

Zadanie na doma

7. mája 2018

1 HTTP klient

Napište v Pythone skript `http_get`, ktorý očakáva ako jediný povinný argument URL a na štandardný výstup vráti web stránku (alebo obrázok alebo niečo iné) ktoré je na danom URL. Nesmiete používať Pythonovské moduly, ktoré implementujú HTTP (t.j. `requests`, `urllib2` a iné). Máte to spraviť cez normálne TCP sockety.

1.1 Bližšia špecifikácia

- Obsah máte dávať na výstup ak je status 200 OK, inak nie.
- Na presmerovacie statusy (301,302,303,307,308) reagujte tak, že pôjdete na to URL, ktoré vám bolo poslané v `Location` headri.
- Musíte implementovať aj `Content-length` aj `Transfer-encoding: chunked`.
- Statusy, ktoré nie sú implementované vypíšte na `sys.stderr` a dajte `sys.exit(1)`.
- Nemáte implementovať `https` URL, iba `http`. Také URL treba vyhodnotiť ako chybu a dať `sys.exit(1)`, a to aj v prípade, že ste boli na také URL presmerovaní serverom.
- Nemáte implementovať iné porty ako 80.

1.2 Na čo si treba dať pozor

- HTTP je zmiešaný textovo-binárny protokol – headre sú text (ASCII), obsah je binárny. Z toho vyplýva, že máte použiť `f=socket.makefile("rwb")`, pre čítanie a zapisovanie v binárnom móde.
- Dáta ktoré budete potom čítať zo socketu prostredníctvom `f` budú typu `bytes` nie `str`. Napriek tomu môžete a máte pri čítaní HTTP statusu a headrov používať `f.readline`. Ja som riadky ihneď prevádzal na typ `str` cez `bytes.decode('ASCII')`, aby som mal pre interné spracovanie už `str`.

- Po ukončení headrov už nesmiete predpokladať, že to čo prúdi zo socketu je text v akomkoľvek kódovaní.
- POZOR: pre vypisovanie obsahu (typ `bytes`) je nutné použiť `sys.stdout.buffer`, ten slúži ako štandardný výstup pre *binárne* dáta. Typ `bytes` vôbec na `sys.stdout` nejde zapisovať, tam môžete zapisovať iba typ `str`.

2 Ako som to robil ja

- Z URL som vytiahol *hostname* a *path* pomocou `re.match`.
- Pripojil som sa a poslal `GET` request, včítane povinného `Host:` headra.
- Prečítal som prvý riadok odpovede, dekodoval ako ASCII a potom pomocou `re.match` vytiahol status a jeho popis.
- Prečítal som všetky headre, dekodoval ako ASCII a uložil si ich do slovníka. POZOR veľkosť písomí nie je v HTTP protokole určená, môžete dostať povedzme `CoNTent-LENGtH` a je to legálne. To som vyriešil tak, že som na každý názov headra dával pred jeho uloženie do slovníka `str.lower()`, čím sa všetky písmená hodia na malé.
- Ak bol status presmerovací, vytiahol som URL z *location* a pokračoval tam (mám to ako nekonečný cyklus, vyskakuje sa z neho pri statuse 200).
- Prečítal som obsah odpovede (dve rôzne vetvy: bolo `content-length` alebo máme `transfer-encoding: chunked`) a ak bol status 200, zapísal som obsah na `sys.stdout.buffer`. Pri presmerovacích statusoch treba obsah ignorovať a zavrieť socket.