

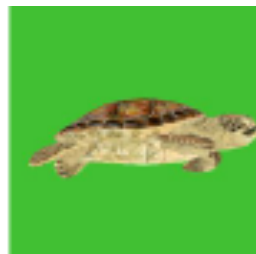
実は簡単!?

AIを攻撃してみよう

KOTOKAZE

この違いは？というお話です

Original:
turtle



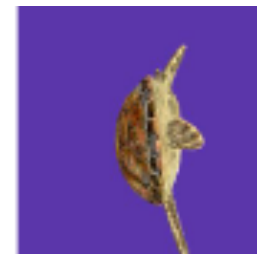
97% /
0%



96% /
0%



96% /
0%



20% /
0%

亀を
ジグソーパズル
と分類している

Adv:
jigsaw
puzzle



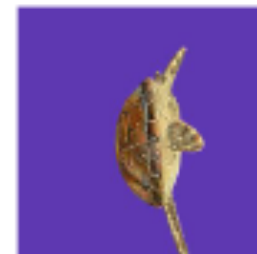
0% /
100%



0% /
99%



0% /
99%



0% /
83%

そもそもAIって？



学習



推論



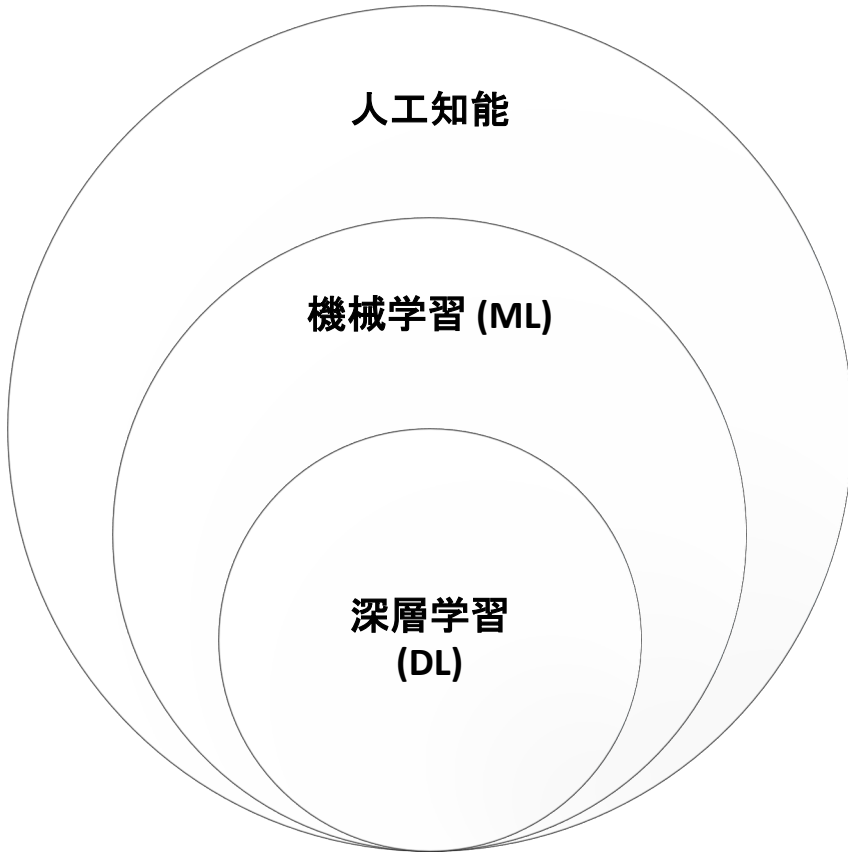
認識



記憶

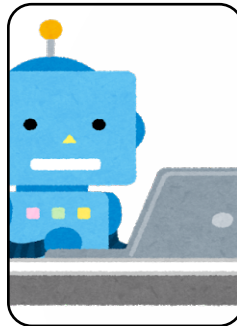
の4つをコンピュータ上で行うもの

AI の分類



ML

- 要設定



DL

- 自動チューニング

機械学習 (ML) の分類

機械学習

教師あり

- 分類
- 回帰

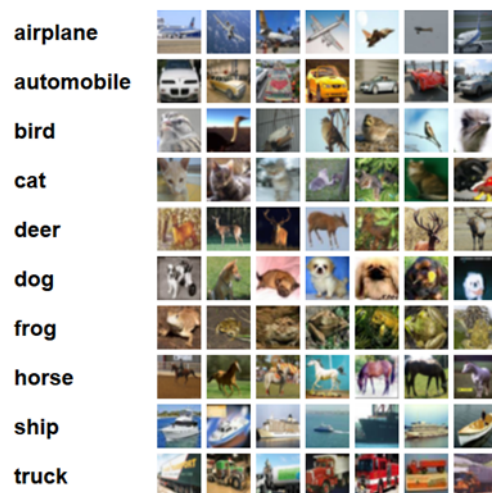
教師なし

強化学習

教師あり学習

分類

学習データ

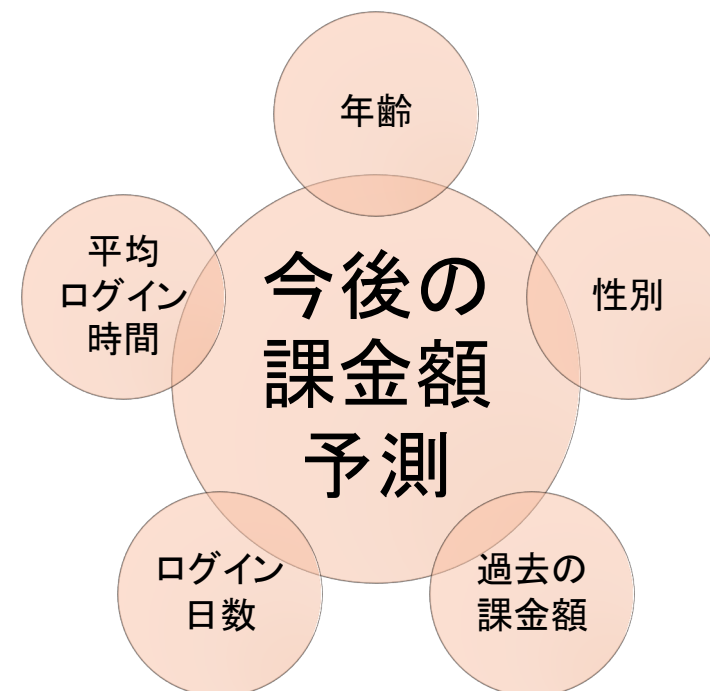


判定したいもの

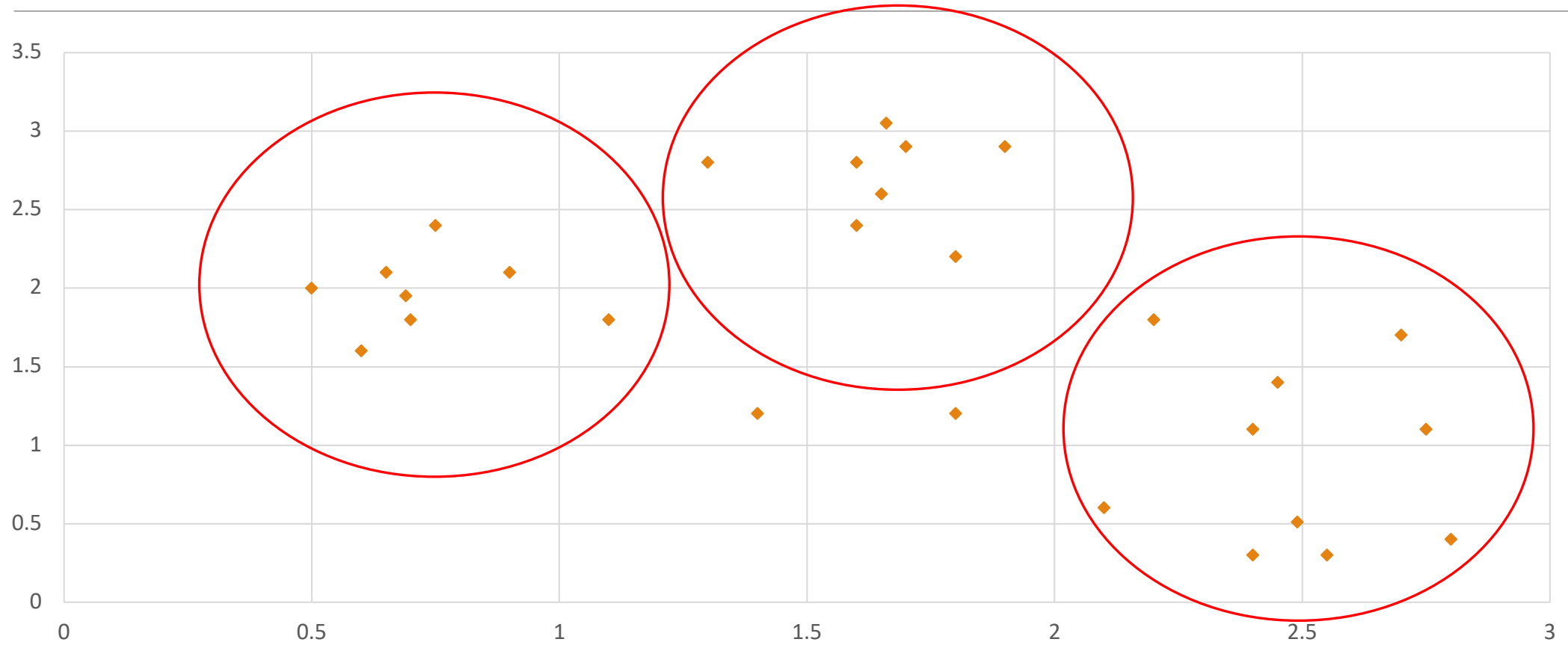


ラベリング済み画像

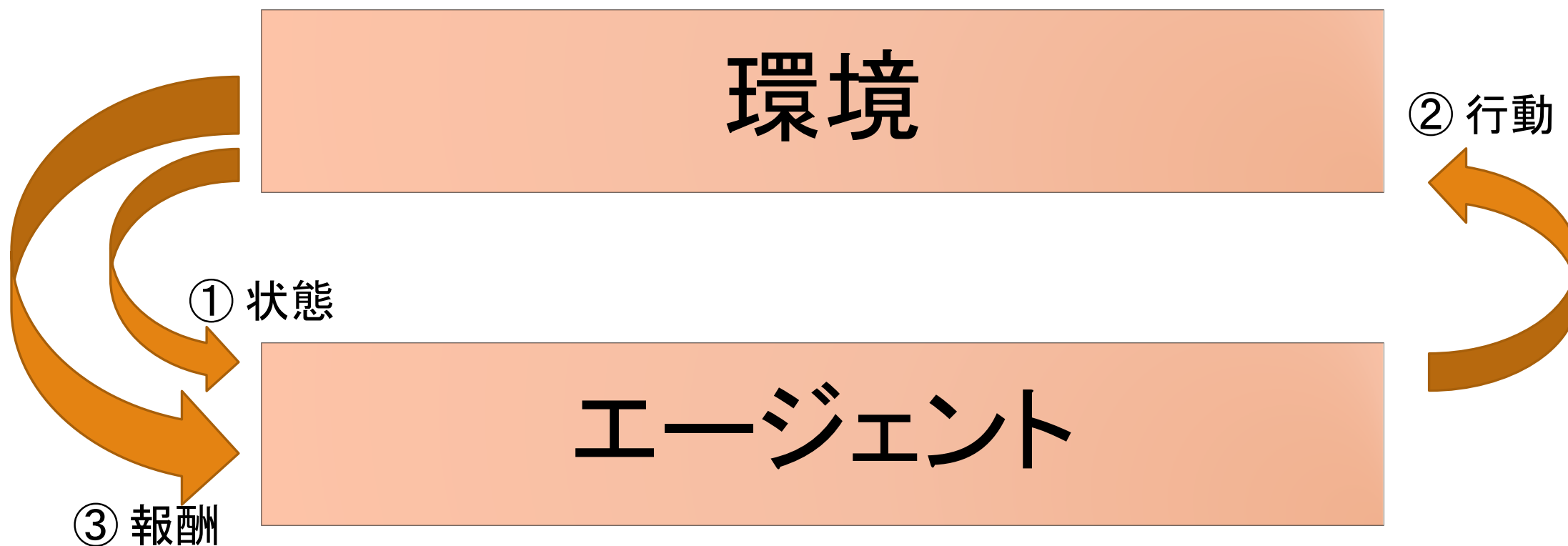
回帰



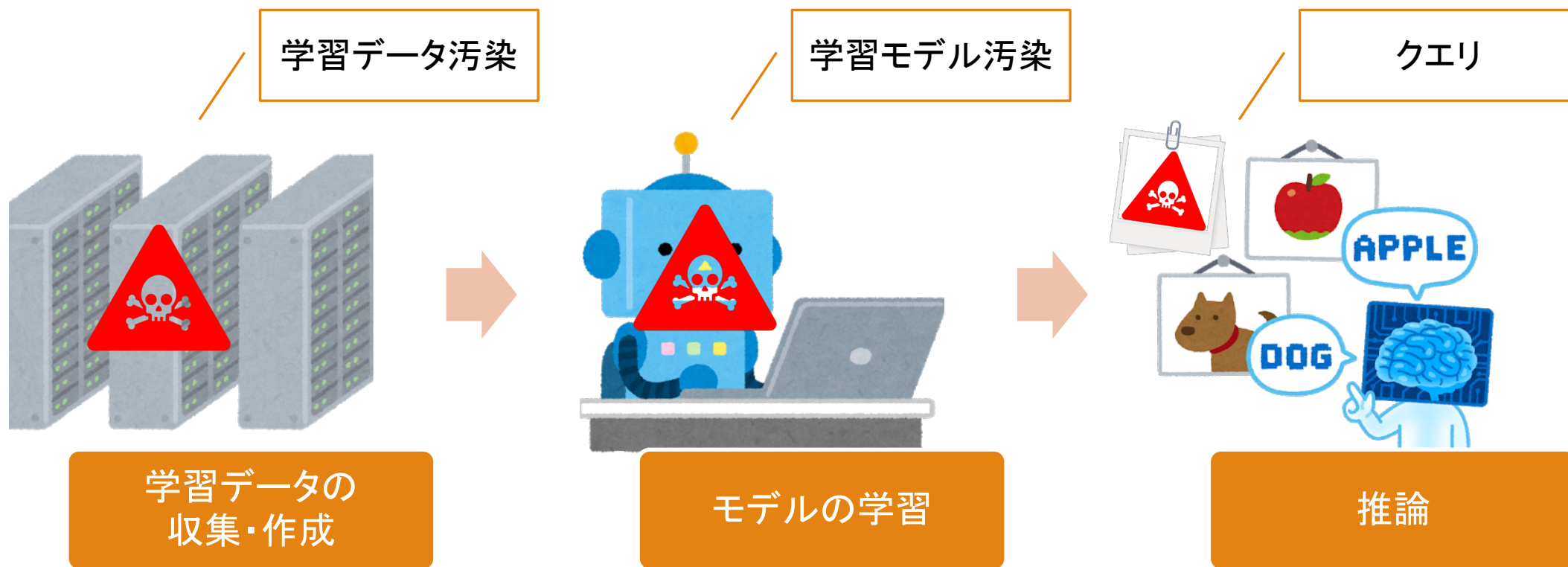
教師なし学習: データの自動分類



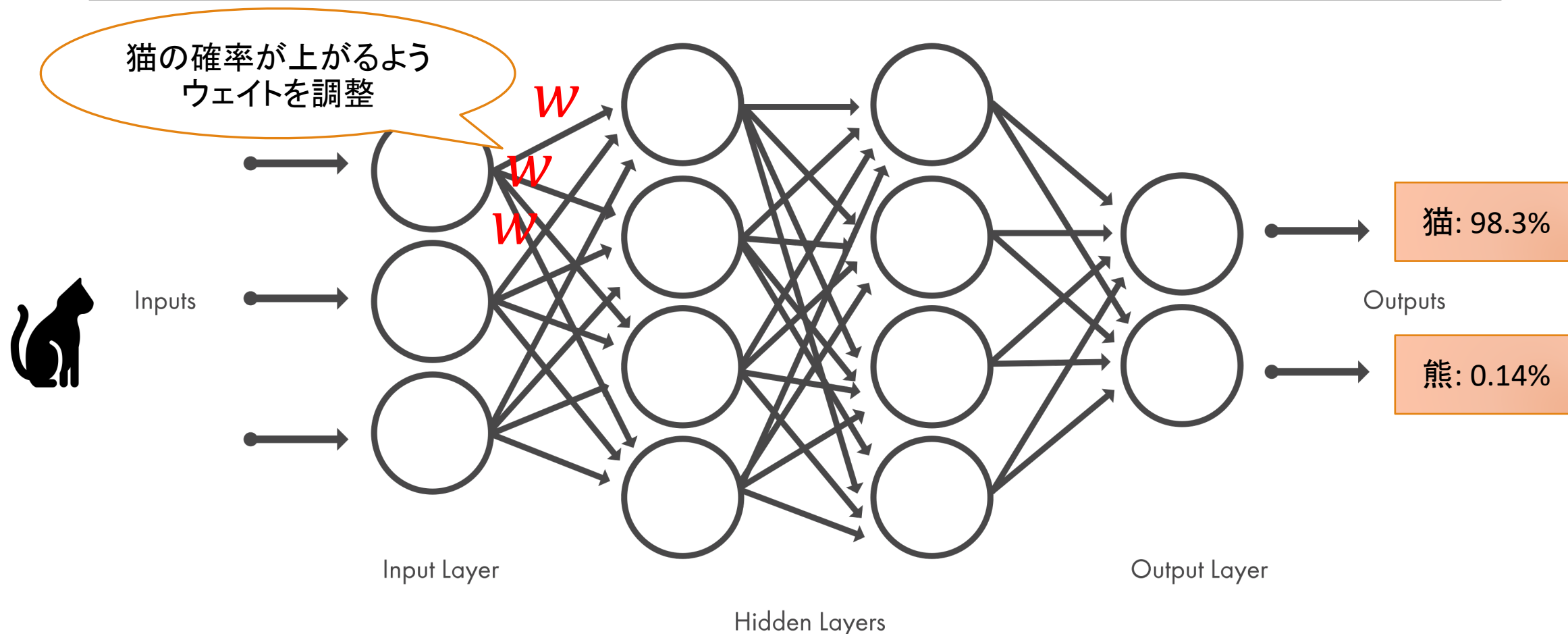
強化学習



モデルへの攻撃の余地



画像分類器 (CNN) の学習の仕組み



敵対的サンプル (Adversarial Examples)

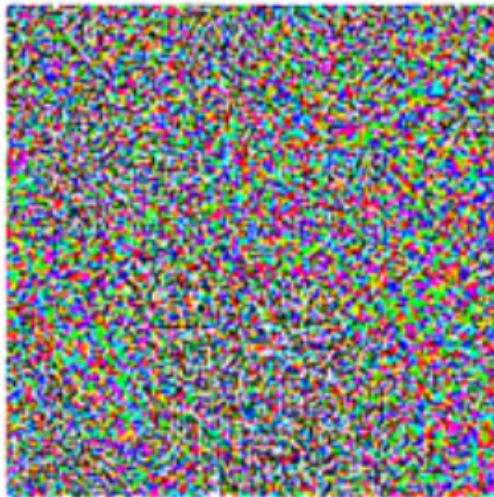


x

“panda”

57.7% confidence

$+ .007 \times$



$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence

$=$



$x +$

$\epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”

99.3 % confidence

AEs 攻撃の仕組み

