

3.1 & 3.2

First of all let's start by simulating the process of filling our hard drives. As said in the exercise text in a realistic setting we would use a hash function, but here it suffices to assign files randomly with uniform probability to any hard drive; this is true because an "asymmetric" hash (i.e. one that preferentially fills e.g. HD 0) would never be actually utilized, hence asymptotically the uniform distribution will behave similarly to a realistic hash.

Let's say we're working with 10 1TB HDDs (as said in the exercise text) and with files of `filesize` GB; then our "Montecarlo" algorithm to simulate the storage process works as follows:

- we create an empty length-10 vector x , so that its element represents the amount of stored data in the drives at any given time;
- we initiate a loop where for as long as no hard drive has reached 1 TB (1000 GB) we sample uniformly an integer i in $[0, 9]$, then increment $x[i]$ by `filesize`;
- when one drive i.e. one x component reaches 1000 GB we stop.

Once we have this vector we can use a barplot to represent the insides of our simulated drives.

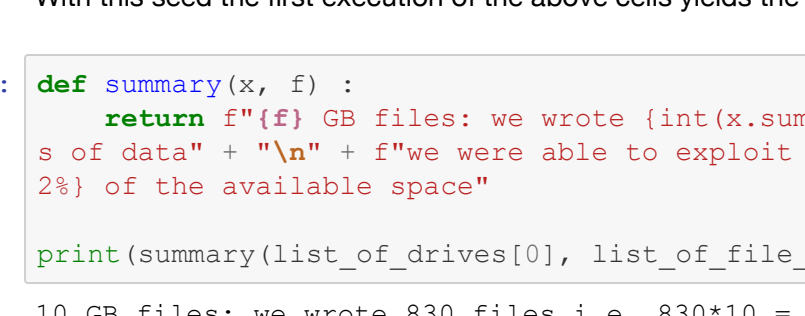
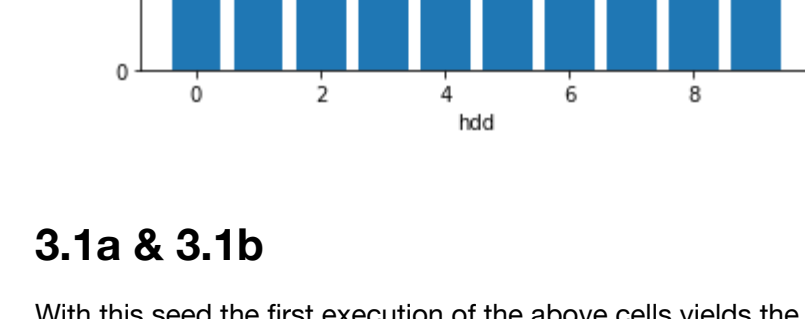
```
In [6]: #init
def populate_drives(filesize):
    x = np.zeros(10)
    while (x < 1000).all(): # I TB = 1000 GB
        i = np.random.randint(10) += filesize
        return x

def populate_and_plot_drives(filesize, verbose = True):
    x = populate_drives(filesize)
    if verbose:
        print(f"x in GB (filesize = {filesize}) GB: ", x)
    fig, ax = plt.subplots()
    ax.bar(np.arange(10), x)
    # we add to the plot 2 informations: how many files were successfully written + how much storage out
    # of the 10 TBs we were able to use
    ax.set_title(f"({int(x.sum()/filesize)} files ({filesize} GB each), {x.sum()/1e4:.2%} occupied)") #
    x_sum() = tot. n. of GB saved, x.sum()/1e4 = occupied fraction of the 10 available TBs
    ax.set_xlabel("nB")
    ax.set_ylabel("GB")
    return x

list_of_file_sizes = [10, 1, 4e-3] # the assignment asks to make plots for filesize = 10 GB, 1 GB, 4 MB
list_of_drives = []

for i in list_of_file_sizes:
    list_of_drives.append(populate_and_plot_drives(i))

x in GB (filesize = 10 GB): [ 790.  920.  800.  860.  770. 1000.  890.  860.  670.  740.]
x in GB (filesize = 1 GB): [ 9345. 9345. 9345. 9345. 9345. 9345. 9345. 9345. 9345. 9345.]
x in GB (filesize = 0.004 GB): [ 996.804 994.82 995.46 998.64 995.52 1000. 997.932 997.0
997.636 997.864]
```



3.1a & 3.1b

With this seed the first execution of the above cells yields the following results.

```
In [7]: def summary(x, f):
        return f"({f} GB files: we wrote {int(x.sum()/f)} files i.e. {int(x.sum()/f)*f} = {x.sum()/1e4:.2%} GB
of data" + "\n" + f"we were able to exploit {x.sum()/1e4}/10000 = {x.sum()/1e4:.4f} = {x.sum()/1e4:.2%}
of the available space"

print(summary(list_of_drives[0], list_of_file_sizes[0]))

10 GB files: we wrote 830 files i.e. 830*10 = 8300.00 GBs of data
we were able to exploit 8300/10000 = 0.8300 = 83.00% of the available space
```

If the simulation cell is executed again multiple times (without resetting the seed) one usually finds values between 80% and 92% (mostly centered around ~85%); this relatively large spread is analyzed below.

3.2a & 3.2b

With this seed the first execution of the simulation cell yields the following results.

```
In [8]: print(summary(list_of_drives[1], list_of_file_sizes[1]))

1 GB files: we wrote 9345 files i.e. 9345*1 = 9345.00 GBs of data
we were able to exploit 9345/10000 = 0.9345 = 93.45% of the available space
```

This time the results are "tighter", in the sense that if the cell is re-executed without a seed reset one typically finds values between 92% and 97% (mostly centered around ~95-96%), and gives us the first hint that smaller file size equals more efficient exploitation of space. *Why does this happen?*

3.3

With this seed the first execution of the simulation cell yields the following results.

```
In [10]: print(summary(list_of_drives[2], list_of_file_sizes[2]))

0.004 GB files: we wrote 2493931 files i.e. 2493931*0.004 = 9975.72 GBs of data
we were able to exploit 9976/10000 = 0.9976 = 99.76% of the available space
```

With a 4 MB file size it's basically impossible to find numbers outside the range 99.5-99.8% by manually re-executing the code cell. *Why is this the case?*

We also notice that the smaller the file size the flatter the obtained distribution. Once again we ask: *but why?*

Increasing efficiency explanation

We can explain the above results using a simple semi-quantitative argument based on the law of large numbers as follows.

Let's temporarily imagine that: the number of files we write/the number of iterations the algorithm runs, is not a random variable, but a predetermined constant (for example fixed in advance by us). If this is the case then the number of files we are able to store inside drive i obeys a *binomial* distribution with $p = 0.1$; this follows from the fact that at each algorithm iteration every hard drive has a 10% chance of being "chosen" by our uniform generator; therefore every iteration is equivalent to one Bernoulli trial per hard disk (hence the binomial, by "accumulation" of Bernoulli trials). This fact is useful because it lets us give an order of magnitude estimate of how large n needs to be to actually fill almost all of the available space.

Indeed let's say that we chose n such that on average each hard drive is almost full (in order to reproduce the 80-90%+ occupied storage values we observed empirically); this implies that if n_i is the number of files in drive i we then have

$$E[X_i] = E[n_i f] \approx 1000 \text{ GB}$$

where we defined X_i as the number of occupied GBs in drive i and f as the constant size of a single file.

Since

$$n_i \sim \text{Binom}(n, p = 0.1) \quad \forall i = 1, \dots, 10$$

we obtain by linearity:

$$E[X_i] = E[n_i f] = npf = \frac{nf}{10} \approx 1000 \text{ GB} \quad \forall i = 1, \dots, 10$$

Hence:

$$\frac{nf}{10} \approx 1000 \text{ GB} \Rightarrow n \approx \frac{10^4 \text{ GB}}{f}$$

If we want to be more precise we can write:

$$n \lesssim \frac{10^4 \text{ GB}}{f}$$

since all X_i but one are strictly $< 10^4 \text{ GB}$ when the algorithm stops.

The result that $n \propto 1/f$ approximately is quite intuitive and reasonable, since to cover the same fixed distance with smaller steps more time is needed.

For completeness let us remark that in the unluckiest case possible (at least in principle) every file will end up in the same HD, which means that the algorithm will only last $10^4/f$ iterations (since a single hard drive can only fit 1000 GBs of data); hence

$$\frac{10^4 \text{ GB}}{f} \leq n \leq \frac{10^4 \text{ GB}}{f}$$

but this result is purely of academic interest; for all practical purposes it suffices to say n will always be slightly less than $10^4/f$.

Even though this estimate for the upper bound on n was obtained using a somewhat simplistic model it actually yields a pretty close number to the "experimental" results (as shown in the next cell); it also correctly reproduces the scaling of the problem (notice that when we go from $f = 10 \text{ GB}$ down to $f = 1 \text{ GB}$ we save ~ 1 o.o.m. more files, and similarly when f drops to 4 MB we gain another 3 extra o.o.m.%).

```
In [12]: def estimate_n(filesize):
        return 1e4/filesize

for x, f in zip(list_of_drives, list_of_file_sizes):
    observed_n = x.sum()/f
    predicted_n = estimate_n(f)
    print(f"relative difference between observed ({observed_n}) and predicted ({predicted_n}) ({f}) GB: {n
p.abs(predicted_n - observed_n)/predicted_n:.2%}")

relative difference between observed 830.0 and predicted 1000.0 (10 GB): 17.00%
relative difference between observed 9345.0 and predicted 10000.0 (1 GB): 6.55%
relative difference between observed 2493931.000008268 and predicted 2500000.0 (0.004 GB): 0.24%
```

Now that we can somewhat satisfactorily explain/predict the observed values of n we can explain why the efficiency increases, why the histogram gets flatter with decreasing f and why the above approximation gets better with smaller f 's.

The above equations say that the smaller f the longer the algorithm will last; indeed as we already remarked for these values of f it is almost guaranteed to be close to its theoretical maximum. We also notice that even in the worst case scenario we can be sure that, if $f_1, f_2, f_3 = 10 \text{ GB}, 1 \text{ GB}, 0.004 \text{ GB}$ respectively then we will have $n_1 \leq 10n_2$ and $n_2 \leq 1000n_3$ - so that 4 MB files make n inevitably huge.

Using the above we obtain:

$$1000 \leq n_1 \leq 1000$$

$$1000 \leq n_2 \leq 10^4$$

$$2.5 \cdot 10^5 \leq n_3 \leq 2.5 \cdot 10^6$$

This means that the smaller f the more random numbers we sample from the uniform distribution, and therefore the closer we get to having a *flat* distribution (by the law of large numbers; this already explains the observed shapes! But if a small f makes the empirical distribution closer to the theoretical uniform one then it means that the smaller f the closer the bin heights; this in turn guarantees that all bins have a relatively very close to the one of the full drive, which stopped the algorithm. This finally gives two important results:

- With a smaller f all of the X_i 's (i.e. GBs stored in drive i) become almost equal to each other, and in particular to the one which reaches 1 TB; this means that as f becomes smaller, f becomes higher efficiency/percentage of overall occupied storage.
- Due to the above the smaller f the more "constant" the X_i 's become, which means their std gets much smaller since they all get closer to their ideal unknown mean. Notice that as they all get closer to 1 TB the assumption $E(X_i) \approx 1 \text{ TB}$ becomes more and more correct, which makes our $n \approx 10^4/f$ result less of a simplistic approximation and more of an actual prediction; this explains why as f gets smaller this approximation gets better and better.

To recap: with smaller values of f more files will need to be written before any hard drive is full, which means we're sampling a lot more uniform random numbers and therefore obtaining a distribution very close to uniform. While this makes all bins approximately as tall as the one that reached 1 TB, hence we get closer to filling 100% of the available space (while also making our formula more correct).

Notice that this increase in efficiency goes increases monotonically with the decreasing f , which seems to suggest that going even further than 4 MB is desirable; this is actually not true due to other reasons. In particular notice that a typical computer is able to process 4 MB much more efficiently than single byte, hence 4 MB offers a pretty good compromise between storage efficiency and IOPS - while at the same time making individual blocks much lighter to e.g. move around the internet (which is relevant if the hard drives are far away) compared to, say, 1 GB (which would offer comparable efficiency).

3.4

Let's organize our results in a small dataframe for visual clarity's sake.

```
In [13]: df = pd.DataFrame(np.vstack((list_of_file_sizes, list(map(np.mean, list_of_drives))), list(map(np.std, l
ist_of_drives))))
df.columns = ["filesize", "mean", "std"]
display(df) # obviously the mean is just tot. GB written / 10, i.e. we obtain x.sum() once again but div
ided by 10
```

	filesize	mean	std
0	10.000	830.0000	90.664216
1	1.000	934.5000	33.767588
2	0.004	997.5724	1.528796

We already explained why the mean scales as $1/f$: smaller files \Rightarrow more files are written \Rightarrow more numbers are sampled from the uniform distribution \Rightarrow the resulting histogram is closer to flat \Rightarrow the rectangles' heights get closer \Rightarrow each drive's content becomes closer to the theoretical uniform one then it means that the smaller f the closer the bin heights; this in turn guarantees that all bins have a relatively very close to the one of the full drive, which stopped the algorithm. This finally gives two important results:

This suffices to explain the decreasing σ behaviour: since a smaller and smaller f makes the drives closer and closer to full we obtain that smaller $f \Rightarrow$ almost constant \bar{x} components \Rightarrow smaller σ

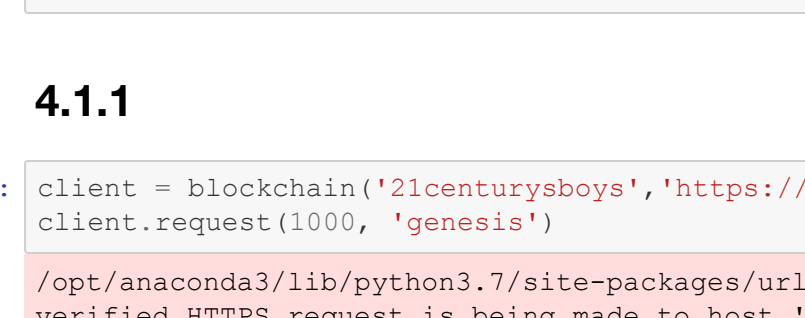
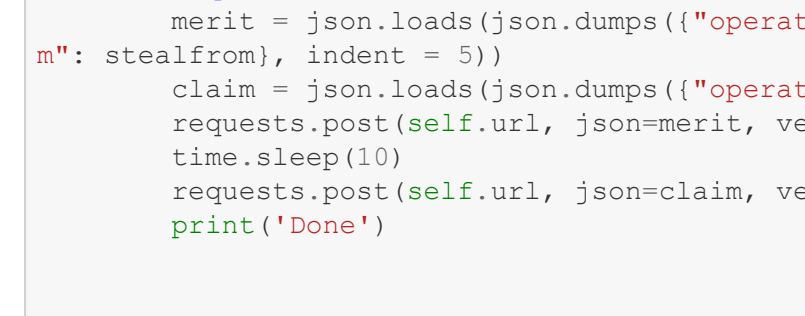
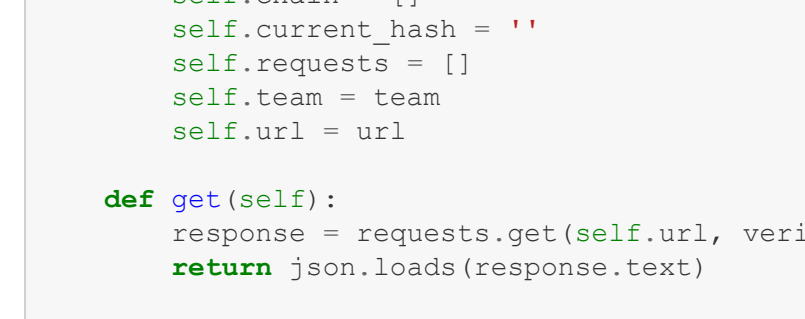
Let's now repeat the experiment for many times and plot how the random variable "% of occupied space" is distributed.

```
In [14]: #init
def obtain_drive_statistic(filesize, n_experiments = 1000, percent = True):
    y = np.zeros(n_experiments)
    for i in range(n_experiments):
        y[i] = populate_drives(filesize).sum()
    if percent:
        return y/1e4
    else:
        return y

def plot_drive_statistic(filesize, n_experiments = 1000):
    y = obtain_drive_statistic(filesize, n_experiments)
    fig, ax = plt.subplots()
    ax.hist(y)
    ax.set_title(f"({n_experiments} experiments, {filesize} GB filesize)")
    ax.set_xlabel("% of occupied space")
    return y

In [15]: # since the histogram for 4 MB takes a bit of time to speed up debugging we save the results
save = True

if not os.path.exists("df1.csv"):
    l = [plot_drive_statistic(f) for f in list_of_file_sizes]
    df1 = pd.DataFrame(np.array(l), columns = [str(f) for f in list_of_file_sizes])
    if save:
        df1.to_csv("df1.csv", index = False)
    #display(df1)
else:
    df1 = pd.read_csv("df1.csv")
    for i in range(3):
        fig, ax = plt.subplots()
        ax.hist(df1.iloc[:,i])
        ax.set_title(f"1000 experiments, {df1.columns[i]} GB filesize")
        ax.set_xlabel("% of occupied space")
```



We obtained gaussian-like distribution. This is not surprising due to the central limit theorem; indeed we expect that as the number of experiment grows this distribution becomes closer and closer to a gaussian with mean equal to the "true" mean (with the same holding for the standard deviation). Due to what we already explained the "true" mean scales inversely with f , whereas the "true" std scales proportionally with f ; this explains why the last histogram is so much "lighter" horizontally and more on the right wrt the first, etc.

```
In [16]: # even though 1000 experiments aren't that many they nonetheless allow us to obtain more solid estimate
s for the random variable "% of occupied space"

def experiments_summary(x, f):
    n = x.mean()
    return f"occupied storage ({f} GB files): {n:.2%} +/- {n:.2%}"

for i in range(3):
    print(experiments_summary(df1.iloc[:,i], list_of_file_sizes[i]))

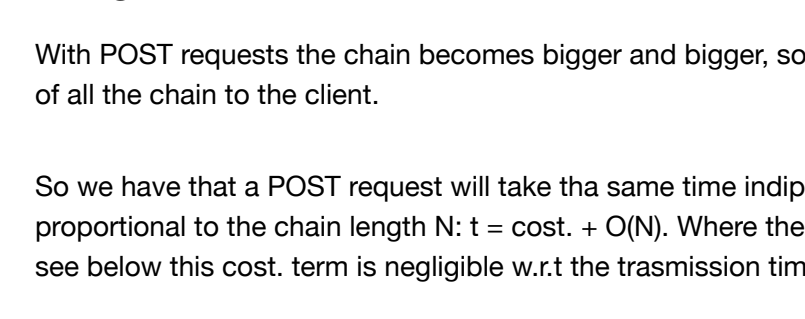
occupied storage (10 GB files): 84.89% +/- 4.37%
occupied storage (1 GB files): 95.26% +/- 1.48%
occupied storage (0.004 GB files): 99.69% +/- 0.10%
```

```
In [17]: # if we want we can do fancy things like a kde

for i in range(3):
    sns.kdeplot(100*df1.iloc[:,i], label = list_of_file_sizes[i])

plt.xlabel("% of occupied storage")
plt.legend(loc = "best");

# indeed the 4 MB gaussian is much tighter than the others, and almost centered around 100%
```



Appendix: more rigorous derivation of the n average upper bound estimate computed above

A more careful statistical modelling of the problem can let us obtain results analogous to what we already derived, but in a less "hand-way" fashion; let's start with some notation.

Let n_i ($i = 0, \dots, 9$) be the number of files in hard drive i . If the whole process stopped when e.g. drive 0 is full (instead of the first which reaches 1 TB) then the vector (n_0, \dots, n_9) would be distributed according to a negative multinomial with $\tilde{p} = (0.1, \dots, 0.1)$ (since each drive has a 10% prob. of being "chosen" at any given iteration) and $r = 1000 \text{ GB}/f$ (since by definition r is the number of successes in slot 0 after which the process stops). Hence we write:

$$(n_0, n_1, \dots, n_9) \sim \text{NM} \left(r = \frac{1000 \text{ GB}}{f}, \tilde{p} = (0.1, 0.1, \dots, 0.1) \right)$$

The above relation holds independently of which drive decides when the process stops.

In our problem, though, it's not necessarily drive 0 which stops the algorithm; this means we need to do the following.

- 1) First we marginalize the above wrt every index but one; by symmetry this will give us the distribution of any n_i .
- 2) Then we compute the distribution of the sum of all n_i 's, because this will give us the distribution of $n = \sum_{i=0}^9 n_i$.

According to [this wikipedia page](#) the following two facts hold:

- The marginal of a negative multinomial wrt variable i is a negative binomial with the same r and p_i rescaled in such a way that $p_0 + p_i = 1$. This lets us solve point 1); indeed we obtain

$$n_i \sim \text{NB} \left(r' = r = \frac{1000 \text{ GB}}{f}, p' = p_i = \frac{1}{2} \right)$$

- The random variable "sum of two positive/negative binomial/multinomial distributions with the same p/p' " is a distribution of the same kind, with the same p and the other parameter equal to the sum of the two r 's; this lets us solve point 2); since if combined with the previous bullet point it implies that

$$n \sim \text{NB} \left(r = \frac{10^4 \text{ GB}}{f}, p = \frac{1}{2} \right)$$

Notice that n 's distribution cannot actually be a NB, since the probability of n surpassing the max value is 0 (whereas the NB has no such constraint); this means that what we're doing is not really "cutting" the rightmost part of the distribution, which means that (n) is slightly less than $\langle \text{NB}(r = 10^4 \text{ GB}/f, p = 1/2) \rangle$. Since

$$\left\langle \text{NB} \left(r = \frac{10^4 \text{ GB}}{f}, p = \frac{1}{2} \right) \right\rangle = \frac{pr}{1-p} = \frac{10^4 \text{ GB}}{f}$$

we obtain

$$\langle n \rangle \lesssim \frac{10^4 \text{ GB}}{f}$$

which is the same result we already obtained, but in a more "professional" way. As a bonus we can also produce approximations for other interesting quantities, such as the std.

Assignment 4: Rest APIs & Block Chain Technology

Below there is a simple client implementation with its usage.

```
In [25]: import requests
import json
import time

class blockchain:

    def __init__(self, team, url):
        self.accounts = []
        self.chain = []
        self.current_hash = ''
        self.requests = []
        self.team = team
        self.url = url

    def get(self):
        response = requests.get(self.url, verify = False)
        return json.loads(response.text)

    def request(self, coin, stealfrom):
        merit = json.loads(json.dumps({"operation": "merit", "team": self.team, "coin": coin, "stealfro
m": stealfrom, "indent = 3}))
        claim = json.loads(json.dumps({"operation": "claim", "team": self.team, "indent = 3}))
        requests.post(self.url, json=merit, verify = False)
        time.sleep(10)
        requests.post(self.url, json=claim, verify = False)
        print('Done!')
```

4.1.1

```
In [26]: client = blockchain('21centuryboys', 'https://pansophy.app:8443')
client.request(1000, 'genesis')

/opt/anaconda3/lib/python3.7/site-packages/urllib3/connectionpool.py:1004: InsecureRequestWarning: Un
verified HTTPS request is being made to host 'pansophy.app'. Adding certificate verification is stron
gly advised. See: https://urllib3.readthedocs.io/en/latest/advanced-usage.html#ssl-warnings
InsecureRequestWarning,
```

```
Done
/opt/anaconda3/lib/python3.7/site-packages/urllib3/connectionpool.py:1004: InsecureRequestWarning: Un
verified HTTPS request is being made to host 'pansophy.app'. Adding certificate verification is stron
gly advised. See: https://urllib3.readthedocs.io/en/latest/advanced-usage.html#ssl-warnings
InsecureRequestWarning,
```

4.1.2

```
In [27]: int(24*60*60/10)

Out[27]: 8640
```

4.2.1

This function computes the block hash that will be in the data of the next block in the chain. If one block is modified and we calculate the new hash we see that it is different w.r.t. to the hash (of the same block) computed at the creation of the next block in the chain. So if the 2 hashes are identical we are sure that the block is unchanged and repeating the operation block by block like a 'chain' we can verify that nobody compromised the all blockchain.

4.2.2

It's simple: we get all the block's data using a GET request; we start from the first block and we compute its hash; if the hash is the same reported in the second block we repeat the same operation on the second block, on the third and so on... at the last block we have validated all the chain.

4.2.3

With POST requests the chain becomes bigger and bigger, so the GET response of the server might take much time to transmit the JSON file of all the POST to the client.

So we have that a POST request will take the same time independently from the chain length. Instead the GET request will take a time proportional to the chain length $N: 1 \leq \text{cost} \propto O(N)$. Where the cost term is due for the cryptography of the data. How we can see below this cost term is negligible w.r.t. the transmission time and it is about 100 ms.

```
In [28]: #time client.request(1000, 'genesis')

/opt/anaconda3/lib/python3.7/site-packages/urllib3/connectionpool.py:1004: InsecureRequestWarning: Un
verified HTTPS request is being made to host 'pansophy.app'. Adding certificate verification is stron
gly advised. See: https://urllib3.readthedocs.io/en/latest/advanced-usage.html#ssl-warnings
InsecureRequestWarning,
```

```
Done
CPU times: user 20.4 ms, sys: 6.51 ms, total: 26.9 ms
Wall time: 10.3 s
/opt/anaconda3/lib/python3.7/site-packages/urllib3/connectionpool.py:1004: InsecureRequestWarning: Un
verified HTTPS request is being made to host 'pansophy.app'. Adding certificate verification is stron
gly advised. See: https://urllib3.readthedocs.io/en/latest/advanced-usage.html#ssl-warnings
InsecureRequestWarning,
```

We report a small percentage of the data printed by the comment above just for illustration purposes:

```
CPU times: user 74.1 ms, sys: 26.5 ms, total: 101 ms
Wall time: 712 ms
Out[29]: {'accounts': ['21centuryboys', '1AC Milan', '1.Andreas J. Peters', '1.AntonVilladiVilla', '401016[...]]
```

```
In [23]: import matplotlib.pyplot as plt
import numpy as np

plt.plot(range(1000), 0.1 + 0.1*np.arange(1000), label = 'GET')
plt.plot(range(1000), (0.1 + 10)*np.array([1]*1000), label = 'POST')
plt.legend(loc = 'best')
plt.xlabel('chain length')
plt.ylabel('time request')
plt.tick_params(
    #size="x",
    which="both",
    bottom=False,
    left=False,
    labelbottom=False,
    labelleft=False) # Labels along the bottom edge are off
```


4.2.4

The server can recompute the current account balances just starting from the first block and recomputing all the transactions having care of checking the hash.

4.2.5

One advantage of using the REST API is the facility with which we can create a client object that interacts with the server. We only need two functions (in this case) for our purpose: GET and POST. Furthermore the HTTPS protocol allows for secure (with the SSL cryptography) and authenticated (not in this case) communication between client and server. The JSON data format's advantages are the facility with the server can handle the coming data (for example in a data streaming process) and the client can create the data to send. One disadvantage can be the scalability problem of the GET request discussed in 4.2.3.

```
In [ ]:
```