

Sample Python Extraction Problem

Extract the following sample data into a Pandas dataframe.

Sample Data:

Chilli - 250

Beans long - 250

Bitter gourd - 500

Coccinia - 500

Garlic - 250

Radish white - 250

Carrot ooty - 500

Ladies finger - 500

Onion - 2 kg

Potato - 1 kg

Sambar cucumber - 500

Banana robust - 2 kg

Bitter gourd - 500gms

Brinjal purple big one - 1

Capsicum green - 2-3 pieces

Baby bottle gourd - 1

Banana yellaki - 1/2 kg

Arbi 1/2 kg

Potato - 1kg

Bitter gourd 250 gms

Radish white with leaves 4-5

Delhi carrot 500 gms

English Cucumber 500gms

Bathua 2-3 bunches

Sarson 3 Bunches

Beetroot 500 grams

1 cauliflower

Nati dhaniya 1 bunch

Apple Washington 500gms

Banana yellaki 500 gms

1 Broccoli

Baby corn 1pkt

Mushrooms 1pkt

2bunch coriander 1 methi 1 mint 500 gm delhi carrot 1 kg local tomato 1kg fresh green peas
500gm seedless black grapes250 gm beetroot

500gm lady finger 1 kg fresh n tender green peas 250gm ridge gourd 250 beans long 250gm
capsicum 100gm amla 100gm rose

Onion - 1 kg @ 1
Tomato natti - 1 kg
Yelaki banana - 1 kg

1 Bunch coriander
1 Broccoli
1/2 kg Cucumber salad
1/2kg apple Washington
1/2 kg banana Yelakki
1/2 kg seedless grapes
1/2 kg delhi carrot

Beans haricot-500gms
Delhi carrot-1 kg
Cucumber english-1kg
Cucumber salad-500gms
Ladies finger-500gms
Cauliflower- 1 pcs
Avarakai-500gms
Banana yelaki-1kg(medium ripen)
Papaya-1(ripen)

Several variations and cases need to be handled:

- Some entries have line breaks, others do not
- Number is at the start of the line, other times it's in the middle
- Might or might not be a space between the quantity and the unit of measure (1pkt vs 500 gms)
- Quantities might be vague - e.g. 2-3
- Might have extra information at the end (e.g. ripen)
- Various separators are possible - dashes, commas, etc.
- Extra words e.g. Brinjal purple big one

The solution might need to be based on a combination of RegEx, simple extraction rules and other, more complex NER type approaches.

