

Warsztaty 7: Deepfake i oszustwa AI Przykłady fejków i czerwone flagi

Material do cwiczen. Wszystkie sytuacje sa fikcyjne.
Nie zawieraja prawdziwych nazwisk, marek, numerow telefonow ani linkow.

Cel:
Nie uczymy sie tworzyc deepfake'ow ani oszustw.
Uczymy sie zatrzymywac akcje, nazywac sygnały ostrzegawcze i sprawdzac zrodlo.

Rytual przy kazdym przykladzie:
STOP -> zrodlo -> drugi kanal -> kod/link/pieniadze -> decyzja.

JAK PROWADZIC TO CWICZENIE

1. Pokaz albo przeczytaj przyklad.
2. Popros uczestnikow, zeby nie oceniali od razu: prawda czy falsz.
3. Zadaj trzy pytania:
 - Czego ta tresc chce ode mnie?
 - Jaka presje buduje?
 - Co sprawdzam drugim kanalem?
4. Dopiero potem pokaz czerwone flagi.
5. Na koncu wybierzcie pierwszy bezpieczny krok.

PRZYKLAD 1: GLOS BLISKIEJ OSOBY

Co uczestnik dostaje:
Krotka wiadomosc glosowa z nieznanego numeru.

Tresci nie odtwarzamy jako prawdziwego nagrania. Czytamy opis:
"Hej, to ja. Stoje przy kasie, karta mi nie przechodzi. Telefon mi pada.
Wyslij mi kod BLIK, oddam za chwile. Nie oddzwaniaj, bo nie mam czasu."

Na co zwrocic uwage:

- [] glos brzmi znajomo, ale numer jest nieznan
- [] prosba dotyczy kodu BLIK, czyli pieniedzy
- [] pojawia sie zakaz oddzwania
- [] jest presja czasu i zawstydzienie
- [] brakuje potwierdzenia drugim kanalem
- [] glos moze brzmiec plasko, zbyt rowno albo bez naturalnych pauz

Pytania do sali:

- Czy znajomy glos wystarcza do podania kodu?
- Jaki numer wybieramy do oddzwonienia?
- Co mowimy, jesli ktos naciska?

Bezpieczny krok:

Nie podaje kodu. Rozlaczam sie albo nie odpowiadam w tym kanale.
Oddzwaniem na znany wczesniej numer tej osoby.

PRZYKLAD 2: TELEFON Z INSTYTUCJI

Co uczestnik dostaje:
Telefon od osoby podajacej sie za pracownika waznej instytucji.

Opis:
Rozmowca mowi, ze konto albo sprawa urzedowa jest zagrozona. Kaze zostac na lini i,

wejść w przesłany link, zainstalować aplikację i potwierdzić działanie kodem.
Mówi: "Jeśli się rozłączymy, system zamknie sprawę".

Na co zwrócić uwagę:

- [] rozmówca nie pozwala zakończyć rozmowy
- [] link i aplikacja pochodzą od rozmówcy, nie z oficjalnego źródła
- [] prośba dotyczy kodu, danych albo potwierdzenia operacji
- [] autorytet instytucji ma zmniejszyć czujność
- [] nie ma czasu na samodzielne sprawdzenie numeru
- [] głos może brzmieć jak odczytywany skrypt

Pytania do sali:

- Dlaczego "nie rozłączaj się" jest czerwona flaga?
- Skąd bierzemy oficjalny numer?
- Co robimy z linkiem wysłanym podczas rozmowy?

Bezpieczny krok:

Kończę rozmowę. Nie klikam linku i nie instaluję aplikacji.

Samodzielnie wpisuję adres oficjalnej strony albo dzwonię na oficjalny numer.

PRZYKŁAD 3: REKLAMA INWESTYCYJNA Z "EKSPERTEM"

Co uczestnik widzi:

Reklame w mediach społecznościowych z osobą wyglądającą jak ekspert albo znana twarz.

Opis:

Film obiecuje prosty zarobek i "dostęp tylko dziś". Pod filmem jest przycisk prowadzący do formularza. Formularz prosi o numer telefonu, e-mail i kwotę pierwszej wpłaty.

Na co zwrócić uwagę:

- [] obietnica szybkiego i pewnego zysku
- [] znana twarz albo autorytet ma zastąpić dowody
- [] presja: tylko dziś, ostatnie miejsca, tajna metoda
- [] link prowadzi poza oficjalną stronę instytucji
- [] formularz zbiera dane kontaktowe i informacje o pieniądzu
- [] w filmie usta, głos, gesty albo rytm wypowiedzi mogą nie pasować

Pytania do sali:

- Czy znana twarz jest dowodem, że oferta jest prawdziwa?
- Co sprawdzamy na stronie KNF albo oficjalnej stronie instytucji?
- Dlaczego podanie samego telefonu też może być ryzykowne?

Bezpieczny krok:

Nie klikam reklamy. Nie zostawiam danych.

Sprawdzam ostrzeżenia i oficjalne informacje samodzielnie znalezionym kanałem.

PRZYKŁAD 4: SZOKUJĄCE WIDEO DO UDOSTĘPNIENIA

Co uczestnik dostaje:

Film przesłany w komunikatorze z dopiskiem:

"Zobacz, zanim usuną. Przekaz dalej, bo media milczą."

Opis:

Na filmie osoba publiczna mówi coś bardzo ostrego i zaskakującego.

Nagranie nie ma daty, miejsca, linku do pełnego wystąpienia ani kontekstu.

Na co zwrócić uwagę:

- [] silna emocja: gniew, szok, strach albo oburzenie

- [] presja na szybkie udostępnienie
- [] brak daty, miejsca i pełnego nagrania
- [] brak źródła pierwotnego
- [] gesty, mimika, usta i głos mogą nie pasować do siebie
- [] krótki wycinek może być wyrwany z kontekstu nawet bez deepfake'u

Pytania do sali:

- Co jest tutaj silniejsze: dowód czy emocja?
- Gdzie szukamy pełnego nagrania?
- Czy dopisek "nie wiem, czy prawdziwe" wystarczy przed udostępnieniem?

Bezpieczny krok:

Nie przesyłam dalej. Szukam pełnego źródła, daty, miejsca i niezależnego potwierdzenia.

PRZYKŁAD 5: PILNA ZBIORKA Z PODEJRZANYM LINKIEM

Co uczestnik widzi:

Poruszający post z obrazem, krótkim filmem i linkiem do wpłaty.

Opis:

Post mówi, że pomoc jest potrzebna natychmiast. Link prowadzi do strony, która przypomina znaną organizację, ale adres jest inny. Na stronie brakuje pełnych danych organizatora, regulaminu i niezależnego potwierdzenia.

Na co zwrócić uwagę:

- [] bardzo silna emocja i poczucie winy
- [] "wplac teraz" zamiast spokojnej informacji
- [] podobny, ale nieoficjalny adres strony
- [] brak danych organizatora i regulaminu
- [] komentarze mogą sztucznie wzmacniać presję
- [] obraz może być prawdziwy, ale użyty w fałszywym kontekście

Pytania do sali:

- Jak pomagać bez klikania w link z posta?
- Co musi być na wiarygodnej stronie zbiórki?
- Czy udostępnienie bez sprawdzenia może komuś zaszkodzić?

Bezpieczny krok:

Nie klikam linku z posta. Jeśli chcę pomóc, wchodzę samodzielnie na oficjalną stronę organizacji albo sprawdzam zbiórkę w znanym serwisie.

PRZYKŁAD 6: SMS Z LINKIEM I DOPLATA

Co uczestnik dostaje:

SMS z informacją o niedopłacie kwocie, blokadzie paczki albo wygasającej sprawie .

Opis:

"Dopłata 1,49 zł. Brak płatności dziś spowoduje anulowanie dostawy. Opłać tutaj: [podejrzany link usunięty z ćwiczenia]"

Na co zwrócić uwagę:

- [] mała kwota ma obniżyć czujność
- [] link przychodzi w wiadomości, a nie z aplikacji lub oficjalnej strony
- [] jest presja terminu
- [] tekst może mieć dziwny szyk, odmianę albo interpunkcję
- [] prośba prowadzi do płatności albo danych karty
- [] nadawca może udawać kuriera, urząd albo operatora

Pytania do sali:

- Dlaczego mała kwota nadal jest ryzykiem?
- Gdzie sprawdzamy status paczki albo sprawy?
- Co robimy z podejrzanym SMS-em z linkiem?

Bezpieczny krok:

Nie klikam linku. Sprawdzam status w oficjalnej aplikacji albo na stronie wpisanej samodzielnie. Podejrzanym SMS z linkiem można przekazać na numer 8080.

SZYBKA LISTA: NA CO PATRZEC

1. Presja
Teraz, szybko, ostatnia szansa, nie rozłączaj się, nikomu nie mów.
2. Prośba o działanie
Kod, przelew, link, aplikacja, hasło, dane karty, dokument, szybka wpłata.
3. Tożsamość
Kto mówi? Czy mogę to potwierdzić poza tym kanałem?
4. Źródło
Czy jest oficjalny adres, data, miejsce, pełny kontekst i regulamin?
5. Technika
Usta, głos, gesty, tło, rytm, pauzy, akcent, nienaturalna równość wypowiedzi.
6. Logika
Czy treść ma sens? Czy obraz, głos, data, miejsce i intencja pasują do siebie?
7. Emocja
Czy treść chce mnie przestraszyć, zawstydzić, wzruszyć albo rozzłościć, zanim zdążyę sprawdzić?

ZDANIE DO POWTARZANIA

Nie muszę od razu wiedzieć, czy to deepfake.
Mam zatrzymać akcję, nie podawać kodu i sprawdzić źródło drugim kanałem.

ZRODŁA DO OMOWIENIA Z GRUPĄ

- CERT Polska: zgłaszanie incydentów: <https://incydent.cert.pl/>
- CERT Polska: podejrzanym SMS z linkiem na numer 8080: <https://cert.pl/kontakt/>
- NASK: rozpoznawanie deepfake po błędach logicznych, głosowych i językowych
- Policja: oszustwa BLIK i wykorzystanie AI do fałszywych głosów lub wideo
- CEBRF KNF: fałszywe inwestycje i reklamy wykorzystujące deepfake