

Аналіз атак на лінійні методи машинного навчання

Середович В.В

Науковий керівник: Музичук Ю.А.

ЛЬВІВСЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ ІМЕНІ ІВАНА ФРАНКА

Факультет прикладної математики та інформатики

Кафедра обчислювальної математики

29 травня 2020 р.

Постановка задачі

- Етапи роботи
- Модель

Методи атак

- FGSM
- I-FGSM
- DeepFool

Методи захисту

- Randomization
- Pixel Reflection

Аналіз алгоритмів

- Аналіз атак
- Аналіз захисту

Висновок

Постановка задачі

Змагальний приклад

Нехай існує класифікатор $f(x) : x \rightarrow y$, де $x \in X, y \in Y$, який передбачає значення y для вхідного x . Метою змагального прикладу є знайти такий x^* , який знаходиться в околі x , але хибно визначається класифікатором. Зазвичай максимальний рівень шуму в змагальному прикладі може бути не більше за певну L_p норму $\|x^* - x\|_p < \varepsilon$, де $p = 1, 2, \infty$. В межах даної роботи для визначення рівня пертурбації буде використовуватись саме L_∞ норма.

Етапи роботи

- Реалізувати лінійну модель машинного навчання
- Розглянути різні методи генерування змагальних прикладів
- Застосувати атаки на створену модель та проаналізувати їх ефективність
- Розглянути можливі методи захисту від атак

Модель

В ролі лінійного методу машинного навчання використовувався алгоритм мультикласової логістичної регресії.

Активаци́йна функція softmax

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{k=1}^C e^{z_k}};$$

Вагова функція кросс-ентропії

$$\xi(Y, X) = \frac{1}{m} \sum_{i=1}^m \xi(y^{(i)}, x^{(i)}) = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^C y_j^{(i)} \log(\hat{y}_j^{(i)}) = J(\omega, b)$$

Для того щоб мінімізувати вагову функцію скористаємось алгоритмом градієнтного спуску.

Ідея Fast Gradient Sign Method методу полягає в тому, щоб знайти такий змагальний приклад x^* який максимізує функцію втрати $J(x^*, y)$ до певного L_∞ обмеження. Цього можна досягти один раз використавши зворотне поширення:

$$X^* = X + \varepsilon \cdot \text{sign}(\Delta_x J(x, y)), \quad (1)$$

Для того щоб результат був в $L_\infty \varepsilon$ - околі вхідного зображення, будемо використовувати функція обрізання: $\text{Clip}_{X, \varepsilon}\{X^*\}$ - функція яка виконує по-піксельне обрізання зображення X^* так, X .

$$\text{Clip}_{X, \varepsilon}\{X^*\}(x, y, z) = \min\left\{255, X(x, y, z) + \varepsilon, \max\{0, X(x, y, z) - \varepsilon, X^*(x, y, z)\}\right\} \quad (2)$$

де $X(x, y, z)$ - це значення каналу z зображення X з координатами (x, y) .

I-FGSM

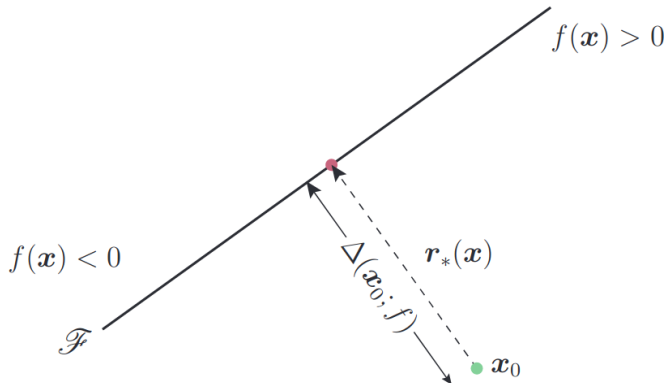
Ми застосовуємо звичайний метод декілька разів з певним невеликим кроком α і після кожного кроку використовуємо функцію обрізання, для того щоб переконатись що значення знаходяться в ε -околі оригінального зображення:

Algorithm 3 $I - FGSM$

- 1: **Input:** Приклад x , значення класу y_{true} , класифікатор f
- 2: **Input:** значення пертурбації ε , кількість ітерації T .
- 3: **Output:** Adversarial x^* з нормою $\|x^* - x\|_\infty \leq \varepsilon$;
- 4: $\alpha = \varepsilon/T$;
- 5: $x^* = x$;
- 6: **for** $t = 0$ **to** $T - 1$ **do**
- 7: $x^* = Clip_{X,\varepsilon}\{x^* + \alpha \cdot sign(\Delta_x J(x, y))\}$;
- 8: **end for**
- 9: **return** $x^* = x_T$.

DeepFool

Для випадку бінарної класифікації легко бачити, що надійність моделі f в точці x_0 , дорівнює відстані від x_0 до площини гіперпараметра $\mathcal{F} = \{x : \omega^T x + b = 0\}$, яка розділяє два класи.



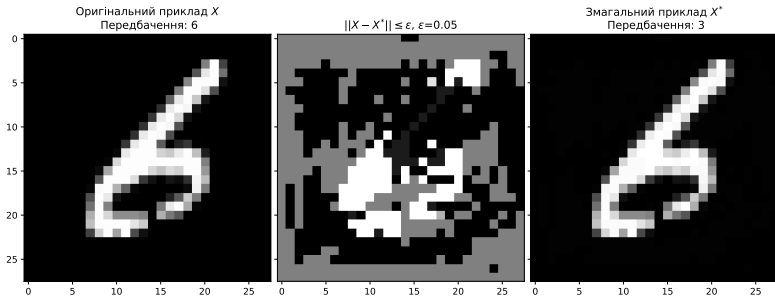


Рис.: Результат роботи алгоритму I-FGSM для $\epsilon = 0.05$

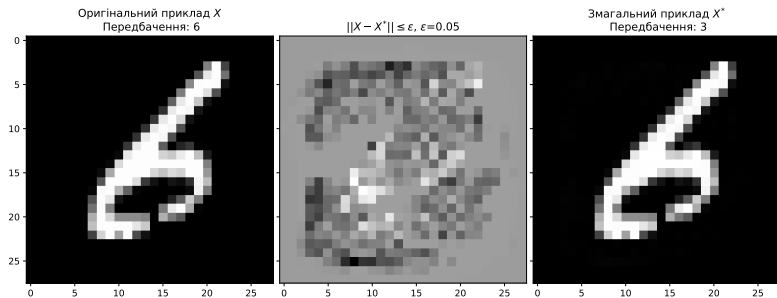
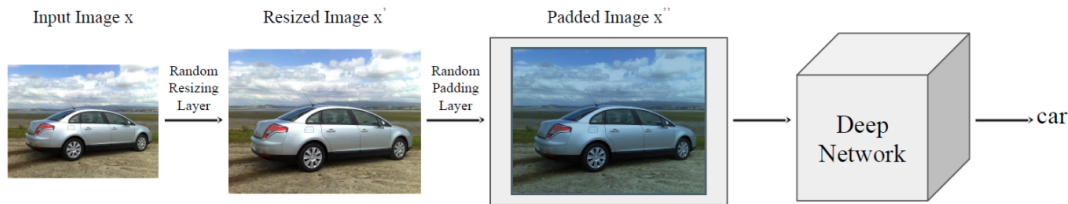


Рис.: Результат роботи алгоритму DeepFool для $\epsilon = 0.05$

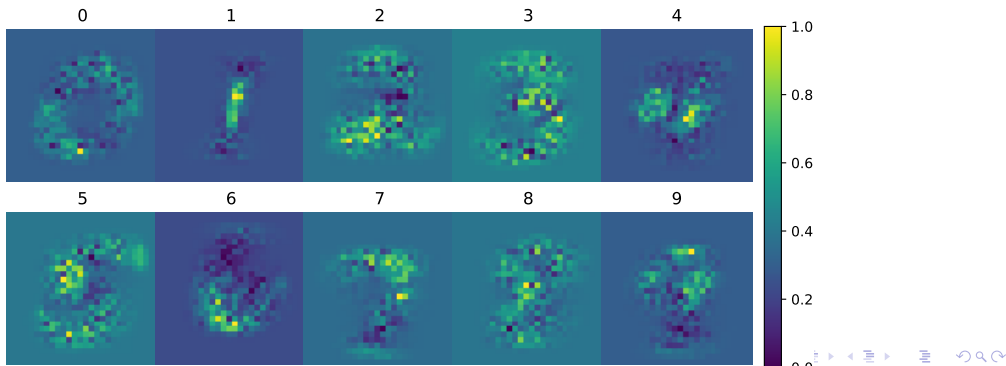
Randomization

- 1 Оригінальне зображення x розмірами $W \times H \times 1$ замінюють на нове x' з випадково вибраними розмірами $W' \times H' \times 1$.
- 2 Другим кроком буде випадкове наповнення деякого простору навколо зображення x' нулями, після чого буде утворене нове зображення x'' з розмірами $W'' \times H'' \times 1$.



Pixel Reflection

Ідея методу полягає в тому, щоб випадково вибрати піксель із зображення і потім, так само випадково замінити його на інший піксель в його малому okolí. Для того щоб покращити цей алгоритм можна визначити найбільш важливі для класифікації класів місця зображення і уникати їх деформації, знижуючи ймовірність внесення туди змін. Для цього можна скористатись активаційною картою класів моделі. Для випадку нашої моделі вони будуть мати вигляд.



Аналіз атак

Пертурбація vs. Успішність атаки

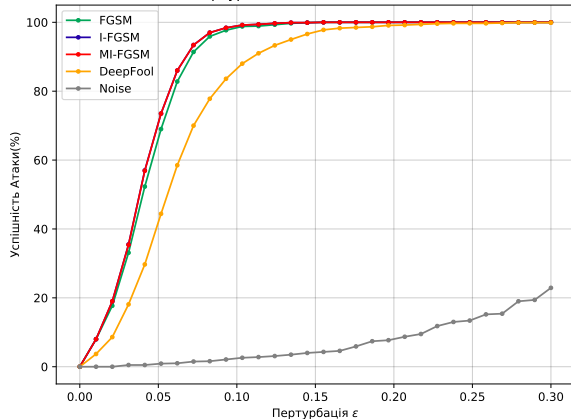


Рис.: Графік залежності успішності атаки від величини пертурбації. Нижньою межею буде виступати випадковий шум.

Пертурбація vs. Точність

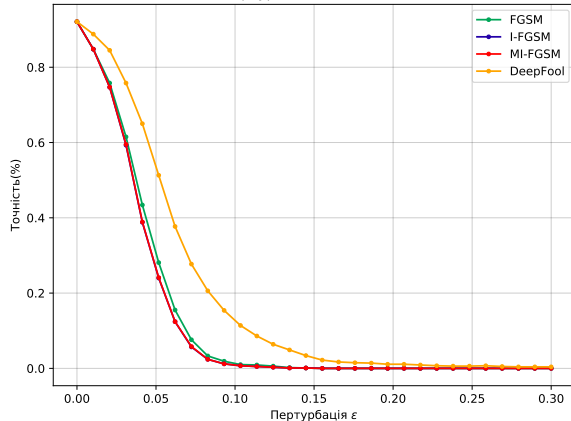


Рис.: Графік залежності успішності класифікації моделі від величини пертурбації. Початкова точність моделі $\sim 92\%$

Randomization

Алгоритм випадкової зміни розмірності опинився неефективним для захисту даної моделі. Навіть при дуже малих деформаціях зображень, класифікатор починав передбачати їх суттєво гірше і точність моделі знижувалась.

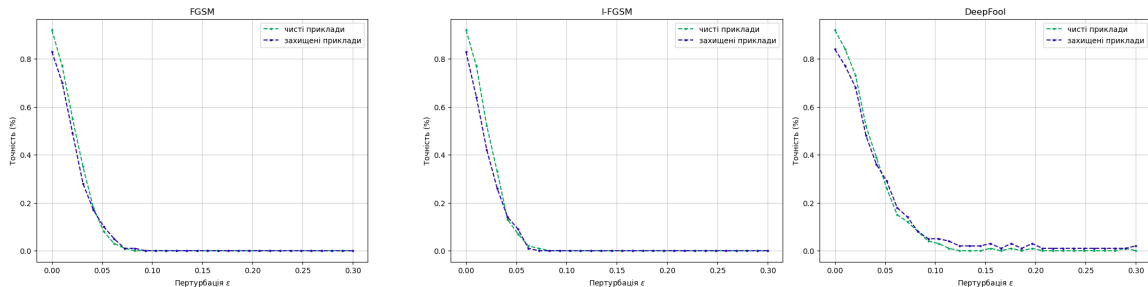


Рис.: Randomization

Pixel Deflection

Навіть при відносно великій кількості зсунутих пікселів модель майже не втрачала точність. При оптимальних параметрах деформації вхідного зображення, точність моделі для чистих зображень не падала зовсім, однак успішність проведених атак знижувалась досить слабо, а отже алгоритм не є надійним методом захисту.

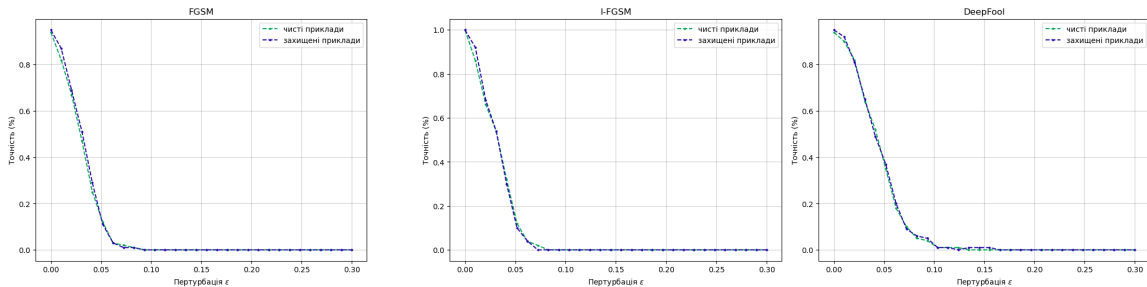


Рис.: Pixel Deflection

Висновок

В межах L_∞ норми найефективнішими методами опинились ітераційна модифікація та momentum оптимізація алгоритму FGSM. Також були розглянуті два методи захисту які базуються на ідеї руйнування структури змагальних прикладів через додавання в зображення контрольованого шуму. Методи захисту які базуються на ідеї випадкових деформацій зображень опинились не дуже ефективними у випадку лінійних моделей машинного навчання.