

# Проверка моделей как средство верификации нейронных сетей

Соколов Павел Павлович

# Зачем нужна верификация?<sup>1</sup>

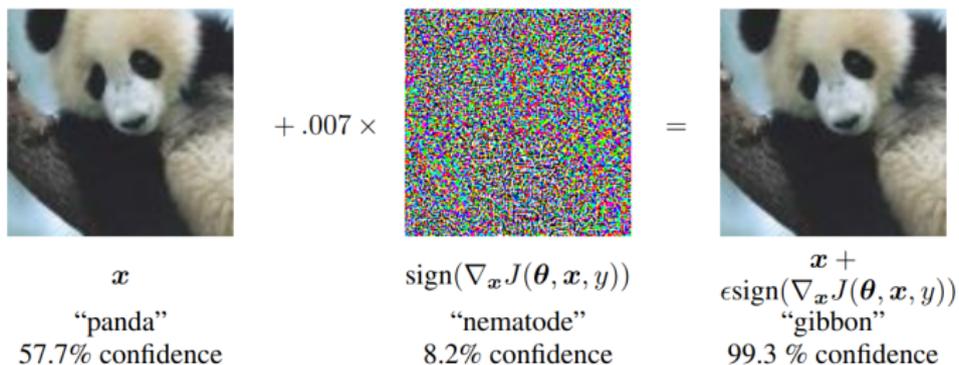
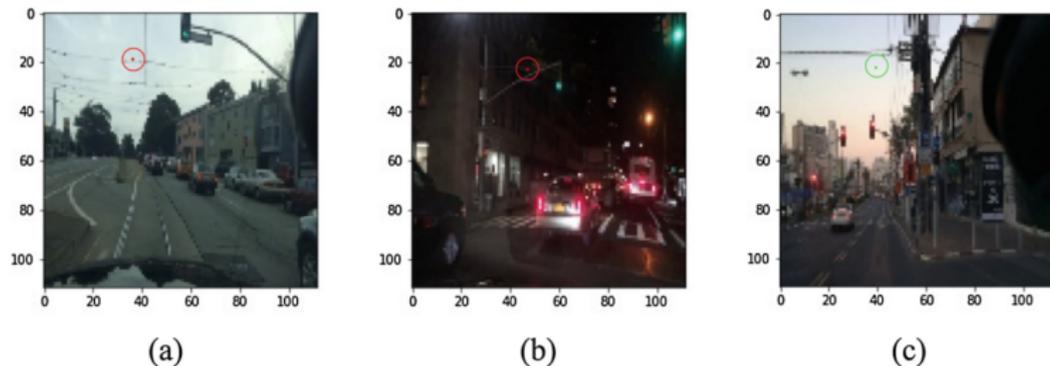


Figure 1: A demonstration of fast adversarial example generation applied to GoogLeNet (Szegedy et al., 2014a) on ImageNet. By adding an imperceptibly small vector whose elements are equal to the sign of the elements of the gradient of the cost function with respect to the input, we can change GoogLeNet’s classification of the image. Here our  $\epsilon$  of .007 corresponds to the magnitude of the smallest bit of an 8 bit image encoding after GoogLeNet’s conversion to real numbers.

---

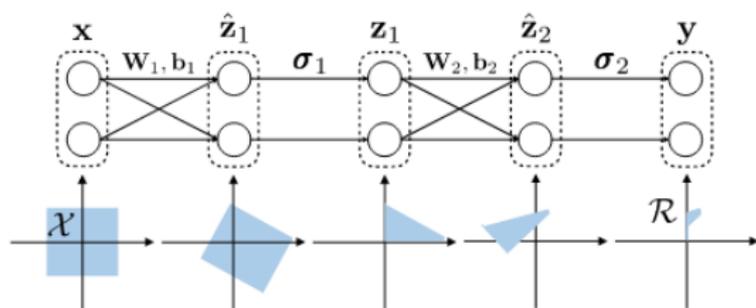
<sup>1</sup>EXPLAINING AND HARNESSING ADVERSARIAL EXAMPLES.  
Ian J. Goodfellow, Jonathon Shlens & Christian Szegedy



**Fig. 4.** Adversarial examples generated on Nexar data demonstrate a lack of robustness. (a) Green light classified as red with confidence 56% after one pixel change. (b) Green light classified as red with confidence 76% after one pixel change. (c) Red light classified as green with 90% confidence after one pixel change. (Color figure online)

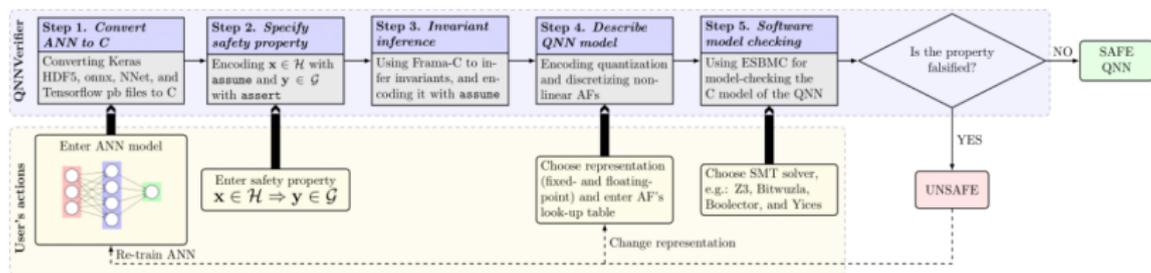
# Задачи верификации нейронных сетей

- ▶ В задаче классификации: робастность — малые изменения входных данных не должны изменять результат.
- ▶ В самоуправляемых устройствах: безопасность — устройство не сталкивается с препятствиями; витальность — устройство доезжает до точки назначения.
- ▶ В языковых моделях — ???
- ▶ ...



**Figure 5.1:** Illustration of reachability methods. The network in the illustration only contains one hidden layer. The input set  $\mathcal{X}$  is first passed through the linear mapping defined by  $\mathbf{W}_1$  and  $\mathbf{b}_1$ . Then it goes through the nonlinear mapping defined by  $\sigma_1$  (ReLU is considered). The corresponding reachable sets are illustrated in the shaded area. The process is repeated for the next layer and the output reachable set is then obtained.

# Робастность классификации — II<sup>4</sup>



<sup>4</sup>QNNVerifier: A Tool for Verifying Neural Networks using SMT-Based Model Checking. Xidan Song, Luiz Sena, Mikel Luján et al.

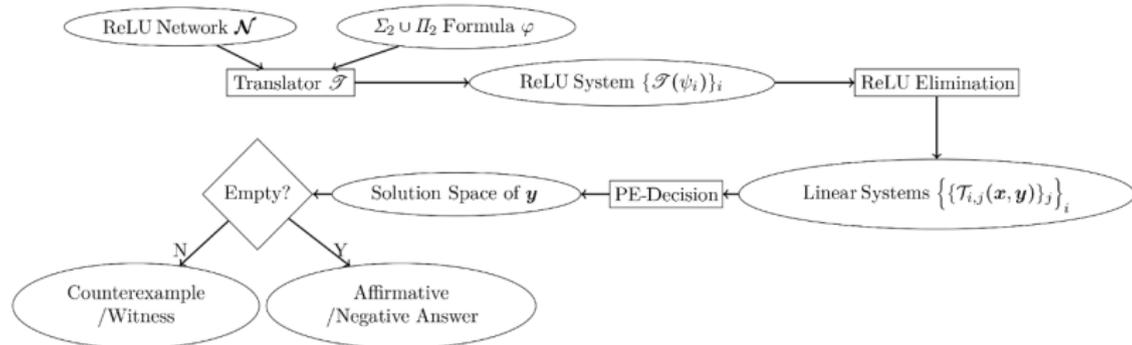


Fig.1. Framework for  $\Sigma_2 \cup \Pi_2$  formula model checking.

## Термы и формулы

$$t ::= v \mid x \mid \bigcirc^k t \mid t[i] \mid Mt \mid t + t$$

$$\phi ::= t_1 \sim t_2 \mid \phi_1 \wedge \phi_2 \mid \neg\phi \mid X\phi \mid \phi_1 U \phi_2 \mid \exists x \in \mathbb{R}^n. \phi \mid \forall x \in \mathbb{R}^n. \phi$$

## Элиминация ReLU

$$\text{ReLU}(t) = E_U t$$

при условии

$$\begin{cases} E_U t \geq 0 \\ (I - E_U)t \leq 0 \end{cases}$$

# Deep statistical model checking<sup>6</sup>

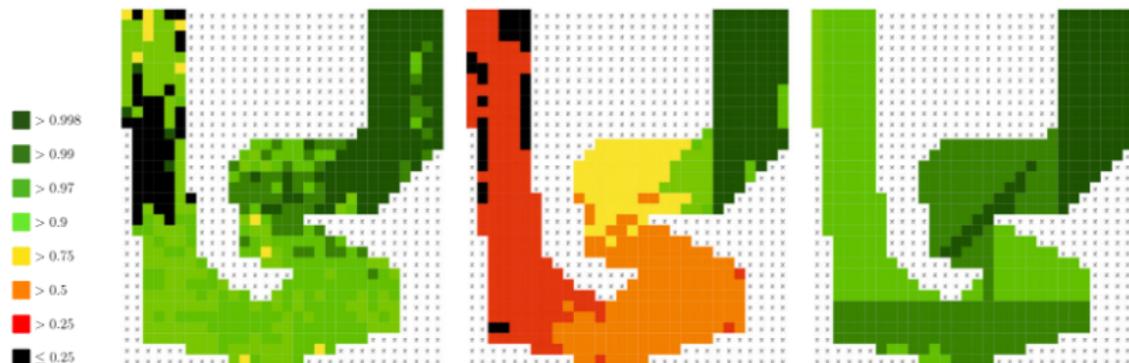


Fig. 4 Goal probability of NN oracle on the Barto-big benchmark trained and executed with 20% noise versus stress-test executed with 50% noise using the same NN (middle) versus optimal policies obtained by probabilistic model checking with 50% noise (right)

---

<sup>6</sup>Analyzing neural network behavior through deep statistical model checking. Timo P. Gros, Holger Hermanns, Jörg Hoffmann et al.

## Определение

Марковский процесс принятия решений — кортеж  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{T}, s_0 \rangle$ , состоящий из конечного множества состояний  $\mathcal{S}$ , конечного множества действий  $\mathcal{A}$ , частичной функции вероятности перехода  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \hookrightarrow \mathcal{D}(\mathcal{S})$  (где  $\mathcal{D}$  — множество распределений вероятности) и начального состояния  $s_0 \in \mathcal{S}$ .

## Определение

(Детерминированная, зависящая от истории) политика действий — функция  $\sigma : \mathcal{S}^+ \rightarrow \mathcal{A}$  такая, что  $\forall w \in \mathcal{S}^*, s \in \mathcal{S} : \sigma(ws) \in \mathcal{A}(s)$ .

(“Беспамятная” политика — такая, что для любых  $w, w'$   $\sigma(ws) = \sigma(w's)$ )

## Определение

Цепь Маркова — кортеж  $C = \langle \mathcal{S}, \mathcal{T}, s_0 \rangle$ , состоящий из множества состояний  $\mathcal{S}$ , функции вероятности перехода  $\mathcal{T} : \mathcal{S} \rightarrow \mathcal{D}(\mathcal{S})$  и начального состояния  $s_0 \in \mathcal{S}$ .

Заметим, что применение политики действий  $\sigma : \mathcal{S}^+ \rightarrow \mathcal{A}$  к Марковскому процессу принятия решений  $\langle \mathcal{S}, \mathcal{A}, \mathcal{T}, s_0 \rangle$  порождает цепь Маркова на множестве состояний  $\mathcal{S}^+$ . (Беспамятная политика позволяет сузить множество состояний до  $\mathcal{S}$ ). Как формулировать и верифицировать свойства такой цепи?

## Синтаксис PCTL

$$\phi ::= \top \mid a \mid \neg\phi \mid \phi \wedge \phi' \mid P_{\geq\theta}(\psi)$$

$$\psi ::= \phi \mid X\phi \mid \phi U^{\leq t}\phi' \mid \phi U\phi'$$

## Техники верификации

На данный момент только численными методами.<sup>7</sup>

---

<sup>7</sup>A Survey of Statistical Model Checking. Gul Agha and Karl Palmskog