



Centrālā statistikas  
pārvalde

# Alternatīvo datu vākšanas rīku izmantošana cenu statistikā

Patēriņa cenu indeksu daļa  
**Dāvis Liberts, Nataļja Dubkova**

15. Maijs 2019



Centrālā statistikas  
pārvalde

**Pieaugošs preču sortiments**

**Tehnoloģiju attīstība**

**Pirkumi interneta veikalos**

**Pakalpojumu digitalizācija**

**Dinamiskās cenas**



Centrālā statistikas  
pārvalde

## Šajā prezentācijā:

- ❑ Automātiskā datu vākšana internetā (web-scraping)
- ❑ Web-scraping datu pielietojums cenu statistikā
- ❑ Citu valstu pieredze
- ❑ CSP paveiktais
- ❑ Izaicinājumi šo datu iegūšanā un pielietošanā



Centrālā statistikas  
pārvalde

## Kā tas sākās

Lai arī cenu statistikā ar tādiem jēdzieniem kā **Web-scraping** (Big Data) saskaramies ikdienā arvien biežāk tikai salīdzinoši nesen, to vēsture ir daudz senāka – un to pirmsākumi meklējami līdz ar pirmajām interneta meklētājprogrammām.

Cilvēkiem radās nepieciešamība iegūt informāciju vienlaicīgi no vairākām tīmekļa vietnēm – savukārt ne katrā no tām bija pieejama lejupielādēšanas iespējas, kā arī manuāla kopēšana ar roku bija "garlaicīga" un neefektīva.

Tas bija aizsākums pirmajiem interneta datu vākšanas robotiem.



Centrālā statistikas  
pārvalde

## Dati mums visapkārt

Mūsdienās šāda datu vākšana ir gājusi tikai plašumā, ko veic neskaitāmi uzņēmumi un organizācijas dažādu pētījumu, mārketinga vai tirgus analīzes nolūkiem.

Var iegūt praktiski jebkurus datus, ja vien tie ir pieejami internetā.

To nodrošina neskaitāmi gan maksas vai bezmaksas rīki, gan individuāli izstrādātas programmas.



Centrālā statistikas  
pārvalde

# Web-scraping darbības princips

## Divi galvenie procesi



**Interneta lapas  
ielāde**

**Datu izgūšana no  
lapas**

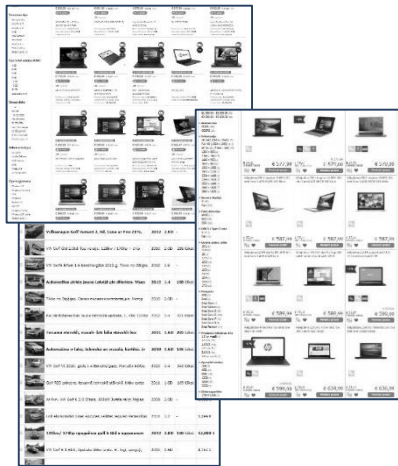


Centrālā statistikas  
pārvalde

# Automātiskā cenu vākšana

## WEB-SCRAPING

## DATU GLABĀŠANA



## DATU VALIDĀCIJA

## DATU ANALĪZE



Centrālā statistikas  
pārvalde

# Web-scraping darbības princips

Automātiski  
tiek ievadīti  
nepieciešamie  
parametri

The screenshot shows a web page for laptops. On the left is a sidebar with filter categories: **Procesora tips** (Intel Celeron, Intel Core i3, Intel Core i5, AMD, Intel Pentium, Intel Atom, Intel Core i7, AMD A series, Pārdoši vairāk...), **Operatīvā atmiņa (RAM)** (4 GB, 6 GB, 8 GB, 12 GB, 16 GB, 32 GB, 40GB DDR3 RAM), **Cietais disks** (8 TB, 650 GB, 128 GB (SSD), 256 GB (SSD), 32 GB (SSD), Nav informācijas, 64 GB (SSHD), 1 TB - (128 SSD), Pārdoši vairāk...), **Videokarte (tips)** (Integrēta, Nvidia GeForce, AMD Radeon, AMD, Integrēta, Nvidia Quadro), and **Operētājsistēma** (Windows 10, Windows 10 Home, DOS, Windows, FreeBSD OS, Mac OS, Windows 10 Pro, Windows XP Home, Pārdoši vairāk...). The main area is a grid of laptop products. One product, an HP 3Q4EE laptop, is highlighted with a red box. Red arrows point from the sidebar and the highlighted product to callout boxes on either side.

Tiek atlasīta  
kāda konkrēta  
prece vai preču  
grupa

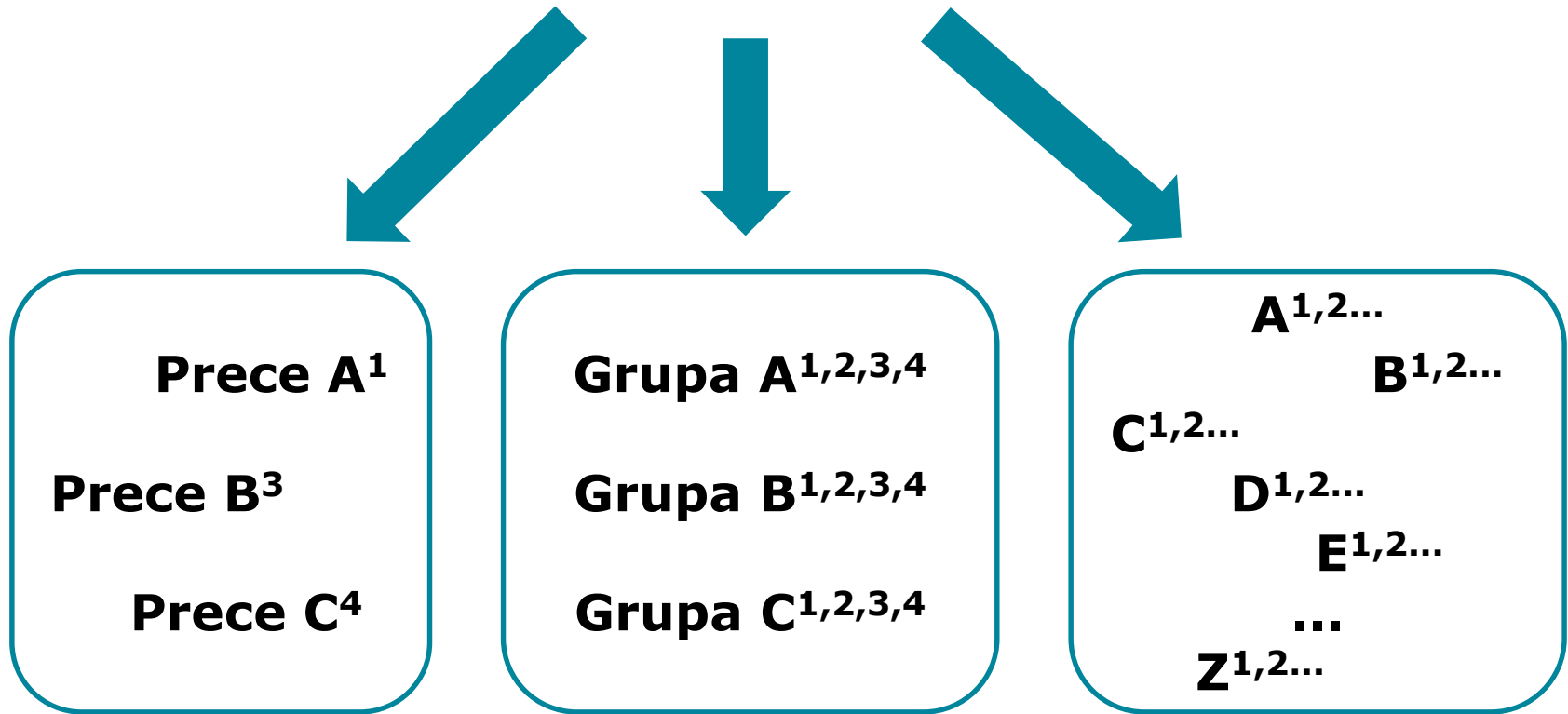
Atlasītie dati tiek strukturizēti un saglabāti tālākai datu analīzei





Centrālā statistikas  
pārvalde

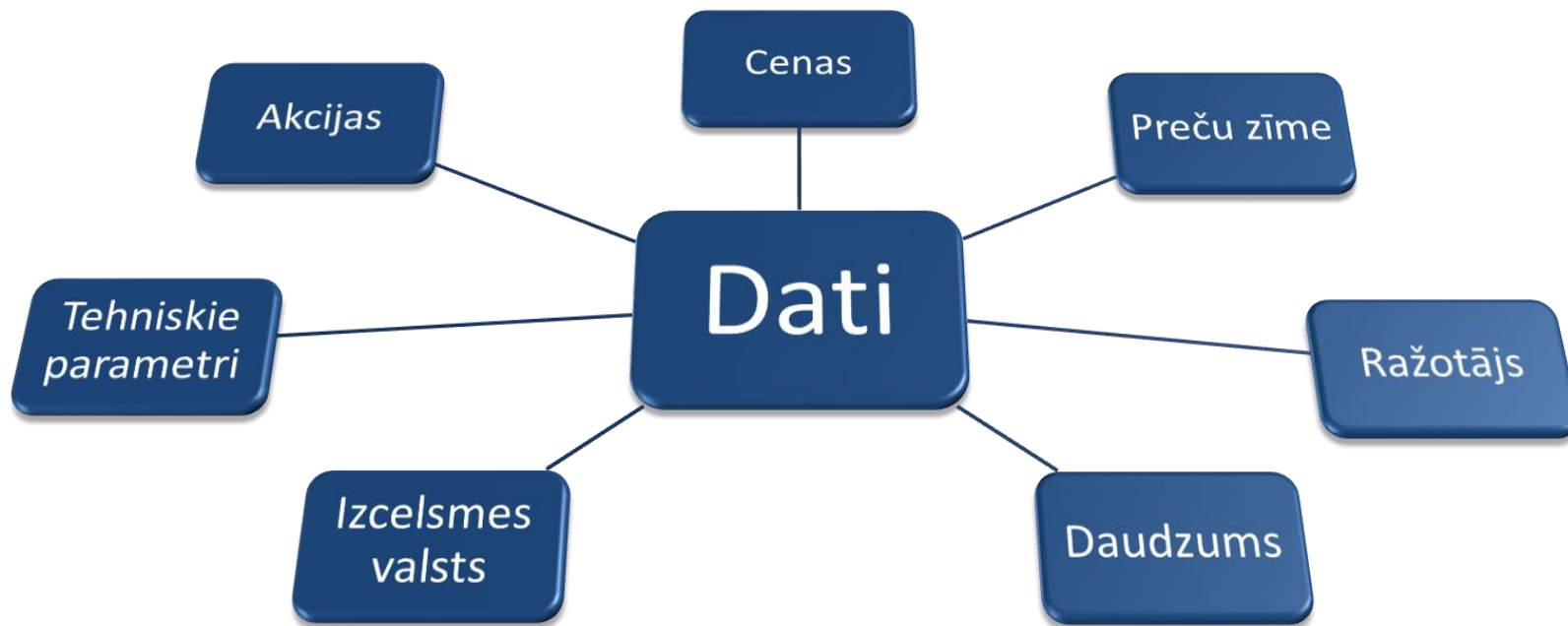
# Web-scraping metodes cenu statistikā





Centrālā statistikas  
pārvalde

# Informācija cenu statistikai





Centrālā statistikas  
pārvalde

## Pielietojums cenu statistikā

- ❑ Automātiskā datu vākšana internetā dod iespēju veikt cenu mērījumus ievērojami biežāk un plašākā apjomā
- ❑ Iespēja optimizēt datu vākšanas procesu
- ❑ Samazina respondentu noslodzi
- ❑ Būtiski vākt cenas tur, kur tiek izdarīti pirkumi
- ❑ Risinājums jauniem izaicinājumi cenu statistikā saistībā ar pieaugošu preču klāstu, tehnoloģiju attīstību un digitalizāciju



Centrālā statistikas  
pārvalde

# Eurostat skatījums

- ❑ Eurostat saskata potenciālu web-scraping izmantošanā cenu statistikā
- ❑ Cenu statistika pāriet no «pildspalvas un papīra» uz inovatīvākām metodēm
- ❑ Digitalizācija arvien vairāk ienāk cilvēku dzīvē un maina pirkumu paradumus - interneta veikali, digitālie pakalpojumi
- ❑ Jauni datu avoti prasa arī jaunas metodes
- ❑ Izaicinājums nākotnē - jaunu datu avotu un metožu interpretācija un komunikācija ar datu lietotājiem



Centrālā statistikas  
pārvalde





Centrālā statistikas  
pārvalde

## Citu valstu pieredze

- ❑ Ir valstis, kas vāc liela apjoma datus par visām preču grupām. Ir valstis, kas sāk ar datu vākšanas automatizāciju, atbilstoši iepriekš noteiktai izlasei
- ❑ Vācot liela apjoma datus par visām pieejamām preču grupām un precēm, liels izaicinājums ir to klasificēšana – tādām nolūkam pēta **machine learning** iespējas
- ❑ Vācijas kolēģi veic pētījumu par dinamiskajām cenām un to izaicinājumiem cenu statistikā – ap 2700 precēm cenu vākšana tiek veikta reizi stundā



Centrālā statistikas  
pārvalde

# CSP paveiktais saistībā ar web-scraping

Eurostat seminārs  
par Web-scraping



Eurostat seminārs  
par Web-scraping



Eurostat seminārs  
par Web-scraping



Uzsākts    Noslēgts

Granta  
projekts  
(Web-scraping  
kā datu avots  
SPCI)



Testa datu  
vākšana  
(turpinās)

2015

2016

2017

2018

2019







Centrālā statistikas  
pārvalde

# Priekšrocības

- ❑ Iespēja iegūt vairāk datu nekā tas iespējams manuāli
- ❑ Iespēja iegūt datus pat pa stundām, dienām vai nedēļām
- ❑ Cenu vākšanas laiks:

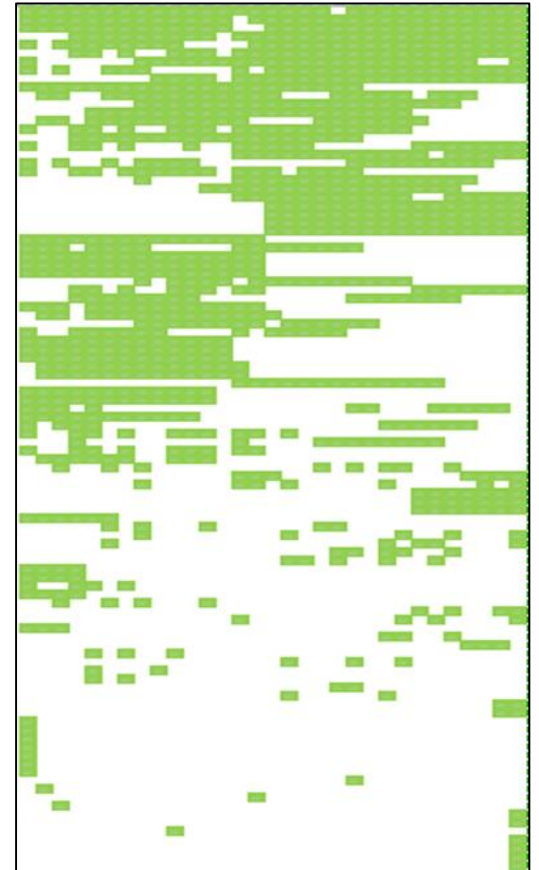
Preču grupa	Manuāli vācot cenas	Web-scraping
3-5 g lietotas automašīnas	~ 3 stundas	~ 10 minūtes
9-11 g lietotas automašīnas	~ 2 stundas	~ 21 minūte
Dzīvokļu īres cenas	~ 5 stundas	~ 30 sekundes

- ❑ Web-scraping viennozīmīgi negarantē labāku datu kvalitāti, bet noteikti daudz lielāku datu apjomu
- ❑ Papildus informācija par tirgus tendencēm



Centrālā statistikas  
pārvalde

# Datu analīzes iespēju piemēri

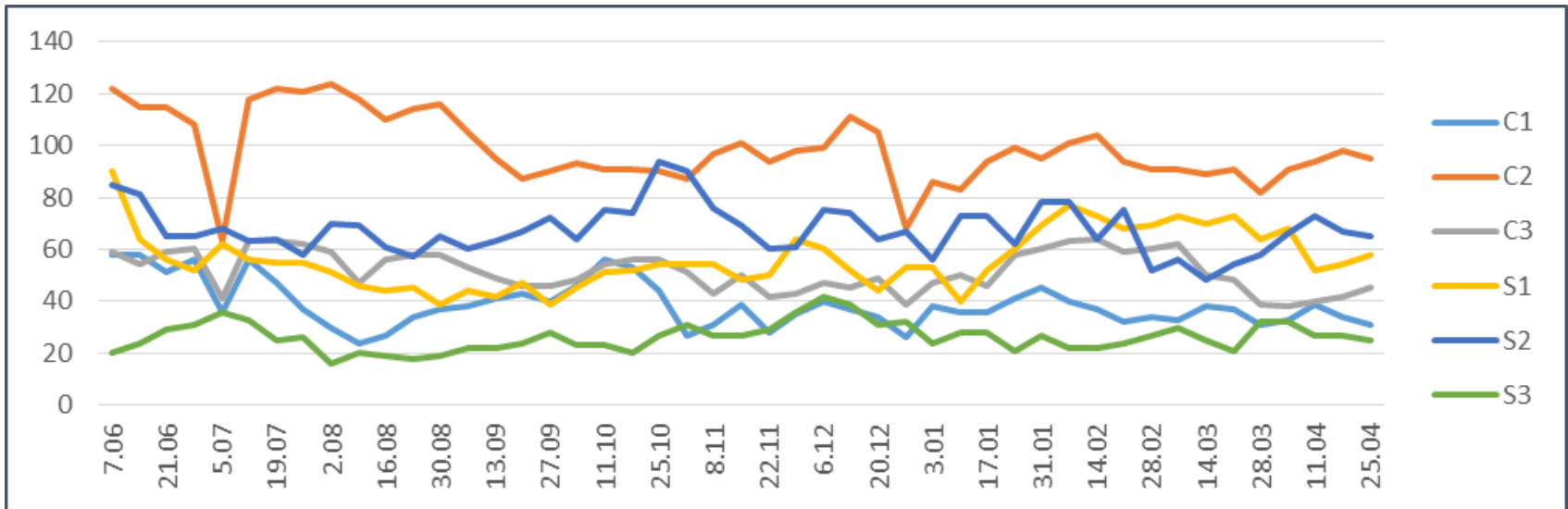


Planšetdatoru pieejamība tirdzniecības vietās



Centrālā statistikas  
pārvalde

# Datu analīzes iespēju piemēri



Dzīvokļu ģres sludinājumu skaits sadalījumā pa segmentiem



Centrālā statistikas  
pārvalde

# Izaicinājumi

- ❑ Tehnoloģiski sarežģīts datu vākšanas process
- ❑ Web-scraping rīks jāpielāgo katram veikalam
- ❑ Nepieciešama pastāvīga procesa uzraudzība un IT atbalsts  
*(tehniskas problēmas, izmaiņas mājas lapās)*
- ❑ Metodoloģiski izaicinājumi  
*(izlases apjoms, cenu vākšanas biežums, datu segmentēšana, datu integrēšana patēriņa cenu indeksu aprēķinā)*
- ❑ Cenu indeksu izstrāde kļūst arvien sarežģītāka - interesantāka



Centrālā statistikas  
pārvalde

# Pārmaiņas cenu reģistrēšanas procesā

## CENU VĀKŠANA MANUĀLI



Cenu reģistratori



Datu analītiķi

## CENU VĀKŠANA AR WEB-SCRAPING



Pastāvīgs IT  
atbalsts



Centrālā statistikas  
pārvalde

**Paldies par Jūsu uzmanību!**

Davis.Liberts@csb.gov.lv  
Natalja.Dubkova@csb.gov.lv

[www.csb.gov.lv](http://www.csb.gov.lv)