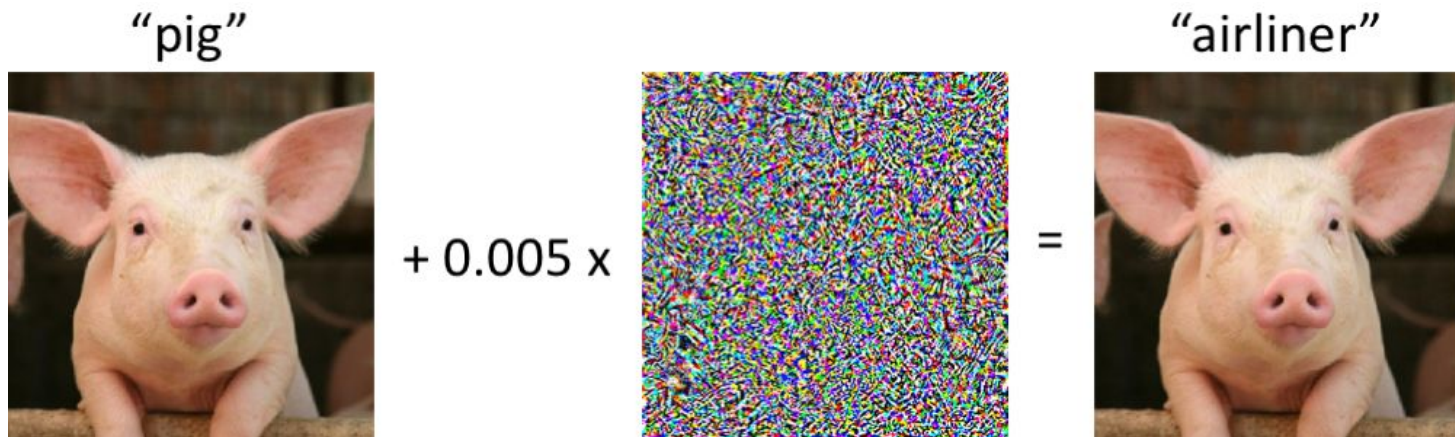


Adversarial Examples Are Not Bugs, They Are Features



Source: [A Brief Introduction to Adversarial Examples – gradient science](#)

План

- Докладчик: вступление, в двух словах о статье
- Исследователь: о соседних результатах и экспериментах

Сеттинг: Adversarial examples

airplane



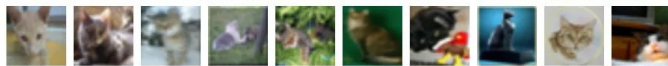
automobile



bird



cat



deer



dog



Adversarial examples, Adversarial attack:

Незаметное возмущение данных,
из-за которого происходит неверная классификация

“A wide variety of models with different architectures trained on different subsets of the training data misclassify the same adversarial example.” - [\[1412.6572\]](#)

Source: [Improving neural networks by preventing co-adaption of feature detectors \(uwaterloo.ca\)](#)

Q: Как возникают Adversarial examples?



“panda”

57.7% confidence

+ 0.007 ×



noise

=



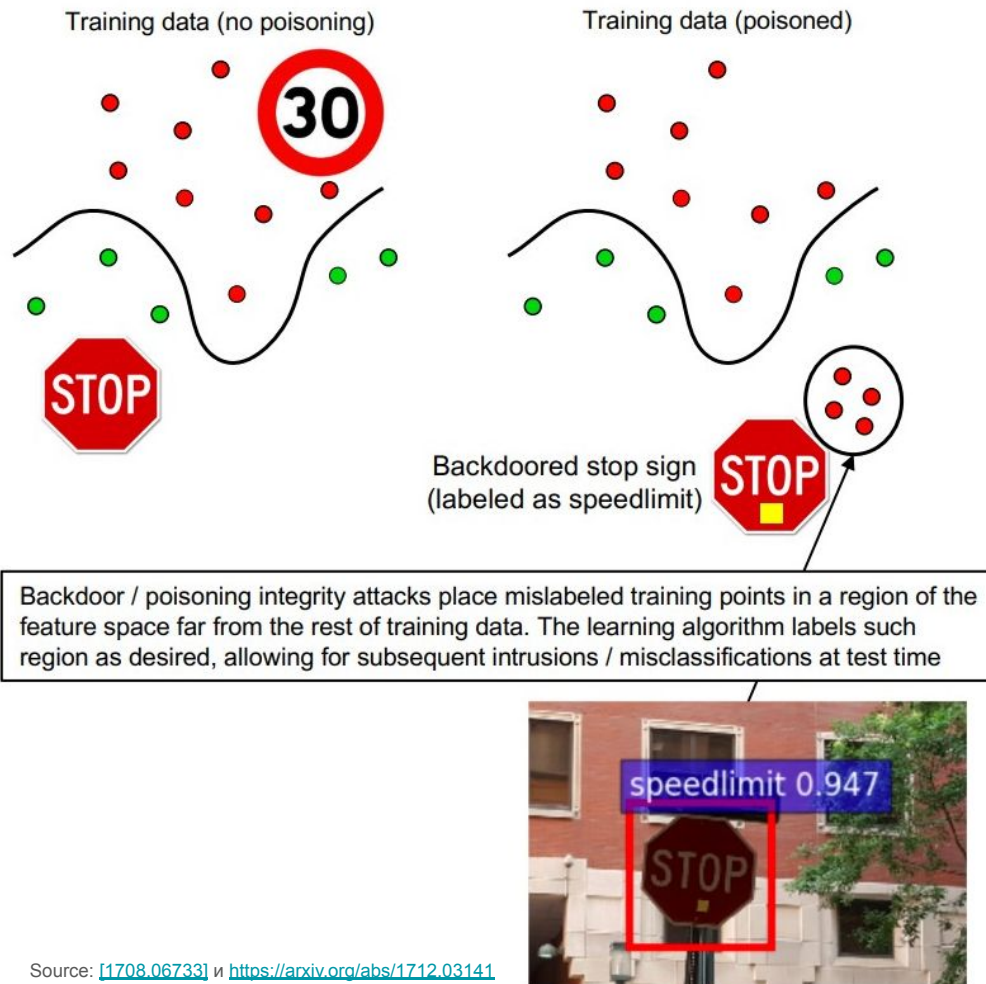
“gibbon”

99.3% confidence

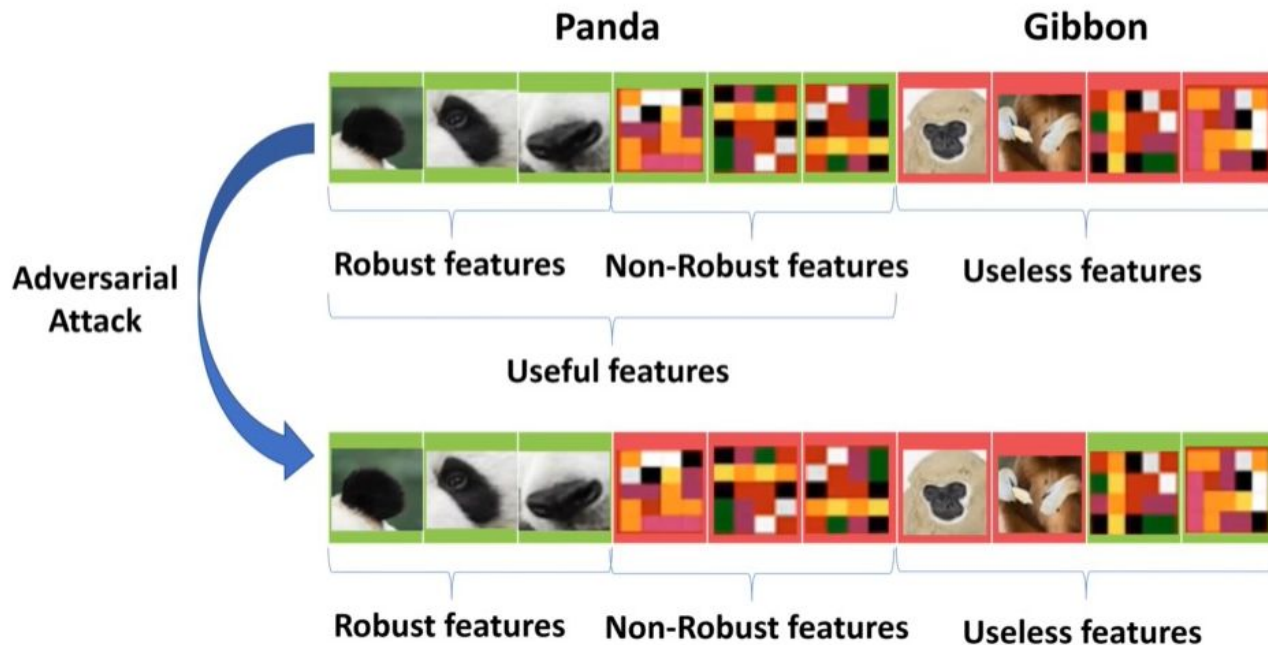
Source: [Deceptive Tricks in Artificial Intelligence: Adversarial Attacks in Ophthalmology](#)

Пример: Data poisoning

- Пре-тренированная модель выпускается в публичный доступ.
- Авторы модели могут активировать уязвимости в модели на основе известных Adversarial examples

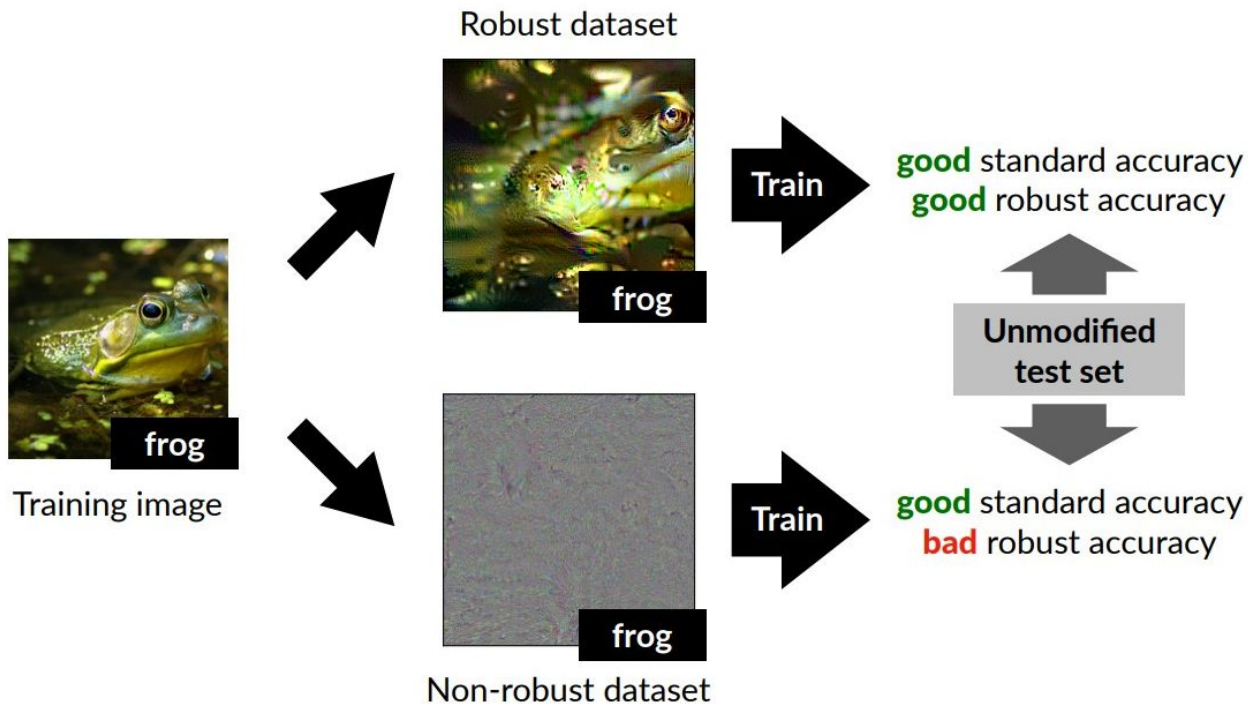


Сеттинг: устойчивые признаки



Proposition: Adversarial examples могут возникать как результат возмущения неустойчивых признаков, не заметных человеческому глазу

Основной тезис(ы)



T1: Уязвимость к Adversarial examples вызвана обучением на неустойчивых признаках

T2: Adversarial examples появляются благодаря хорошей обобщательной способности признаков

Основные определения

Дано:

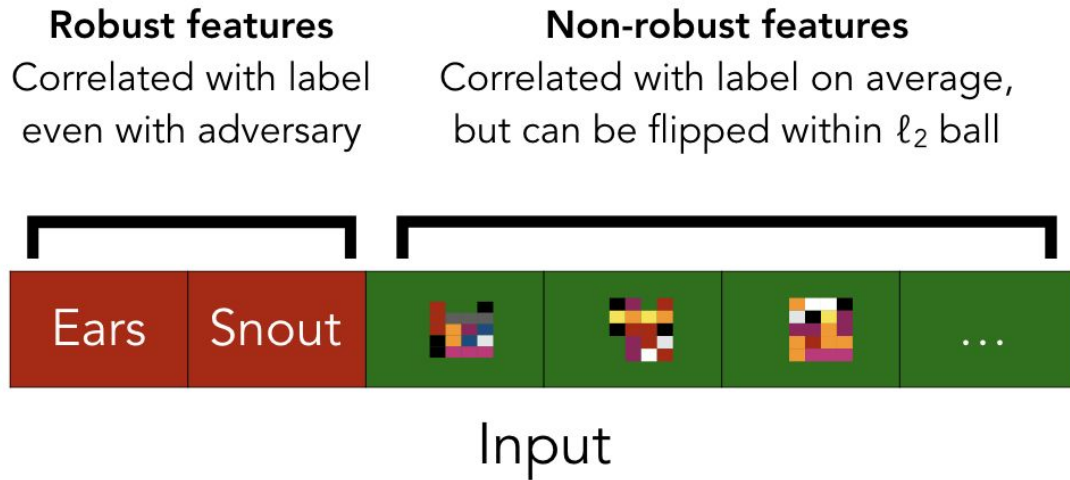
- датасет $(x, y) \sim D$.
- признак $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$

Мы называем признак $f(x)$ **полезным**,
если он скоррелирован с лейблом y :

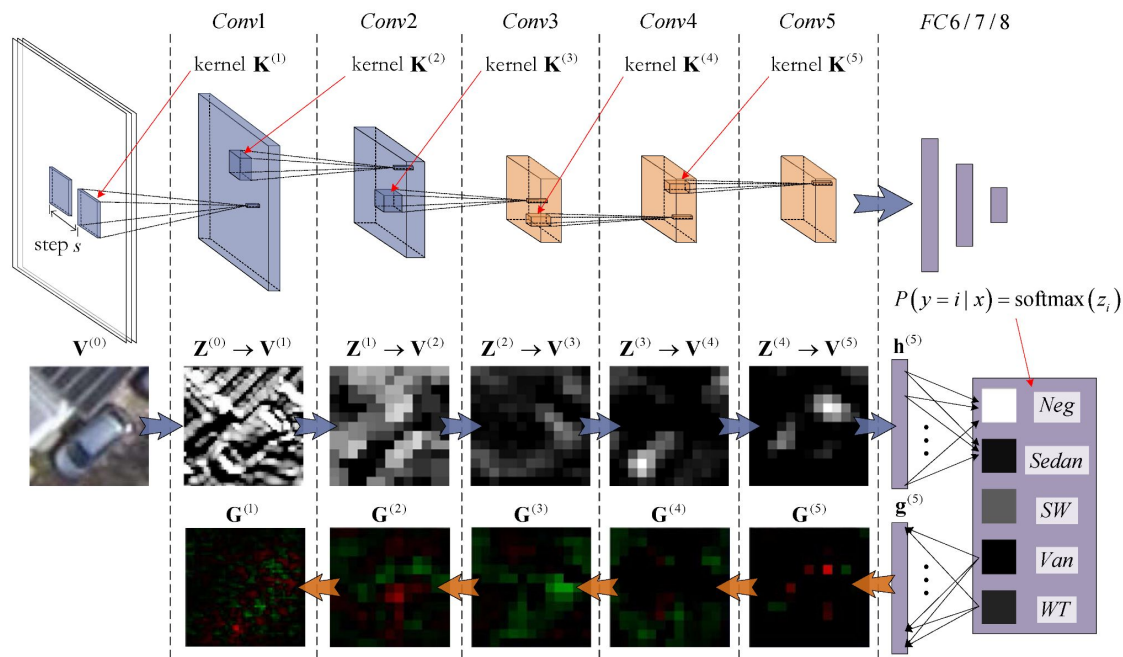
$$\mathbb{E}_{(x,y) \sim D} [y \cdot f(x)] \geq \rho.$$

Мы называем признак $f(x)$ **устойчивым**,
если его полезность устойчива к возмущениям x :

$$\mathbb{E}_{(x,y) \sim D} \left[\inf_{\delta \in \Delta(x)} y \cdot f(x + \delta) \right] \geq \gamma.$$



Построение устойчивого датасета: 1



- Устойчивые признаки условно соответствуют предпоследнему слою нейронной сети, натренированной с помощью **Adversarial training**.

Построение устойчивого датасета: 2

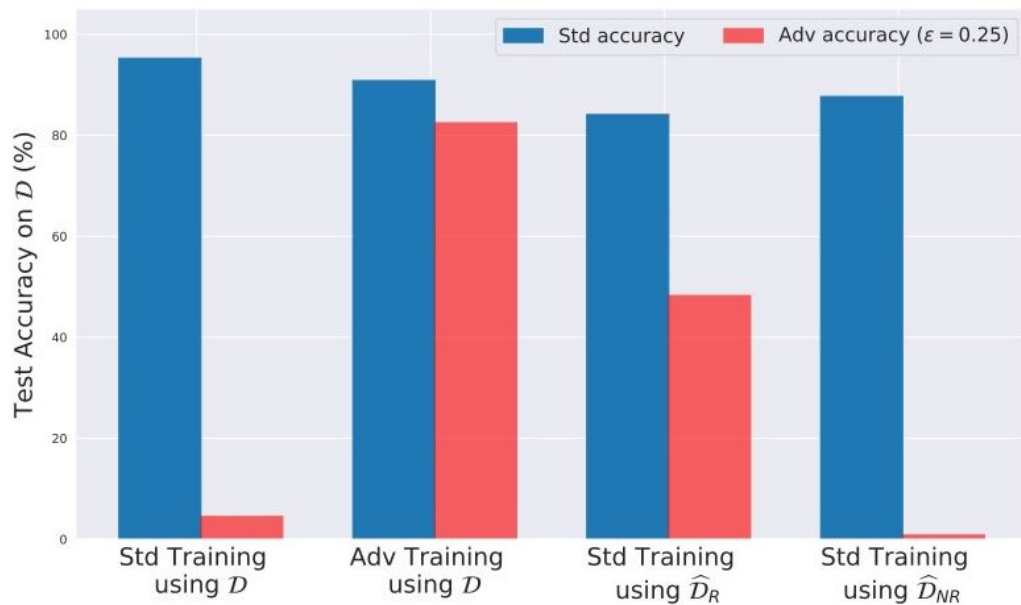
- Для каждого элемента датасета D , строится его “устойчивый аналог”
- Используется **Adversarial training** + projected gradient descent

GETROBUSTDATASET(D)

1. $C_R \leftarrow \text{ADVERSARIALTRAINING}(D)$
 $g_R \leftarrow$ mapping learned by C_R from the input to the representation layer
2. $D_R \leftarrow \{\}$
3. For $(x, y) \in D$
 $x' \sim D$
 $x_R \leftarrow \arg \min_{z \in [0,1]^d} \|g_R(z) - g_R(x)\|_2$ # Solved using ℓ_2 -PGD starting from x'
 $D_R \leftarrow D_R \cup \{(x_R, y)\}$
4. Return D_R



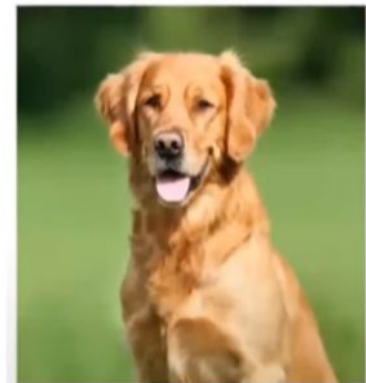
No features no worry



Результаты на CIFAR-10

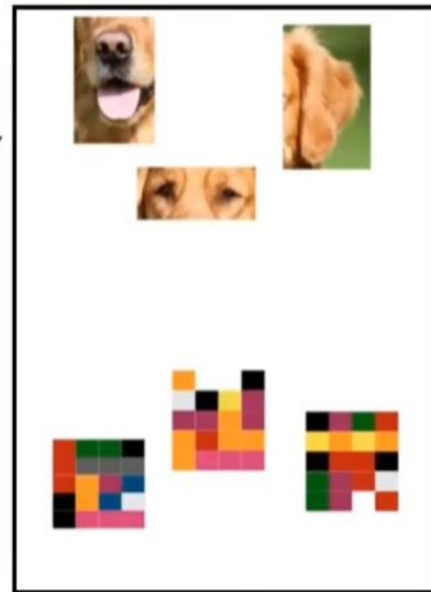
Заключение

- Устойчивость (или неустойчивость) может возникать как свойство данных
- Модели могут находить и использовать устойчивые признаки для повышения точности классификации



dog

features



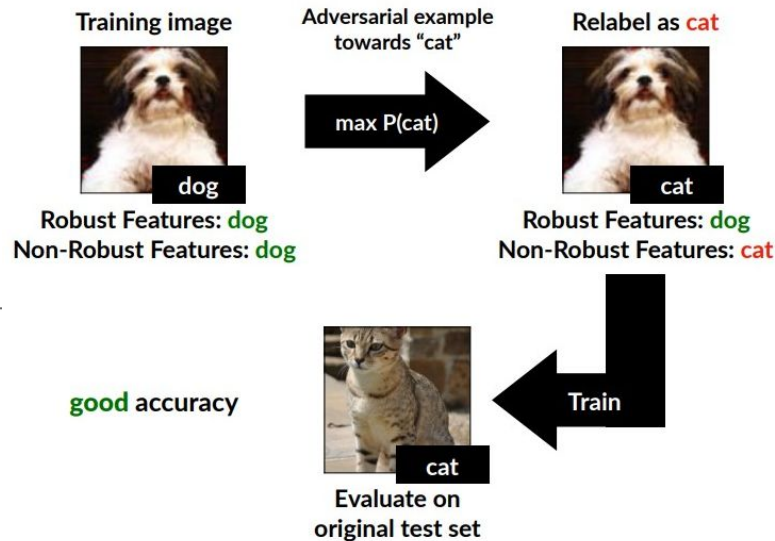
Построение неустойчивого датасета: 1

Idea:

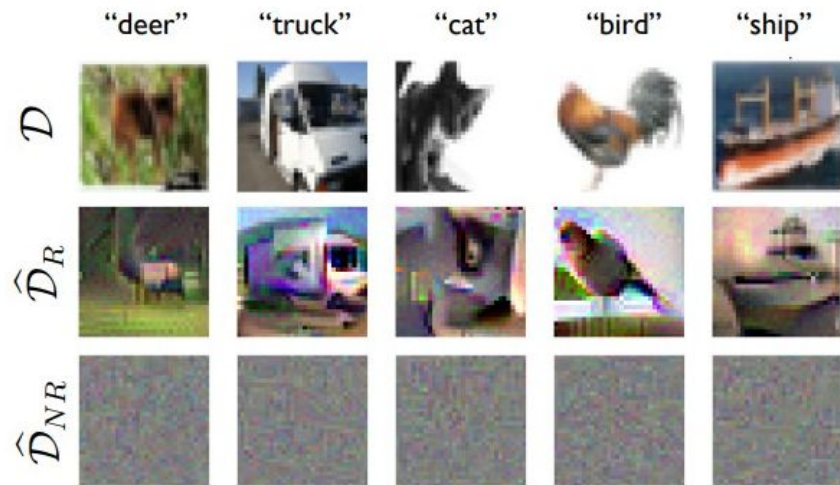
- Избавиться от устойчивых признаков, перемешав датасет и лейблы
- Выучиться на перемешанном распределении

GETNONROBUSTDATASET(D, ϵ)

1. $D_{NR} \leftarrow \{\}$
2. $C \leftarrow \text{STANDARDTRAINING}(D)$
3. For $(x, y) \in D$
 $t \overset{\text{uar}}{\sim} [C]$ # or $t \leftarrow (y + 1) \bmod C$
 $x_{NR} \leftarrow \min_{||x' - x|| \leq \epsilon} L_C(x', t)$ # Solved using ℓ_2 PGD
 $D_{NR} \leftarrow D_{NR} \cup \{(x_{NR}, t)\}$
4. Return D_{NR}



Построение неустойчивого датасета: 2

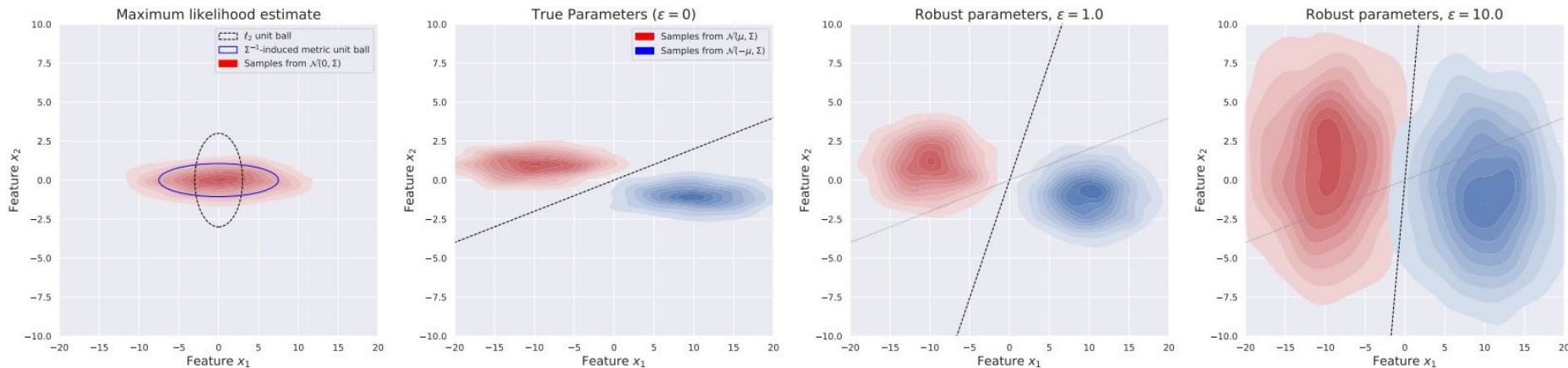


Перемешали лейблы +
+ стартуем со случайного шума

Перемешали лейблы:

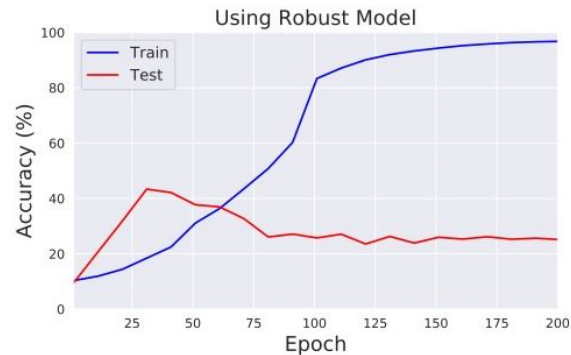
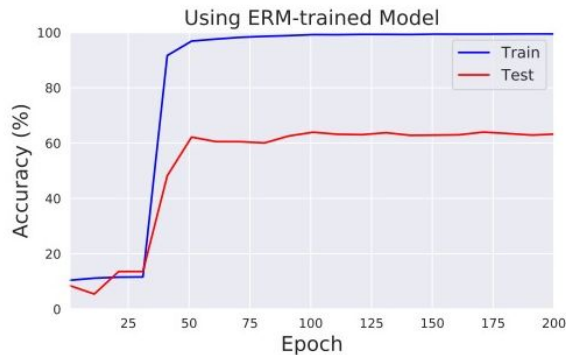


Adversarial training

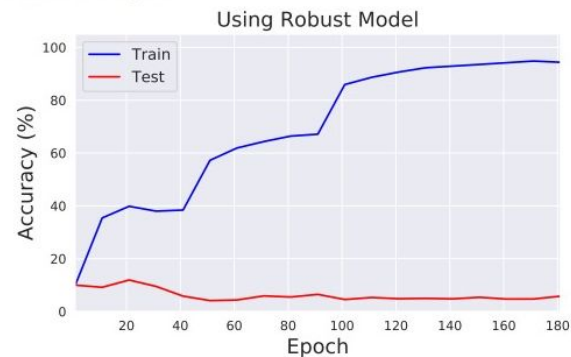
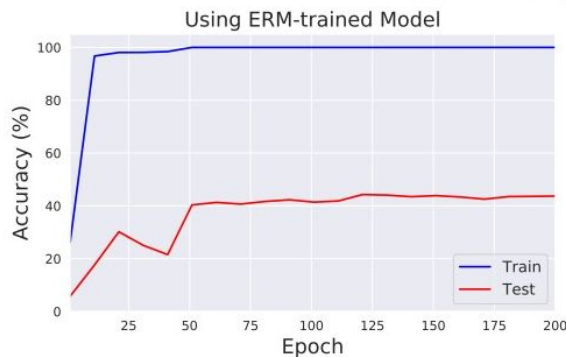


Model	Accuracy	Robust Accuracy	
		$\epsilon = 0.25$	$\epsilon = 0.5$
Standard Training	95.25 %	4.49%	0.0%
Robust Training	90.83%	82.48%	70.90%
Trained on non-robust dataset (constructed from images)	87.68%	0.82%	0.0%
Trained on non-robust dataset (constructed from noise)	45.60%	1.50%	0.0%
Trained on robust dataset (constructed from images)	85.40%	48.20 %	21.85%
Trained on robust dataset (constructed from noise)	84.10%	48.27 %	29.40%

PS: устойчивые признаки могут мешать



(a) Trained using $\hat{\mathcal{D}}_{rand}$ training set



(b) Trained using $\hat{\mathcal{D}}_{det}$ training set