

Even A DNN's mistakes can make sense...

↳ trained to observe patterns.

Metaphor for ML → Clever Hans.

↳ right answers, wrong pattern! We asked the wrong questions...

Do produced probabilities reflect actual frequencies?

Distribution shift: test input might come from diff. distribution v. training inputs.

↳ spurious correlations in training data, selection bias, system feedback (?)

↳ "out-of-distribution" images → confident but wrong assertions.

↳ calibration struggles forOOD & ID images.

Adversarial example: particularly vivid illustration of how learned models may or may not generalize correctly.

→ very special patterns can change model's classification drastically!

→ not easy to defend against (no "bulletproof" method!)

→ can transfer across networks (in many cases) and work in real life.

overfitting hypothesis? Too complex, easy to find inputs that produce crazy outputs.

↳ evidence suggests this is false, since A-E's are not cause of high variance.

linear models hypothesis? NNs ≈ locally linear, extrapolate in somewhat counterintuitive ways when diverging from data.

↳ "Experiment 2": clean "shift" on one side for adversarial direction, therefore (?) mostly linear decision boundary

Human Adversarial Examples? Naïve exemplar: optical illusions.

Adversarial Examples: Features of a model's learning approach. Not bugs.

How to construct an adversarial example?

→ relation: $R(x, x')$

→ attack: $x^* \leftarrow \arg \max_{x'}: R(x, x') \leq \epsilon \quad L_0(x', y)$
 pick image close to x maximize

How to defend?

→ defense: $\theta^* \leftarrow \arg \min_{\theta} \max_{x'}: R(x, x') \leq \epsilon \quad L_0(x', y)$
 robust loss

Fast gradient sign method (FGSM) - simple approx. method for ℓ_2 -norm relation.

$R(x, x') = \|x - x'\|_{\ell_2}$

"first-order" assumption: $L(x', y) \approx L(x, y) + (x' - x)^T \nabla_x L$ ← 1st-order Taylor expansion!

↳ local linear behavior of NN causes this to be effective!

→ attack: $x^* \leftarrow \arg \max_{x'}: \|x - x'\|_{\ell_2} \leq \epsilon \quad (x' - x)^T \nabla_x L$

$x^* = x + \epsilon \text{sign}(\nabla_x L)$

→ can implement Lagrange multiplier: $x^* \leftarrow \arg \max_{x'}: L_0(x', y) - \lambda R(x, x')$

$\delta \leftarrow x' - x$

$\delta \leftarrow \arg \max_{\delta} L_0(x + \delta, y) - \lambda \|\delta\|_{\ell_2}$

can optimize.

Transferability: we can simply use another NN to construct our AE.

↳ zero-shot black-box attacks...

Finite differences gradient estimation: use a moderate number of queries, for each dimension i of x : small num.

→ get $V_i \leftarrow L(x + 10^{-3} e_i, y)$ ← i th canonical vector

→ $\nabla_x L \approx (V - L(x, y)) / 10^{-3}$

Defending: adversarial training.

1. sample minibatch $\{(x_i, y_i)\}$ from \mathcal{D}
2. for each x_i , compute x_i^* (FGSM, etc.)
3. take SGD step: $\theta \leftarrow \theta - \alpha \sum_i \nabla_{\theta} L_0(x_i^*, y_i)$

↳ caveat: decrease overall accuracy on test set ("pay the price").

usually also incorporate original loss gradient

works well against naïve nets.

further "tricks" exist.