
Анализ публикуемых новостей

Голунова Анастасия, DE

Цель

Создать ETL-процесс формирования витрины данных для анализа публикаций новостей из источников:

- <https://lenta.ru/rss/>
- <https://www.vedomosti.ru/rss/news>
- <https://tass.ru/rss/v2.xml>

Инструменты



1. Для хранения и обработки данных по новостям выбрана СУБД Postgres, т.к. данные не большие (не более 1000 новостей в день со всех источников)
2. В качестве оркестратора ETL процесса выбран Apache Airflow
3. Таски дагов написаны при помощи Python оператора



План реализации

Два дага



```
graph TD; A[Два дага] --> B[Инициализирующий]; A --> C[Инкрементальный]; B --> B1[Загрузка всех доступных данных из трех источников в raw слой]; B1 --> B2[Создание таблиц core слоя и их заполнение]; B2 --> B3[Создание витрины данных на основании данных core слоя]; C --> C1[Загрузка данных за текущий день из трех источников в raw слой]; C1 --> C2[Загрузка данных в таблицы core слоя]; C2 --> C3[Обновление витрины данных на основании данных core слоя];
```

Инициализирующий

Расписание: нет (по триггеру)

Загрузка **всех доступных данных** из трех источников в raw слой



Создание таблиц core слоя и их заполнение



Создание витрины данных на основании данных core слоя

Инкрементальный

Расписание: ежедневно в 23:55

Загрузка **данных за текущий день** из трех источников в raw слой

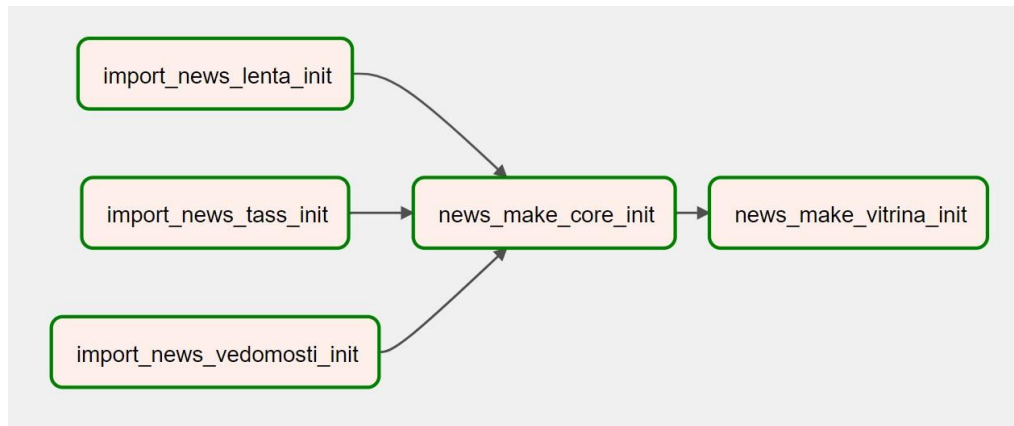


Загрузка данных в таблицы core слоя



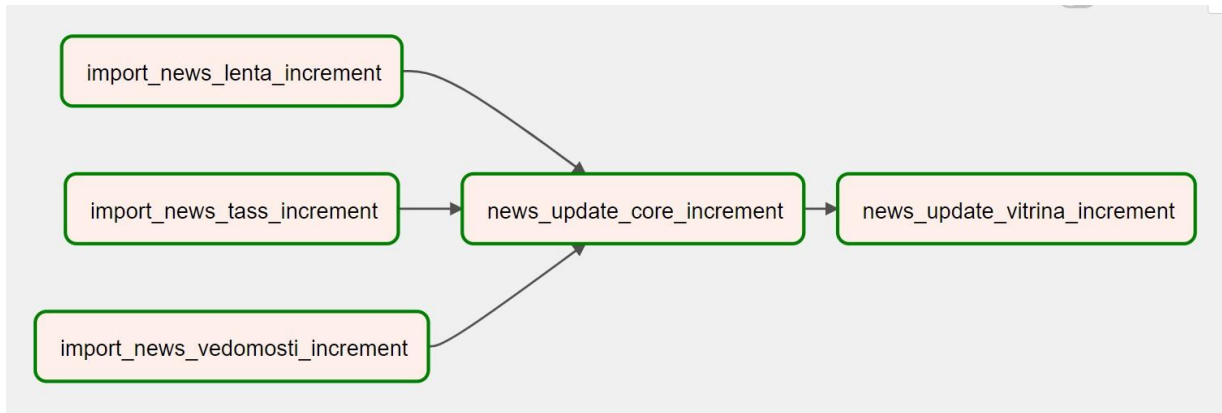
Обновление витрины данных на основании данных core слоя

Инициализирующий даг



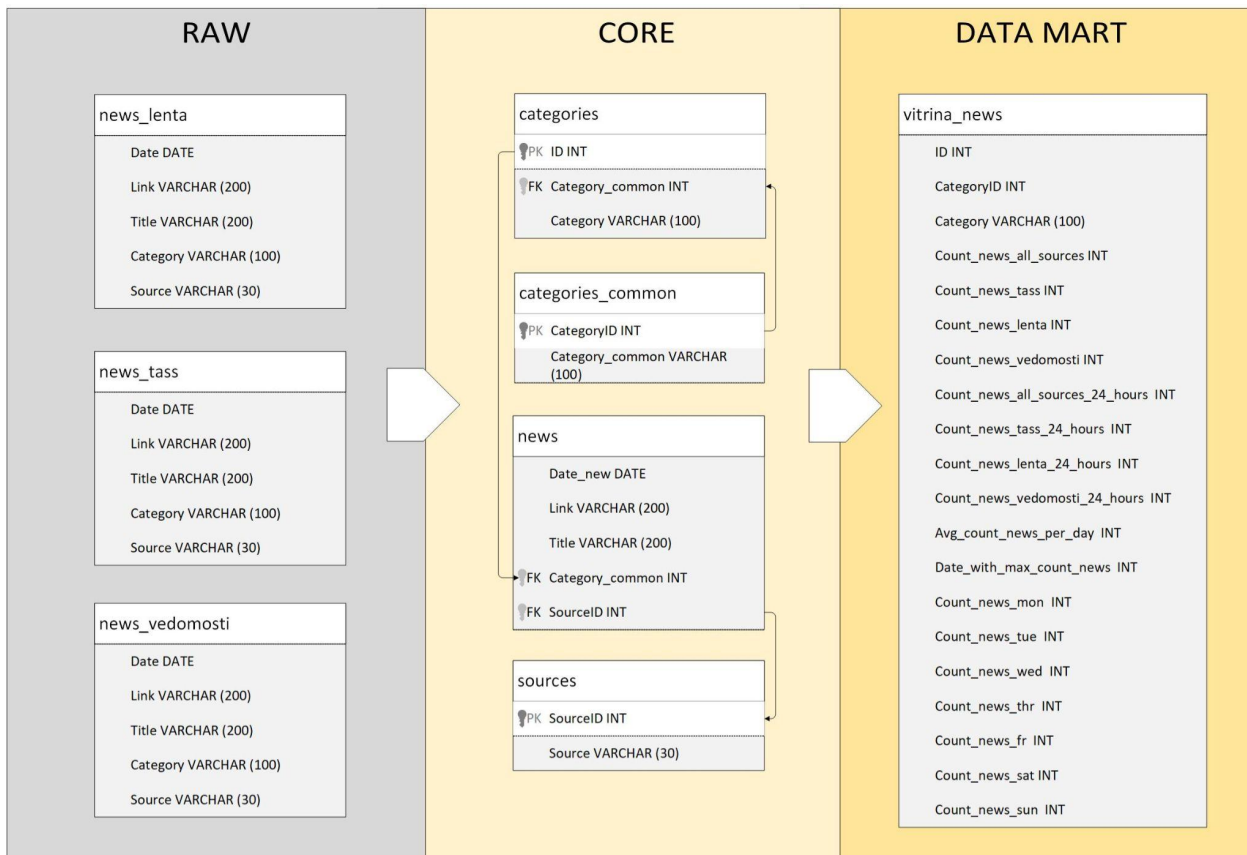
1. `import_news_lenta_init`, `import_news_tass_init`, `import_news_vedomosti_init` - загружают все доступные данные из трех источников, формируют датафреймы и в виде датафреймов отправляют в Postgres с методом `replace`.
2. `news_make_core_init` - создает таблицы `core` слоя и заполняет их данными с `raw` слоя
3. `news_make_vitrina_init` - создает витрину данных на основе данных `core` слоя

Инкрементальный даг



1. `import_news_lenta_increment`, `import_news_tass_increment`, `import_news_vedomosti_increment` - загружают данные за текущий день из трех источников, формируют датафреймы и в виде датафреймов отправляют в Postgres с методом `append`. Время для дага выбрано, так как источник Тасс передает только данные за текущий день.
2. `news_update_core_increment` - актуализирует таблицы core слоя данными с raw слоя
3. `news_update_vitrina_increment` - актуализирует витрину данных на основе данных core слоя

ER cxema



ER схема. Комментарии

1. Схема core слоя - звезда
2. В разных источниках категории новостей отличаются, поэтому создана доп. таблица с унифицированными категориями. НО так как результаты витрины зависят от выбранных общих категорий - лучше их уточнять у аналитиков.

```
SET Category_common =  
CASE  
  WHEN Category LIKE ANY (ARRAY['%экономика%', '%изн[е]с']) THEN 'Экономика/Бизнес'  
  WHEN Category LIKE ANY (ARRAY['%инвестиции%', '%инансы%']) THEN 'Финансы/ Инвестиции'  
  WHEN Category LIKE ANY (ARRAY['%реда обитания%', '%бщество%', '%жизни%', '%енности%', '%абота%']) THEN 'Общество'  
  WHEN Category LIKE ANY (ARRAY['%стран%', '%оссия%', '%9 паралл%']) THEN 'Россия'  
  WHEN Category LIKE ANY (ARRAY['%осква%', '%еверо-Запад%']) THEN 'Москва'  
  WHEN Category LIKE ANY (ARRAY['%олити%', '%иловые%', '%рмия%']) THEN 'Политика'  
  WHEN Category LIKE ANY (ARRAY['%ехнолог%', '%аука%']) THEN 'Технологии'  
  WHEN Category LIKE ANY (ARRAY['%еждународ%', '%мир%']) THEN 'Мир'  
  WHEN Category like '%ультур%' THEN 'Культура'  
  WHEN Category LIKE ANY (ARRAY['%нтернет%', '%овости партнер%']) THEN 'Интернет/ СМИ'  
  WHEN Category like '%утешествия%' THEN 'Путешествия'  
  WHEN Category like '%роисшествия%' THEN 'Происшествия'  
  WHEN Category like '%едвижимость%' THEN 'Недвижимость'  
  WHEN Category like '%порт%' THEN 'Спорт'  
  ELSE 'Другое'  
END
```


Результат

Разработан ETL-процесс формирования витрины данных для анализа публикаций новостей из трех источников.

Скрин

123	doc category	123 count	123 cou	123 cc	123 cour	123 coun	123 cou	123 cour	123 cou	123 avg_count_n	date_wi	123 cc	123 count	123 count	123 cou	123 cou	123 cou	123 a
1	Экономика/Бизнес	80	18	30	32	68	18	30	20	1,6326530612	2023-09-19	8	62	0	0	0	8	0
2	Финансы/ Инвестиции	15	[NULL]	[NULL]	15	14	[NULL]	[NULL]	14	1	2023-09-18	12	3	0	0	0	0	0
3	Общество	154	34	32	88	108	34	32	42	1,6923076923	2023-09-19	28	80	0	0	0	26	0
4	Россия	30	[NULL]	30	[NULL]	30	[NULL]	30	[NULL]	1,0344827586	2023-09-19	0	30	0	0	0	0	0
5	Москва	8	8	[NULL]	[NULL]	8	8	[NULL]	[NULL]	1	2023-09-19	0	8	0	0	0	0	0
6	Политика	259	35	13	211	131	35	13	83	1,8633093525	2023-09-19	66	71	0	0	0	52	0
7	Технологии	20	[NULL]	11	9	17	[NULL]	11	6	1	2023-09-19	3	15	0	0	0	1	0
8	Мир	64	22	42	[NULL]	64	22	42	[NULL]	1,0158730159	2023-09-19	0	64	0	0	0	0	0
9	Культура	18	4	14	[NULL]	18	4	14	[NULL]	2	2023-09-19	0	18	0	0	0	0	0
10	Путешествия	8	[NULL]	8	[NULL]	8	[NULL]	8	[NULL]	1	2023-09-19	0	8	0	0	0	0	0
12	Спорт	22	8	14	[NULL]	22	8	14	[NULL]	2	2023-09-19	0	22	0	0	0	0	0
13	Другое	36	1	25	10	34	1	25	8	1	2023-09-19	7	28	0	0	0	0	0
14	Интернет/ СМИ	10	[NULL]	10	[NULL]	10	[NULL]	10	[NULL]	1	2023-09-19	0	10	0	0	0	0	0
15	Происшествия	7	7	[NULL]	[NULL]	7	7	[NULL]	[NULL]	1	2023-09-19	0	7	0	0	0	0	0