 1. | Introduction | 6 slides |

Adversarial Training with GSE

- Use GSE examples in adversarial training.
- Report robustness-sparsity-explainability-time trade-offs.

S-CFE: Method

- From predictor acceptance → outcome improvement (causal constraints).
- Train on data-level targets, not only model loss.

Model Shifts / Black-Box

- Test KDE / density-gravity plausibility under model change.
- Black-box CFEs: finite-diff or surrogate; consider validity-free variant (accuracy trade-off).

Mixed & Categorical Features

 Design discrete prox/projection (beyond one-hot + APG).

High-Dimensional CFEs

 Swap KDE/GMM for differentiable VAEs/flows; stabilize gradients in q̂.

Learning from p-CFEs @ Scale

 Extend to LLMs/VLMs; connect with theory of learning from perturbations.

Adversarial ⇔ p-CFE

- Conjecture: on robust models, targeted attacks ≈ manifold-aligned p-CFEs.
- Diagnostic: angle between attack and p-CFE directions.



Thank you for your attention!

6. Appendices - GSE

From FISTA to NAG when g = 0

• Set g = 0. Then $\text{prox}_{\alpha g} = \text{Id}$, so the update rule of FISTA becomes a plain gradient step at the look-ahead point \mathbf{y}_t :

$$\mathbf{x}_{t+1} = \mathbf{y}_t - \alpha \nabla f(\mathbf{y}_t). \tag{1}$$

The extrapolation coefficient is

$$\mu_{t+1} := \frac{\beta_t - 1}{\beta_{t+1}} \implies \mathbf{y}_{t+1} = \mathbf{x}_{t+1} + \mu_{t+1}(\mathbf{x}_{t+1} - \mathbf{x}_t).$$
 (2)

Define instead the time-aligned coefficient

$$\mu_t := \frac{\beta_{t-1} - 1}{\beta_t}.\tag{3}$$

• From (2) with index shifted, this gives

$$\mathbf{y}_t = \mathbf{x}_t + \mu_t(\mathbf{x}_t - \mathbf{x}_{t-1}). \tag{4}$$

6. Appendices - GSE

From FISTA to NAG when g = 0 (cont.)

• Introduce the "velocity" v_t:

$$\mathbf{v}_t := \mathbf{x}_t - \mathbf{x}_{t-1}. \tag{5}$$

• Using (4), the look-ahead point is $\mathbf{y}_t = \mathbf{x}_t + \mu_t \mathbf{v}_t$. Plug this into the gradient step (1):

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \mu_t \mathbf{v}_t - \alpha \nabla f(\mathbf{x}_t + \mu_t \mathbf{v}_t). \tag{6}$$

• Now rewrite (6) in velocity form by subtracting x_t from both sides:

$$\mathbf{v}_{t+1} = \mathbf{x}_{t+1} - \mathbf{x}_t = \mu_t \mathbf{v}_t - \alpha \nabla f(\mathbf{x}_t + \mu_t \mathbf{v}_t), \quad \mathbf{x}_{t+1} = \mathbf{x}_t + \mathbf{v}_{t+1}. \tag{7}$$

• Conclusion: Equations (7) are exactly the Nesterov Accelerated Gradient (NAG) updates, where $\mu_t = \frac{\beta_{t-1}-1}{\beta_t}$ provides the momentum parameter.



Figure: Second, third and fifth coordinates of r are set to 0, the other two are perturbed.

Define the selection matrix

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}, \qquad A^{\top} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}. \tag{16}$$

• Perform a QR decomposition of A^{\top} : find orthogonal H and upper–triangular R such that

$$H^{\top}H = I, \qquad HA^{\top} = \begin{vmatrix} R \\ 0 \end{vmatrix}.$$
 (17)

6. Appendices - GSE

Projected NAG example (cont.)

- Since the columns of A^{\top} are already orthonormal up to permutations/signs, one valid choice is obtained by permuting rows; H is not unique.
- Split $H = [YZ]^T$ so that the columns of Y span range(A) and the columns of Z span its orthogonal complement. A concrete valid choice is

$$Y = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \qquad Z = \begin{bmatrix} 0 & 1 \\ 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 0 \end{bmatrix}. \tag{18}$$

• Hence any $r \in \mathbb{R}^5$ can be written as

$$\mathbf{r} = Y \mathbf{r}_y + Z \mathbf{r}_z, \qquad \mathbf{r}_y \in \mathbb{R}^3, \ \mathbf{r}_z \in \mathbb{R}^2.$$
 (19)

Coordinates, permutation, and reduced problem

If

$$\mathbf{r} = \begin{bmatrix} a \\ b \\ c \\ d \\ e \end{bmatrix}, \quad \text{then} \quad \mathbf{r}_{y} = \begin{bmatrix} b \\ c \\ e \end{bmatrix}, \quad \mathbf{r}_{z} = \begin{bmatrix} a \\ d \end{bmatrix}.$$
 (20)

Indeed.

$$Y \mathbf{r}_{y} + Z \mathbf{r}_{z} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} b \\ c \\ e \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} d \\ a \end{bmatrix} = \begin{bmatrix} a \\ b \\ c \\ d \\ e \end{bmatrix} = \mathbf{r}. \tag{21}$$

• Stacking (r_y, r_z) and applying $H^{\top} = [Y \ Z]$ gives a fixed permutation of the entries of r:

$$H^{\top} \begin{bmatrix} \mathbf{r}_{y} \\ \mathbf{r}_{z} \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} b \\ c \\ e \\ a \\ d \end{bmatrix} = \begin{bmatrix} a \\ b \\ c \\ d \\ e \end{bmatrix}. \tag{22}$$

New (reduced) problem: with Z as above,

$$\min_{\boldsymbol{z} \in \mathbb{R}^2} \mathcal{L}(f_{\theta}(\boldsymbol{x} + Z\boldsymbol{z}), t) + \mu \|Z\boldsymbol{z}\|_2.$$
 (23)

• Why $ZZ^{\top} = P_V$, i.e., projection matrix onto ker A?

• The matrix product that reorders coordinates equals P_V , and applying it to a vector (e.g., a gradient) yields

$$P_V \nabla f(\mathbf{r}_t) \tag{25}$$

by the definition of P_V (cf. Eq. (2.14)), which zeros entries outside V.

6. Appendices - GSE

GSE results on targeted adversarial attacks

Table: Targeted attacks performed on ResNet20 classifier for CIFAR-10, and ResNet50 and ViT_B_16 classifiers for ImageNet. Tested on 1k images from each dataset, 9 target labels for CIFAR-10 and 10 target labels for ImageNet.

| | | Best case | | | Average case | | | | Worst case | | | | | | | |
|----------------------|-----------------------------------|-----------------------|-----------------------------|-----------------------------|-----------------------------|------------------------------|-----------------------|-----------------------------|-----------------------------|----------------------|------------------------------|-----------------------|------------------------------|-----------------------------|----------------------------|------------------------------|
| | Attack | ASR | ACP | ANC | ℓ_2 | $d_{2,0}$ | ASR | ACP | ANC | ℓ_2 | $d_{2,0}$ | ASR | ACP | ANC | ℓ_2 | $d_{2,0}$ |
| CIFAR-10 ResNet20 | GSE (Ours) StrAttack FWnucl | 100% 100% 100% | 29.6 78.4 283 | 1.06 4.56 1.18 | 0.68 0.79 1.48 | 137 352 515 | 100% 100% 85.8% | 86.3 231 373 | 1.76 10.1 2.52 | 1.13 1.86 2.54 | 262 534 564 | 100% 100% 40.5% | 162 406 495 | 3.31 15.9 4.27 | 1.57 4.72 3.36 | 399 619 609 |
| ImageNet ResNet50 | GSE (Ours) StrAttack FWnucl | 100% 100% 31.1% | 3516 6579 9897 | 5.89 7.18 3.81 | 2.16 2.45 2.02 | 5967 9620 11295 | 100% 100% 7.34% | 12014 15071 19356 | 14.6 18.0 7.58 | 2.93 3.97 3.17 | 16724 20921 26591 | 100% 100% 0.0% | 21675 26908 N/A | 22.8 32.1 N/A | 3.51 6.13 N/A | 29538 34768 N/A |
| ImageNet ViT_B_16 | GSE (Ours) StrAttack FWnucl | 100% 100% 53.2% | 916 3550 5483 | 3.35 7.85 4.13 | 2.20 2.14 2.77 | 1782 5964 6718 | 100% 100% 11.2% | 2667 8729 6002 | 7.72 17.2 9.73 | 2.87 3.50 3.51 | 4571 13349 7427 | 100% 100% 0.0% | 5920 16047 N/A | 14.3 27.4 N/A | 3.60 5.68 N/A | 9228 22447 N/A |

6. Appendices - GSE

Quantitative evaluation

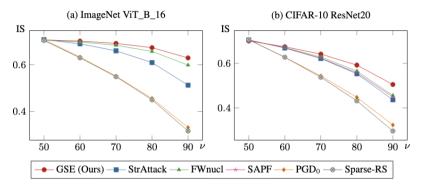


Figure: IS vs. percentile ν for targeted versions of GSE vs. five other attacks. Evaluated on an ImageNet ViT_B_16 classifier (a), and CIFAR-10 ResNet20 classifier (b). Tested on 1k images from each dataset, 9 target labels for CIFAR-10 and 10 target labels for ImageNet.

7. Appendices - S-CFE

Constraining the Sparsity

- Regularize using the indicator function of the sparsity constraint
 - \hookrightarrow Improved control over sparsity

$$I_{\parallel \mathbf{x}' - \mathbf{x} \parallel_0 \le m}(\mathbf{x}') := egin{cases} 0, & ext{if } \parallel \mathbf{x}' - \mathbf{x} \parallel_0 \le m \\ +\infty, & ext{otherwise}. \end{cases}$$

- New $g(\mathbf{x}') := I_{\mathcal{A}}(\mathbf{x}') + \beta I_{\|\mathbf{x}' \mathbf{x}\|_0 \le m}(\mathbf{x}')$ is an indicator function
 - \hookrightarrow Proximal operator coincides with the projection onto the intersection

$$\{\|\mathbf{x}'-\mathbf{x}\|_0\leq m\}\cap\mathcal{A}.$$

7. Appendices - S-CFE

Proximal operator of an indicator function

For any indicator function $I_S(\mathbf{y}) = \begin{cases} 0, & \text{if } \mathbf{y} \in S, \\ +\infty, & \text{if } \mathbf{y} \notin S. \end{cases}$, its proximal operator is the projection onto the set S:

$$\operatorname{prox}_{I_{S}}(\mathbf{x}) = \arg\min_{\mathbf{y}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^{2} + I_{S}(\mathbf{y}) \right\} = \arg\min_{\mathbf{y} \in S} \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^{2} = P_{S}(\mathbf{x}).$$

• Therefore, when $g_p(y)$ is a sum of indicator functions, its proximal operator is the projection onto the intersection of the sets defining those indicators (provided that the intersection is nonempty).

8. Learning from perturbations

Learning from adversarial perturbations

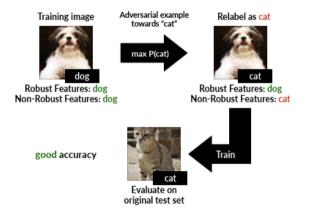


Figure: Training on a dataset which appears mislabeled to humans (via adversarial examples) results in good accuracy on the original test set (Ilyas et al. 2019).