

Medius[®]

With passion.



Enostavna gradnja podatkovnih tokov za usmerjanje, analizo ter učenje v “BigData” svetu

MakeIT 2022, 3.6.2022



David Šenica

- Podatkovni znanstvenik
- Inženir za strojno učenje



Marko Polak

- Arhitekt
- Vodja projektov



O Mediusu



Medius



20

Years of
Innovations

30+

Certified
Engineers

Java

& Open Source
Experts

AAA

Gold
Creditworthiness

ISO

9001 for
SW Development

Internacionalne nagrade



Medius



United Nations
Public Service
Award

tmforum

TM Forum
Excellence
Award



EuroCloud Best
Cloud Service
by Startup



EUREKA-ITEA
Achievement
Award



Best Blockchain
Startup at CV
Competition



- Big Data in strojno učenje
- Kontekstno odvisni podatkovni tokovi
- Opis demo primera Spletne trgovine in komponente sistema
- Live Demo



Big Data



Medius

Big Data

(Velika količina podatkov)



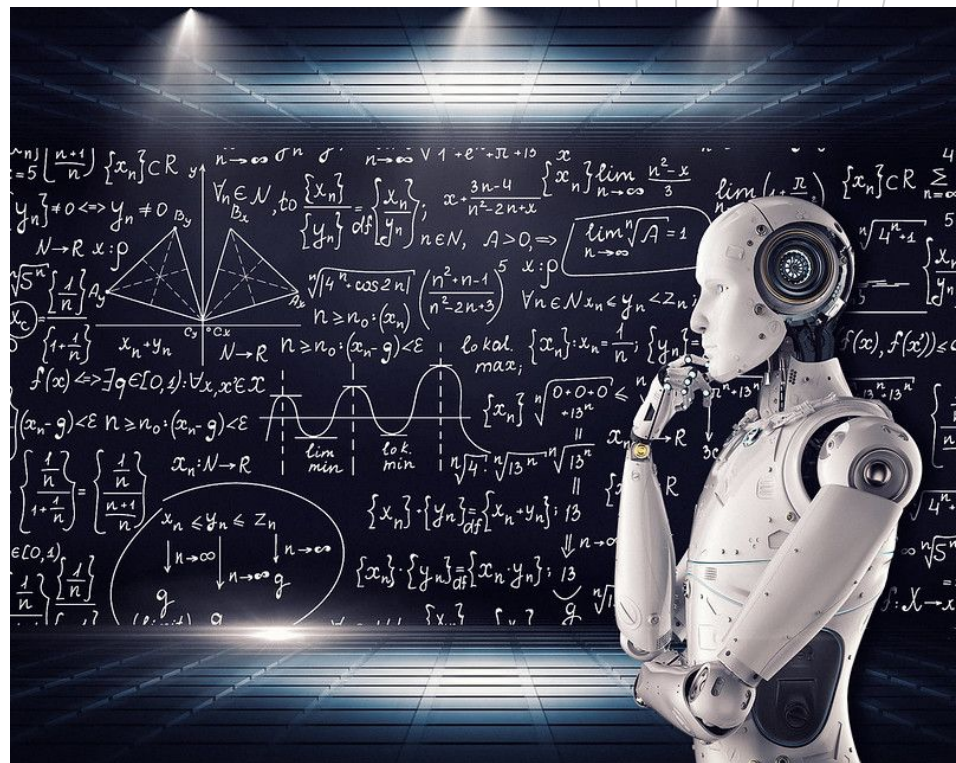
- Tradicionalno: velika količina podatkov (**prevelika**)
- Izzivi
 - zajem
 - shranjevanje
 - analiza
 - iskanje
 - deljenje
 - vizualizacija
 - ...



- Bolj moderno:
 - prediktivna analiza
 - ugotavljanje vedenjskih vzorcev uporabnikov (UBA - User behavior analytics)
 - ali kakršnakoli napredna analiza, ki prinaša dodano vrednost iz podatkov



- reševanje problemov
 - napovedi odpovedi
 - detekcija prevar
 - ...
- podpora našim poslovnim procesom
 - segmentacija kupcev
 - sentimentalna analiza
 - ...

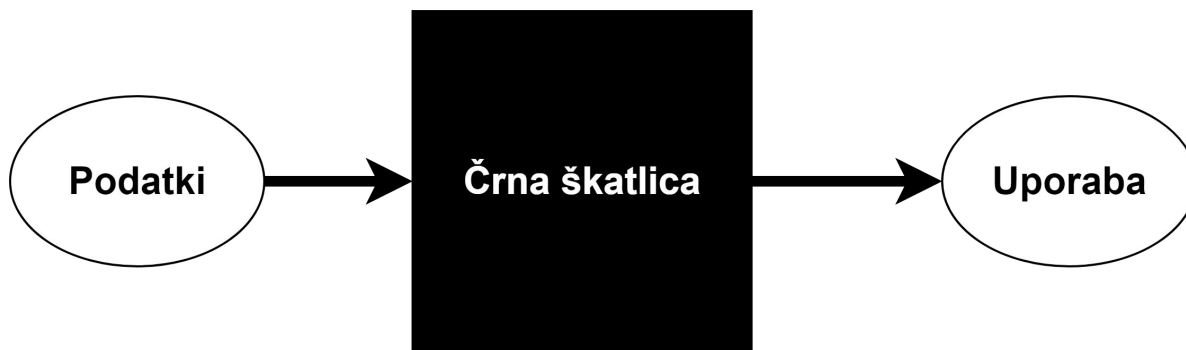


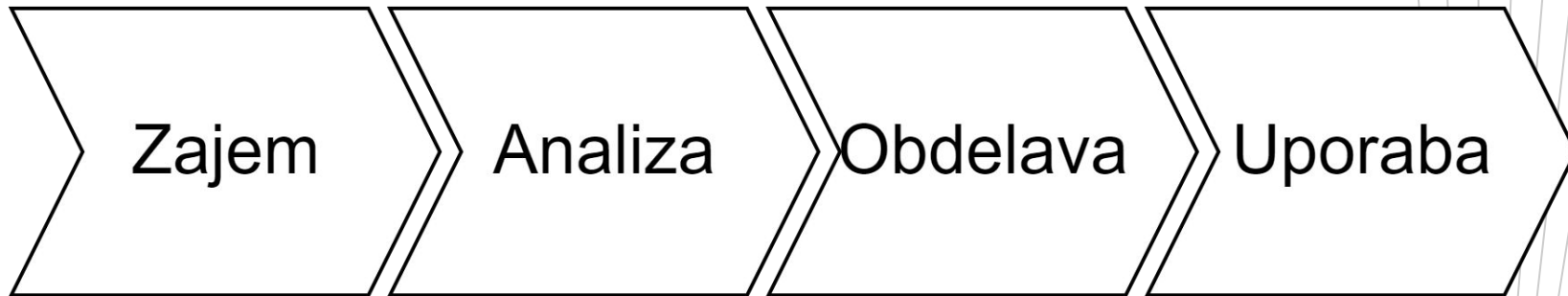
Big Data in strojno učenje - uporabnost

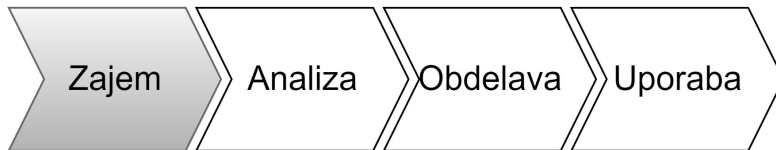


Medius

- optimizacija poslovanja
- hitrejša in boljše strateško odločanje
- personalizacija uporabniške izkušnje
- dvig vrednosti podjetja (revenue)
- itd.





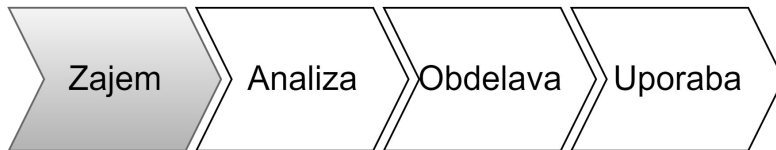


- statični podatki v relacijskih bazah/DWH/NoSQL baze
- komunikacija med aplikacijami(ESB, apache proxy, Istio,...)
- js za sledenje uporabnikom na spletni strani
- namensko spremljanje predefiniranih dogodkov in pošiljanje na zaledni sistem(vgraditev v aplikacije)
- shranjevanje v datoteko
- ...

Big Data - zajem podatkov - protokol



Medius



- standard out
- datoteka
- rest/soap
- namenski vmesnik
- ...



- označevanje(tagging)
- klasifikacija
- kategorizacija
- kontekstualizacija
- ...



Big Data - obdelava



Medius

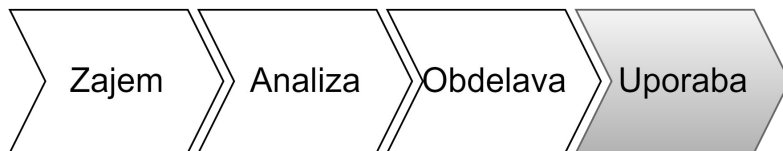


- standardizacija
- transformacija
- povezovanje nestrukturiranih podatkov in strukturiranih podatkov
- učenje in priprava modelov
- ...

Big Data - izpostavljanje rezultatov



Medius



- shranjevanje podatkov v bazo
- paketne obdelave(v nočnem času)
- ...
- namenski servis (predikcijski, evalvacijski,...)
 - vsem dostopni
 - evalvacija/predikcija v realnem času



Sedaj smo pripravljeni, da se podamo v Big Data svet in v uporabo strojnega učenja :D

A pa smo res?





Gartner predvideva, da bo 85% projektov strojnega učenja propadlo.

**Nočemo biti na tem
vlakcu!**

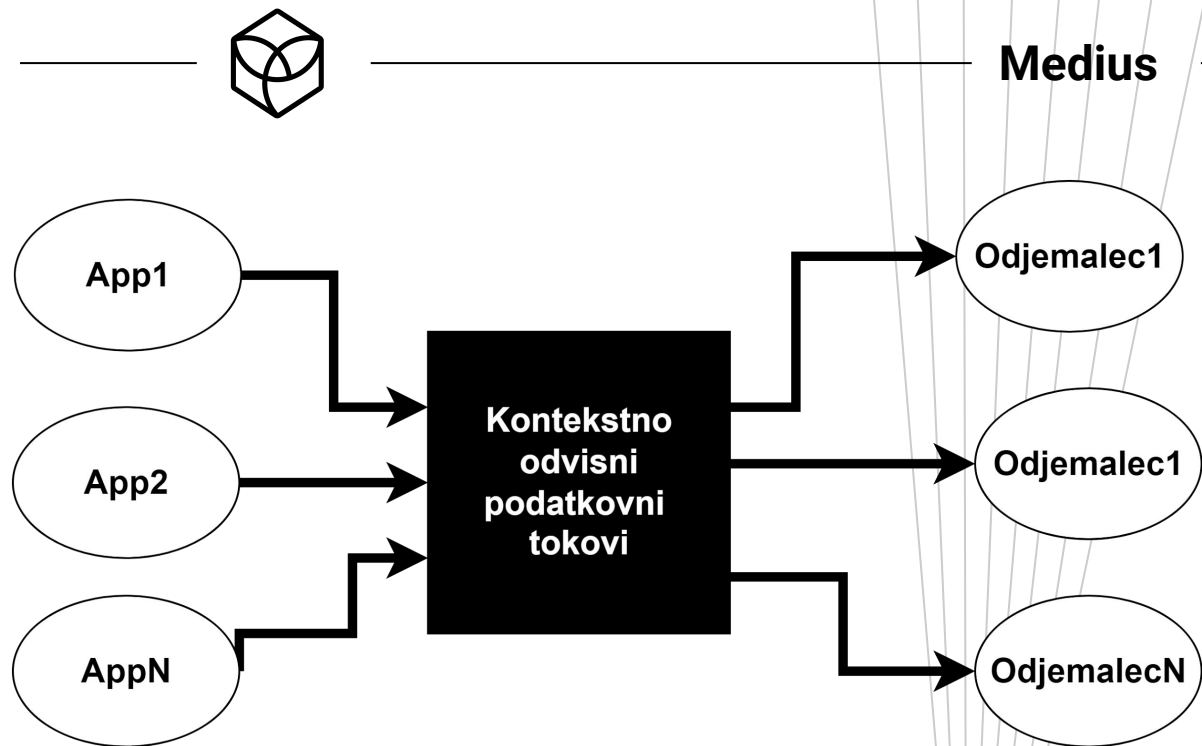




Kontekstno odvisni tokovi podatkov

Kontekstno odvisno tokovi podatkov

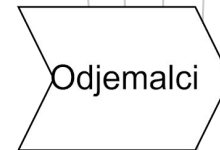
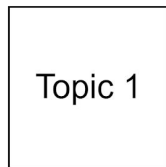
- Zajem podatkov s čim večih virov
- Uporabnost podatkov za čim več odjemalcev



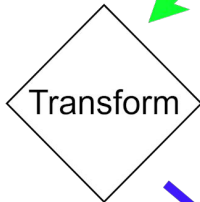
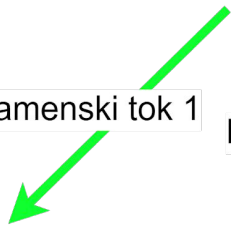
Kontekstno odvisno tokovi podatkov



Medius



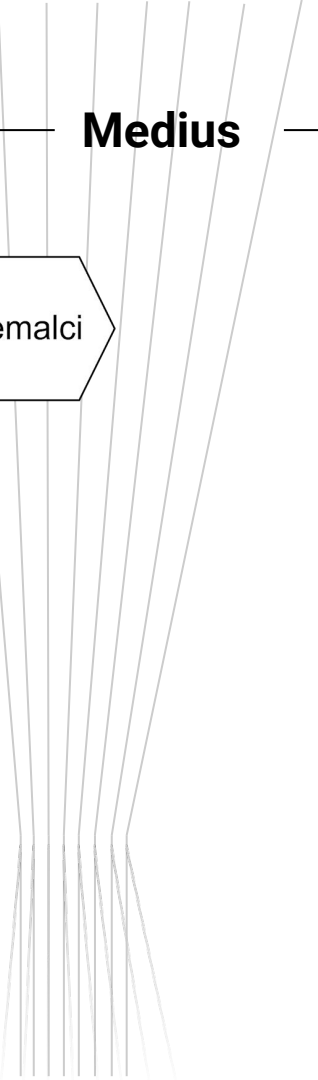
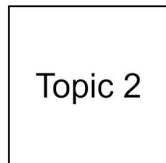
Namenski tok 1



Namenski tok 2



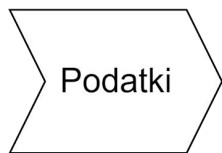
Obogaten tok



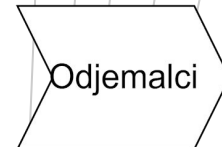
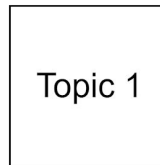
Kontekstno odvisno tokovi podatkov



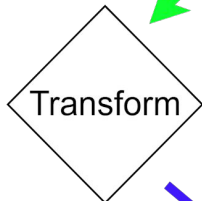
Medius



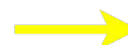
Glavni tok podatkov



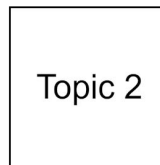
Namenski tok 1



Namenski tok 2



Obogaten tok



- divide and conquer
- analiza
- obogatitev
- povezovanje
- transformiranje
- ...



Zakaj?

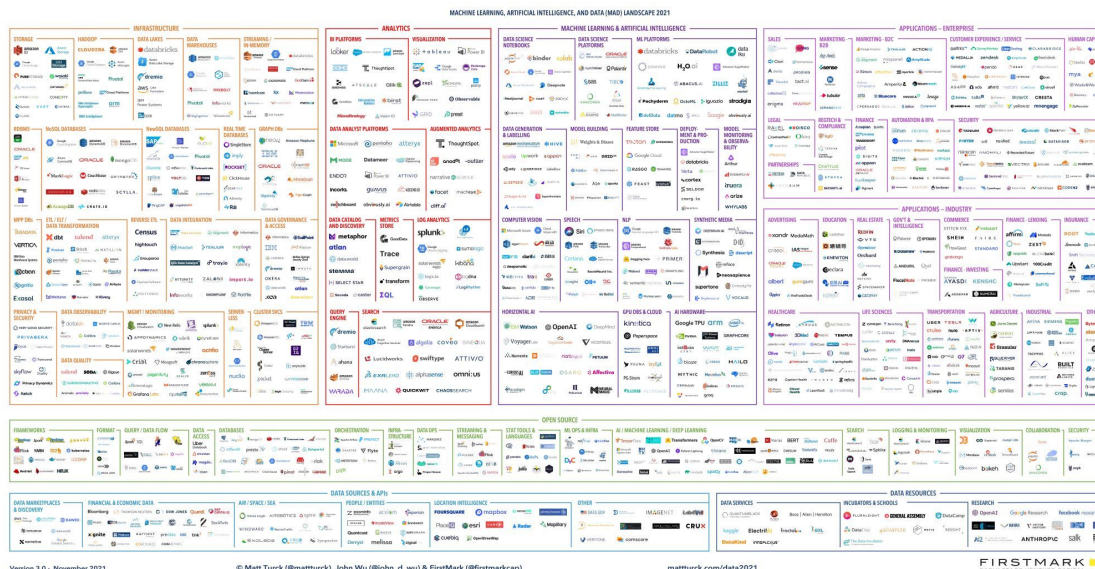
- ker je bistveno, da damo podatkom uporabno vrednost (z ali brez strojnega učenja)
- da imamo podatke pripravljene za trenutek, ko jih bomo rabili - podatki so osnova za izvajanje strojnega učenja
- da imamo možnost obogatiti podatke - glede na kontekst
 - kontekst podatkov
 - kontekst uporabe (kaj želimo)

Kontekstno odvisni tokovi podatkov - komponente sistema



Medius

- Katere komponente izbrati?
- Kaj upoštevati pri izbiri?
- Naj bo zabavno in enostavno :D

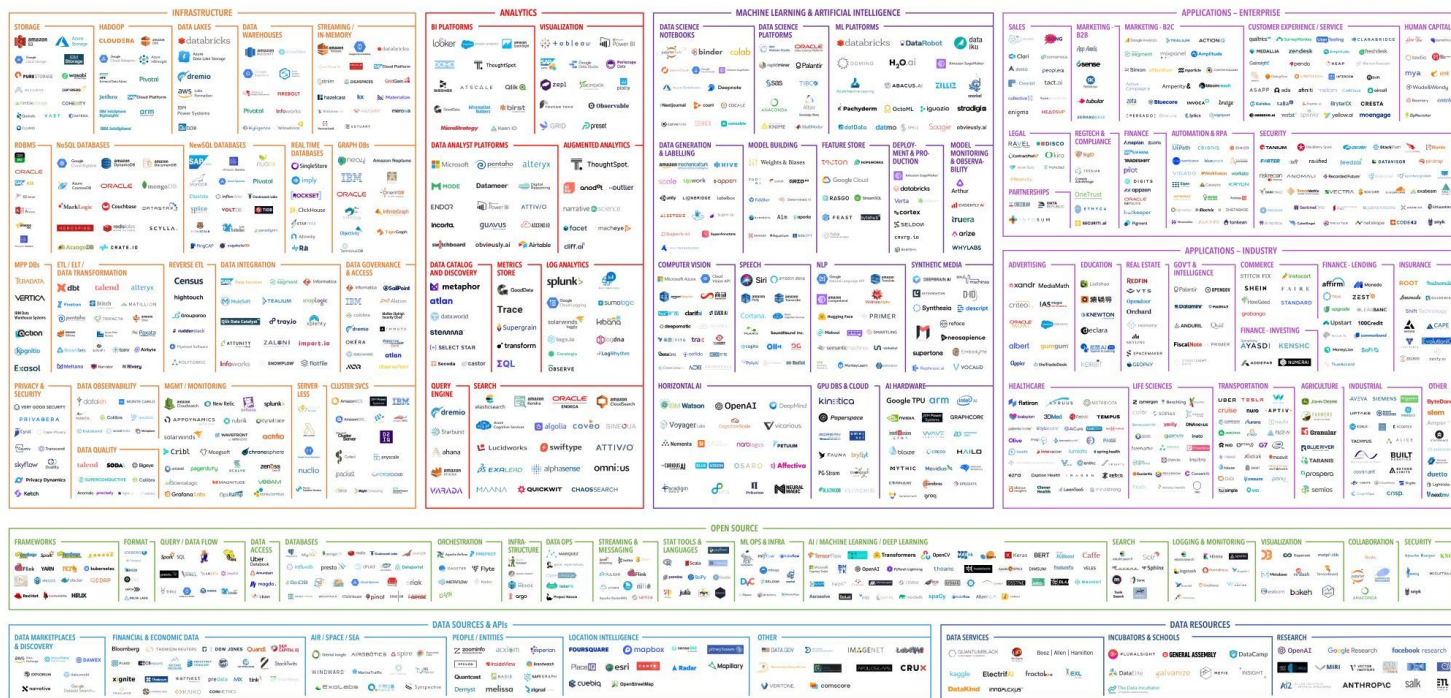


Kontekstno odvisni tokovi podatkov - komponente sistema



Medius

MACHINE LEARNING, ARTIFICIAL INTELLIGENCE, AND DATA (MAD) LANDSCAPE 2021





Enostaven primer tokov v spletni trgovini



- Spletna trgovina - nakupovanje
 - želimo uporabnikom ponuditi boljšo uporabniško izkušnjo s tem, da jim ponudimo izdelke, ki bi si jih želeli kupiti
 - podpreti želimo obstoječo komunikacijo s pripravo naročila in dostavo

Kontekstno odvisno tokovi podatkov - zahteve/želje



Medius

- enostavnost zajema podatkov
- veliko virov
- visoka prepustnost
- skalabilnost
- veliko odjemalcev
- možnost obogatitve
- NRT procesiranje
- hiter odziv na izbrane podatke/rezultate obdelav



- iz ogromne količine podatkov vzamemo samo pomembne za naš problem(kontekst)
- tokovi lahko nadomestijo trenutno komunikacijo med aplikacijami
- če načrtujemo nove aplikacije - lahko arhitekturo enostavno prilagodimo tako, da upošteva ta koncept
- uporabna vrednost je takojšnja - ne šele po končanju ML projekta, ki lahko traja več mesecev

Kontekstno odvisno tokovi podatkov - komponente



Medius

- Kafka
- Kafka streams
- Quarkus
- Apache Spark
- Elastic
- Logstash
- Kibana

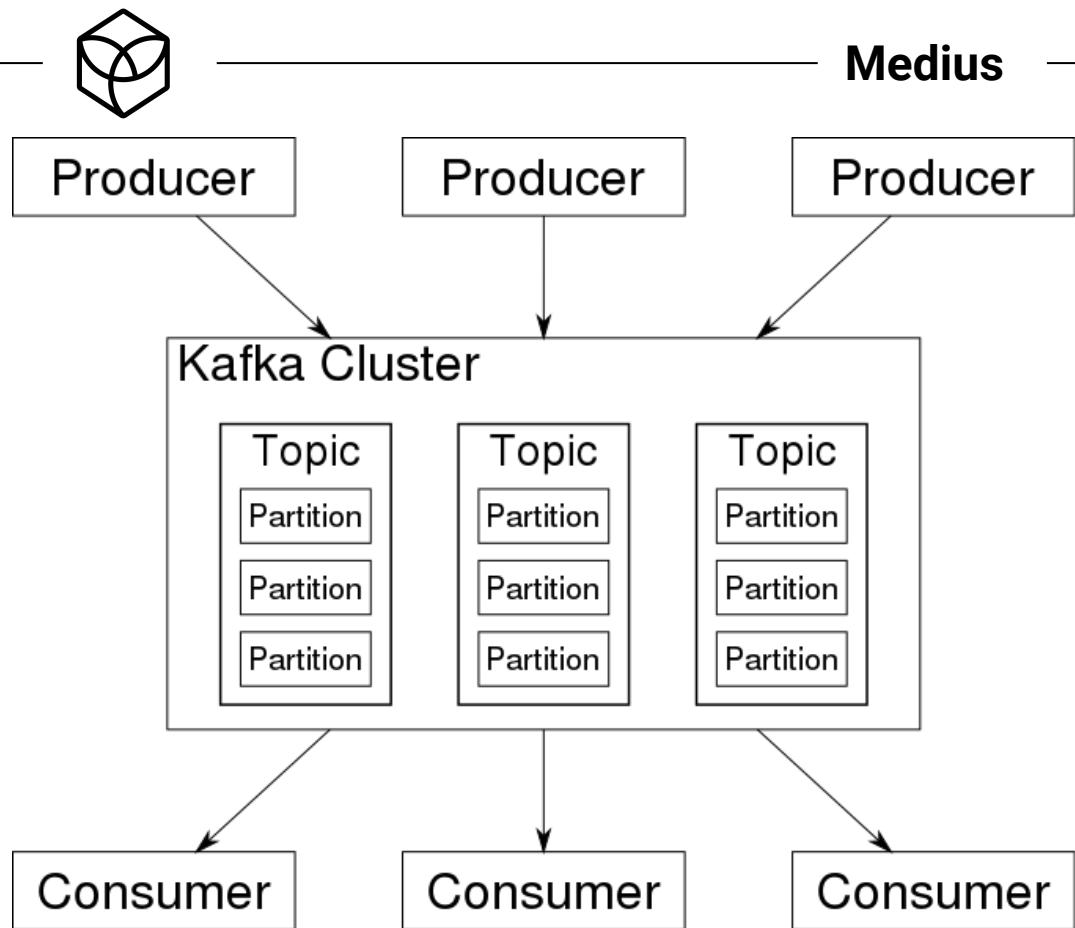




Kontekstno odvisni tokovi podatkov - Kafka

Kontekstno odvisni tokovi podatkov - Kafka

- distribuirana platforma za pretok dogodkov
- visoka prepustnost
- skalabilno (stopnjevalni sistem - scalable system)
- trdoživo hranjenje podatkov
- visok nivo dostopnosti





Kontekstno odvisni tokovi podatkov - Kafka streams



- javanska knjižnica za analiziranje in procesiranje podatkov shranjenih v kafka topicu
- enostavno delo s tokovi podatkov
- domač java API
- procesiranje zaporednih dogodkov
- windowing
- join-i
- agregacije



Kontekstno odvisni tokovi podatkov - Quarkus

Kontekstno odvisni tokovi podatkov - Quarkus



Medius

- programsko ogrodje za razvijanje Java aplikacij (Cloud Native, Container First)

What Makes Quarkus Different?



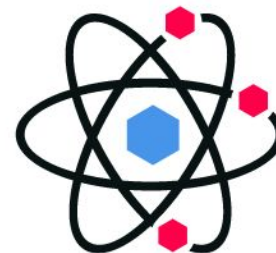
Developer Joy



Kubernetes-native



Best of Breed Libraries
and Standards



Imperative and reactive
code



Kontekstno odvisni tokovi podatkov - Apache Spark

Kontekstno odvisni tokovi podatkov - Apache Spark



Medius

- analiza podatkov in strojno učenje na tokovih
- porazdeljeno obdelovanje
- SQL analiza





Kontekstno odvisni tokovi podatkov - Elastic



- Distribuirana zastonjska iskalna in analitična platforma.
Namenjena za uporabo s katerikoli tipom podatkov (tekstovni, numerični, geoprostorskim, strukturiranim in nestrukturiranim).



elastic



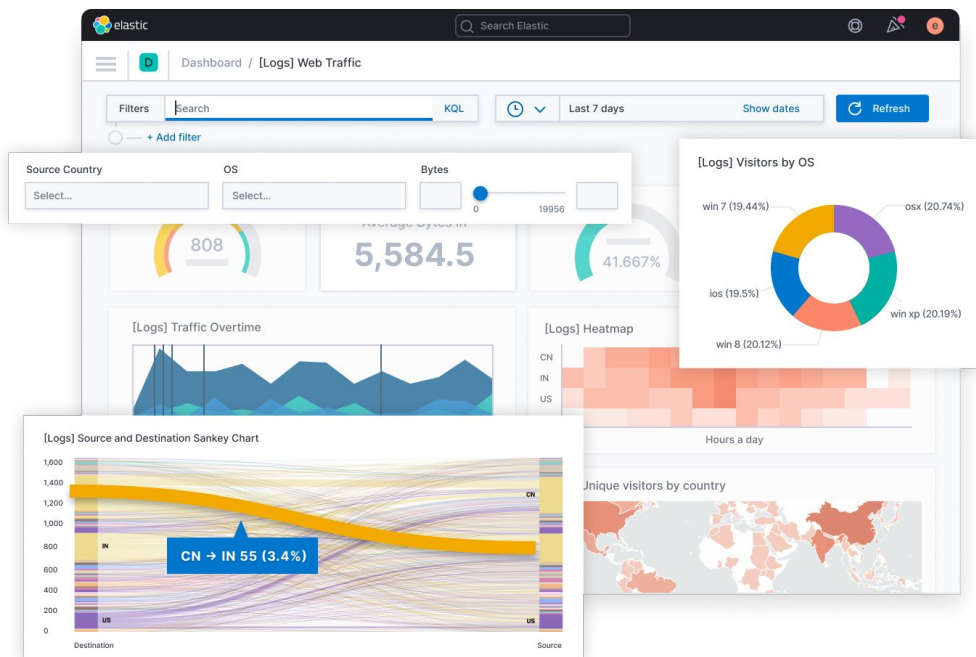
Kontekstno odvisni tokovi podatkov - Kibana

Kontekstno odvisni tokovi podatkov - Kibana



Medius

Zastonjska odprtokodna aplikacija za iskanje in vizualizacijo



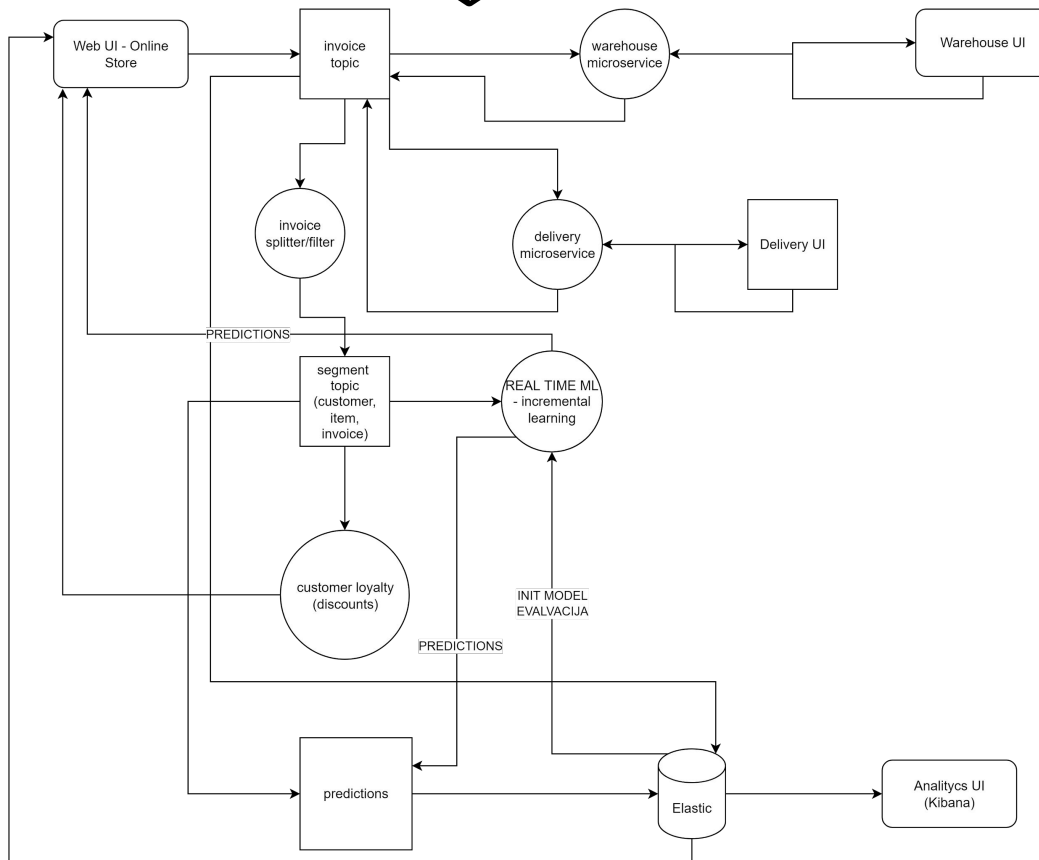


Enostaven primer tokov v spletni trgovini - arhitektura

Spletna trgovina - arhitektura



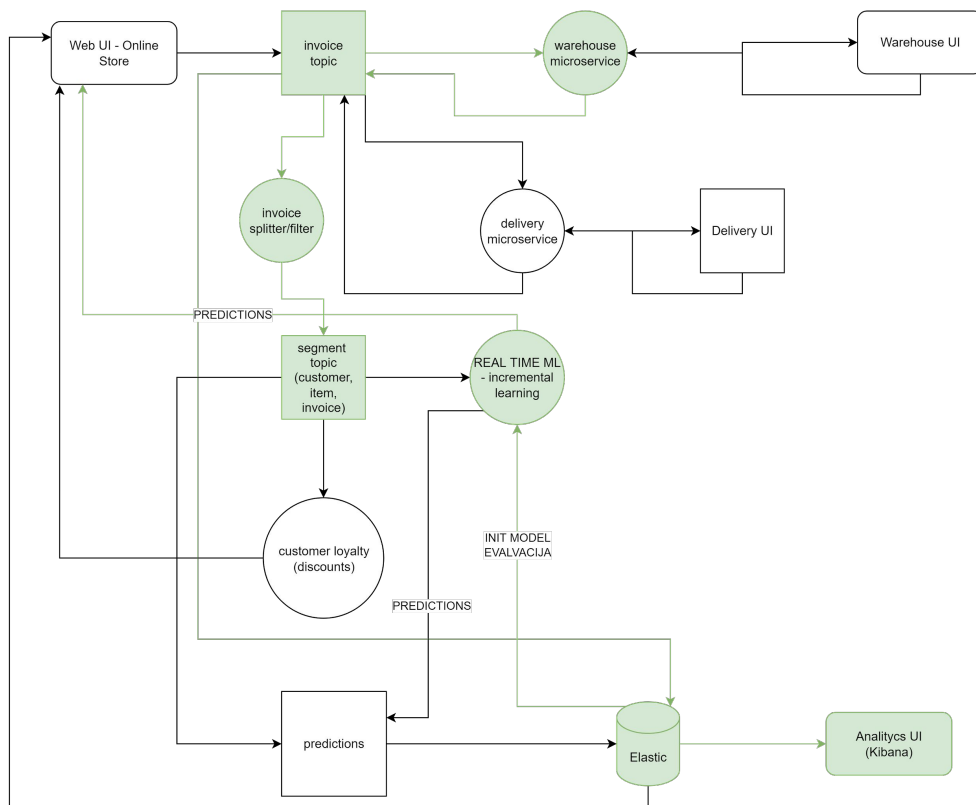
Medius



Spletna trgovina - arhitektura demo



Medius





Invoice (račun)

- customerID
- invoiceDate
- country
- invoicePrice
- invoiceNo
- status
- items:
 - description
 - quantity
 - unitPrice
 - stockCode

Spletna trgovina



Medius

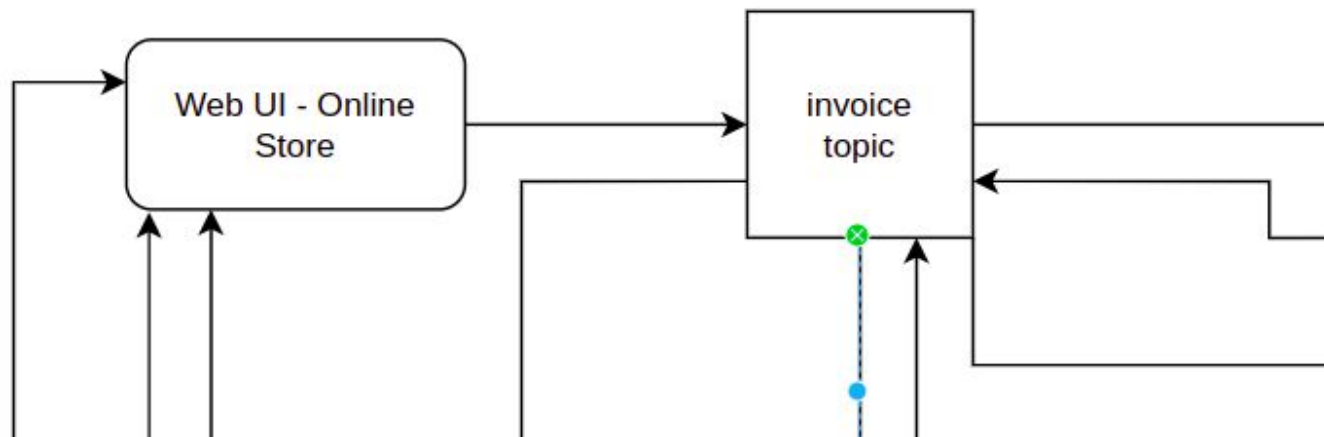
DEMO



Spletna trgovina



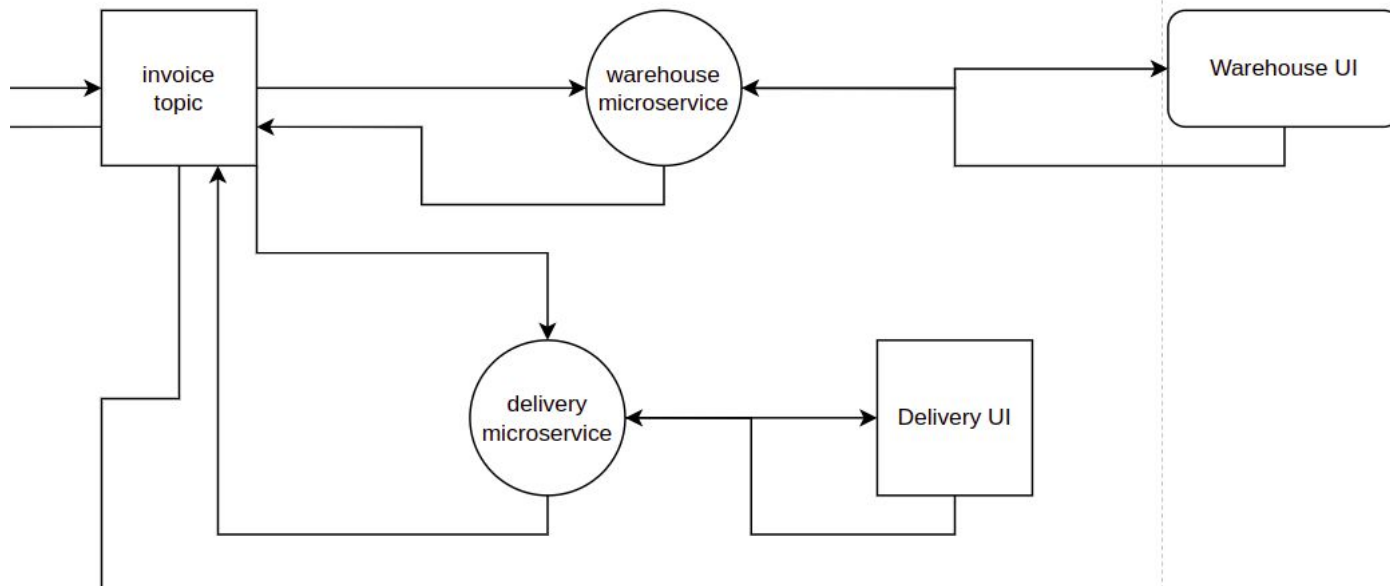
Medius



Spletna trgovina



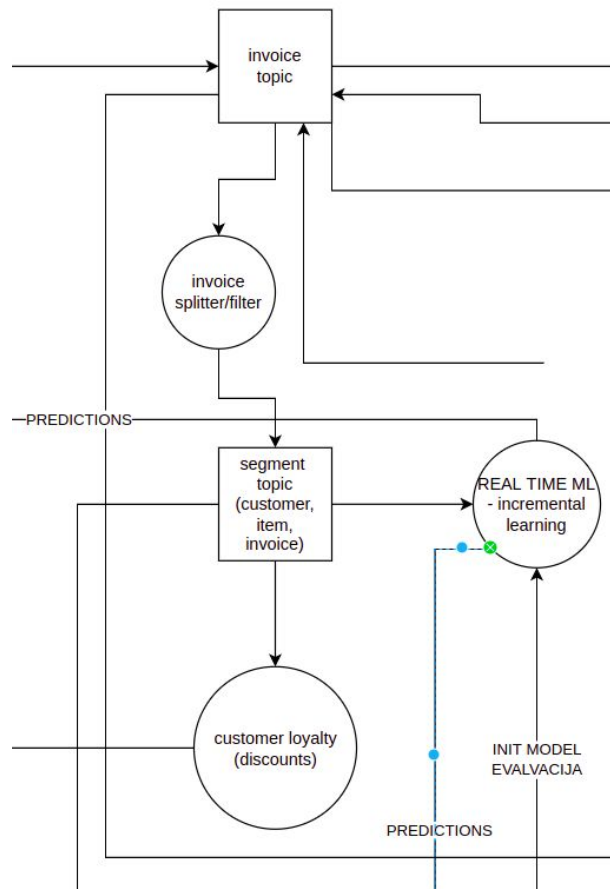
Medius



Spletna trgovina



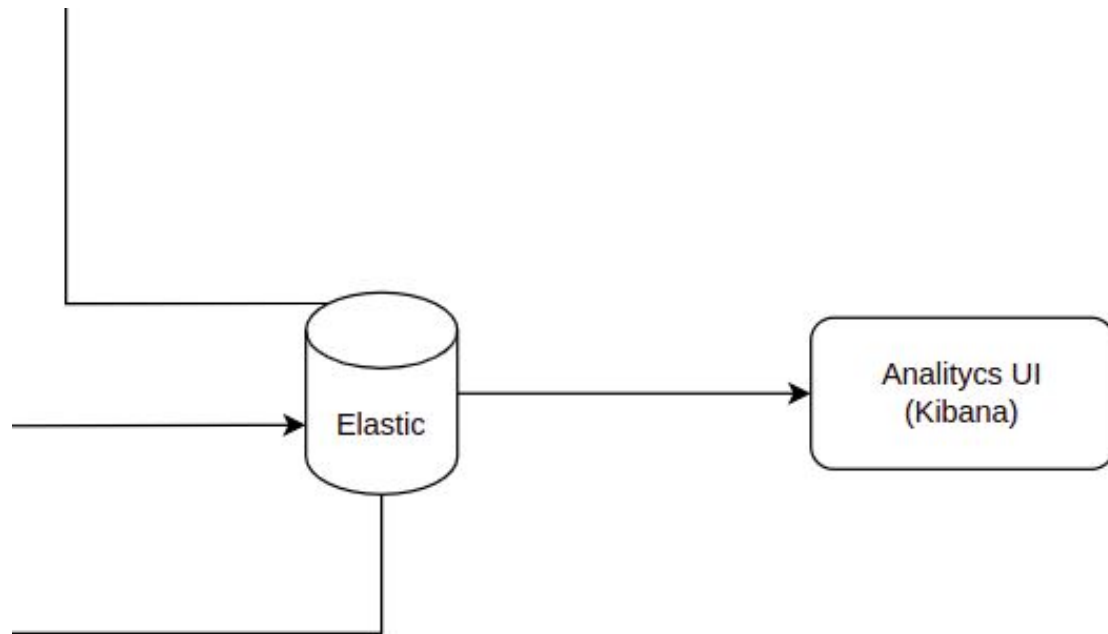
Medius



Spletna trgovina



Medius





Pregled delovanja kafke

Topics						
Topics		Partitions		Replications		Consumer Groups
Name	Count	Size	Total	Factor	In Sync	Consumer Groups
__consumer_offsets	≈ 12244	1.572 MB	50	1	1	
invoices	≈ 53507	115.979 MB	1	1	1	invoices_group_logstash Lag: 0 ▼ 3
makeit-depo-depo-invoices-changelog	≈ 18351	968.82 KB	1	1	1	
segments	≈ 385517	38.628 MB	1	1	1	



Warehouse

```
@Produces
public Topology depoStream() {
    StreamsBuilder builder = new StreamsBuilder();
    ObjectMapperSerde<Invoice> invoiceSerde = new ObjectMapperSerde<>(Invoice.class);

    Map<String, String> changeLogConfig = new HashMap<>();

    StoreBuilder<KeyValueStore<String, Invoice>> stateStore = Stores.keyValueStoreBuilder(
        Stores.persistentKeyValueStore(streamConfig.storeName()),
        Serdes.String(),
        invoiceSerde
    )
        .withLoggingEnabled(changeLogConfig)
        .withCachingEnabled();

    final Topology topology = builder.build();
    topology.addSource(streamConfig.sourceName(), Serdes.String().deserializer(), invoiceSerde.deserializer(), streamConfig.inTopic());
    topology.addProcessor(streamConfig.processorName(), DepoInvoiceProcessor::new, streamConfig.sourceName());
    topology.addStateStore(stateStore, streamConfig.processorName());
    log.info("Stream started");
    return topology;
}
```



Invoice splitter

```
@Produces
public Topology buildStream()
{
    StreamsBuilder builder = new StreamsBuilder();
    ObjectMapperSerde<Invoice> invoiceSerde = new ObjectMapperSerde<>(Invoice.class);
    ObjectMapperSerde<Segment> segmentSerde = new ObjectMapperSerde<>(Segment.class);

    builder.stream(streamConfig.inTopic(), Consumed.with(Serdes.String(), invoiceSerde)) KStream<String, Invoice>
        .filter((k, v) -> v.getStatus() == InvoiceStatus.PENDING)
        .flatMapValues(v -> {
            List<Segment> segments = new ArrayList<>();
            for (var i : v.getItems())
            {
                Segment segment = new Segment();
                segment.setInvoiceId(v.getInvoiceNo());
                segment.setItemId(i.getStockCode());
                segment.setCustomerId(v.getCustomerID());
                segments.add(segment);
            }
            return segments;
        }) KStream<String, Segment>
        .to(streamConfig.outTopic(), Produced.with(Serdes.String(), segmentSerde));
    return builder.build();
}
```

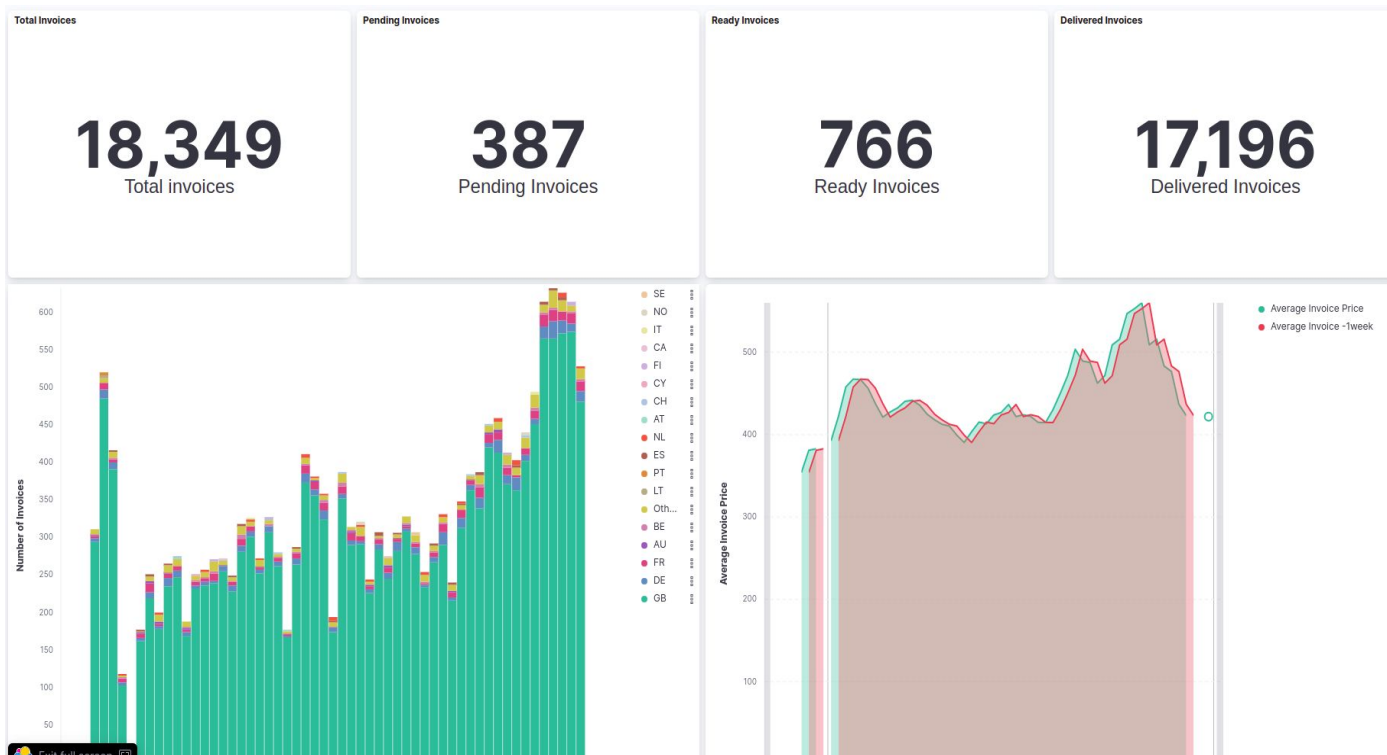


Primer dokumenta na elasticsearchu

Actions	Field	Value
...	_id	581586
...	_index	invoices
...	_score	-
...	@timestamp	May 30, 2022 @ 07:55:03.716
...	@version	1
...	country	GB
...	customerID	13113
...	event.original	> { "invoiceNo": "581586", "customerID": "13113", "country": "GB", "invoiceDate": "2021-12-11T12:49:00", "items": [{ "stockCode": "22061", "description": "LARGE CAKE STAND HANGING STRAWBERRY", "quantity": 8, "unitPrice": 2.95, "stockCode": "23275", "description": "SET OF 3 HANGING OWLS OLLIE BEAK", "quantity": 24, "unitPrice": 1.25, "stockCode": "21217", "description": "RED RETROSPOT ROUND CAKE TINS", "quantity": 24, "unitPrice": 8.95, "stockCode": "20685", "description": "DOORMAT RED RETROSPOT", "quantity": 10, "unitPrice": 7.08 }] }
...	invoiceDate	Dec 11, 2021 @ 13:49:00.000
...	invoiceNo	581586
...	invoicePrice	339.2
...	items.description	LARGE CAKE STAND HANGING STRAWBERRY, SET OF 3 HANGING OWLS OLLIE BEAK, RED RETROSPOT ROUND CAKE TINS, DOORMAT RED RETROSPOT
...	items.quantity	8, 24, 24, 10
...	items.stockCode	22061, 23275, 21217, 20685
...	items.unitPrice	2.95, 1.25, 8.95, 7.08
...	status	DELIVERED



Primer kibana dashboarda



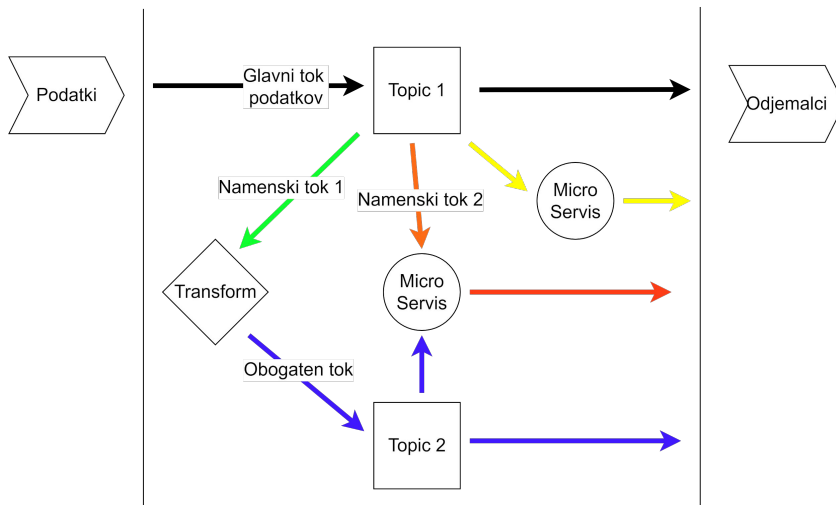


KODA:

<https://github.com/dsenica/medius-makeit-2022>



Kontekstno odvisni tokovi podatkov - Zaključek





Hvala za pozornost!



Vprašanja?