



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ

Data Analysis National Olympiad - DANO

# КАТЕГОРИАЛЬНЫЕ ПЕРЕМЕННЫЕ

Елена Семерикова,  
Валерия Стефаненко

Москва, 2021





# КОДИРОВКА ДАННЫХ

— это такое представление переменных, которое разбивает их на категории с общими признаками или идентифицирует уникальность объекта.



# РОССТАТ:

Средняя начисленная  
зарплата  
за октябрь 2019 г., руб.



Отношение  
заработной платы  
женщин к заработной  
плате мужчин,  
в процентах

Добыча полезных ископаемых	<b>61 246</b>	<b>80 663</b>	<b>75,9%</b>	
Деятельность в области информации и связи	<b>53 887</b>	<b>78 980</b>	<b>68,2%</b>	
Деятельность профессиональная, научная и техническая	<b>55 187</b>	<b>77 719</b>	<b>71,0%</b>	
Строительство	<b>53 296</b>	<b>61 263</b>	<b>87,0%</b>	
Транспортировка и хранение	<b>41 431</b>	<b>55 943</b>	<b>74,1%</b>	
Обрабатывающие производства	<b>40 669</b>	<b>52 500</b>	<b>77,5%</b>	
Торговля оптовая и розничная; ремонт автотранспортных средств и мотоциклов	<b>41 160</b>	<b>52 122</b>	<b>79,0%</b>	
Деятельность в области культуры, спорта, организации досуга и развлечений	<b>35 079</b>	<b>47 265</b>	<b>74,2%</b>	

# ЯНДЕКС.GO: данные по поездкам в такси



Час пик	Дамми переменная
Нет	0
Да	1
Нет	0
Нет	0
Нет	0
Нет	0
Да	1
Нет	0
Нет	0
Нет	0
Нет	0
Нет	0
Нет	0
Да	1

## ДАММИ-ПЕРЕМЕННАЯ, или ФИКТИВНАЯ ПЕРЕМЕННАЯ

разбивает переменные на 2 категории «0» и «1» с общими признаками для удобства их сравнения и исследования уникальных для таких категорий признаков



# ЦИАН: Поиск жилья



циан

Ещё фильтры

Тип сделки  От собственника

До метро: Неважно | Пешком | **Транспортом** | Не более 45 минут

Площадь, м<sup>2</sup>: Общая от до | Кухня от до | Жилая от до

Планировка: Неважно | Смежная | Изолированная |  Схема планировки

Высота потолков: Неважно | От 2,5 м | От 2,7 м | От 3 м | От 3,5 м | От 4 м

Санузел: Неважно | Совмещённый | Раздельный |  Два и более

Балкон/Лоджия: Неважно | **Балкон** | Лоджия

Ремонт: Неважно | Без ремонта | Косметический | Евроремонт | Дизайнерский

[Сбросить фильтры](#) [Показать объекты](#)

№	Доступность метро	On_foot	Наличие балкона	Balcony
1	<b>пешком</b>	<b>1</b>	Да	1
2	транспортом	0	Да	1
3	<b>транспортом</b>	<b>0</b>	Нет	0
4	пешком	1	Да	1
5	пешком	1	Нет	0
6	транспортом	0	<b>Нет</b>	<b>0</b>
7	транспортом	0	<b>Да</b>	<b>1</b>



# ЯНДЕКС.GO: данные по поездкам в такси



УТРО = «1»  
ДЕНЬ = «2»  
ВЕЧЕР = «3»

Время суток	Period	Period_1	Period_2	Period_3
Вечер	3	0	0	1
день	2	0	1	0
утро	1	1	0	0
утро	1	1	0	0
утро	1	1	0	0
день	2	0	1	0
день	2	0	1	0
утро	1	1	0	0
утро	1	1	0	0
день	2	0	1	0
день	2	0	1	0
утро	1	1	0	0
день	2	0	1	0
вечер	3	0	0	1
вечер	3	0	0	1
день	2	0	1	0

# РАЗМЕР ПРЕДПРИЯТИЙ

в рамках государственной  
программы поддержки

от количественных к качественным  
переменным

Размер предприятия	Обороты, руб	Категория	Business_size
до 15 человек	120 млн	микropредприятия	1
до 100 человек	800 млн	малые	2
до 250 человек	2 млрд	средние	3

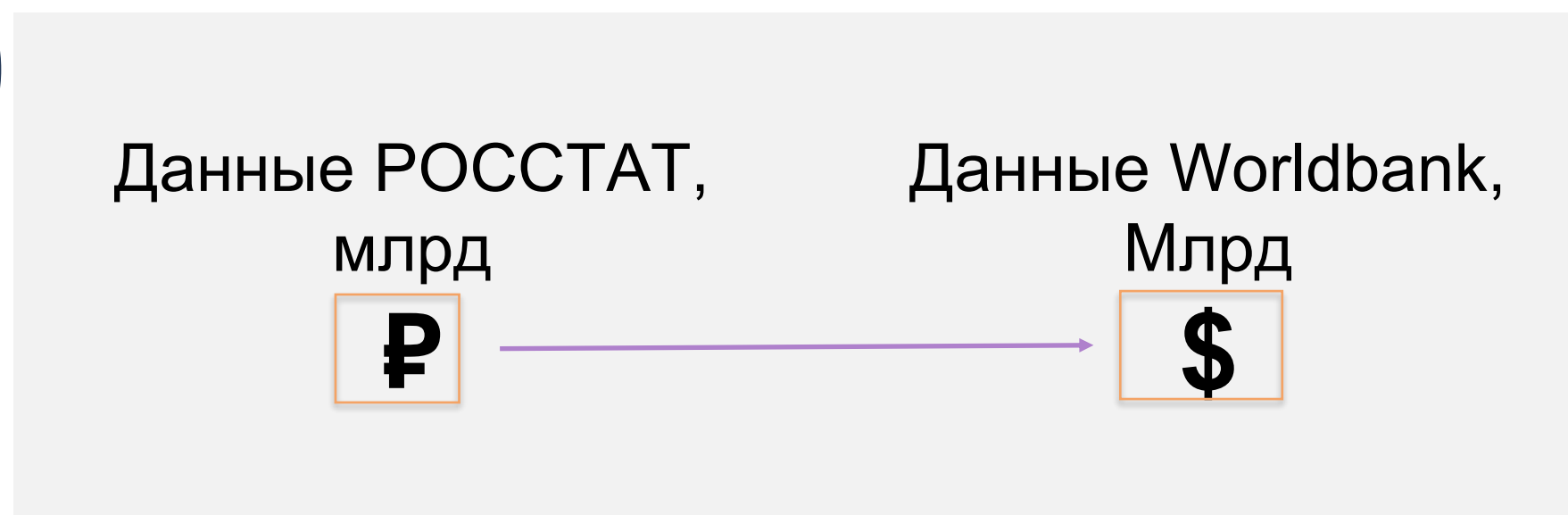






# ВВП СТРАНЫ 2019

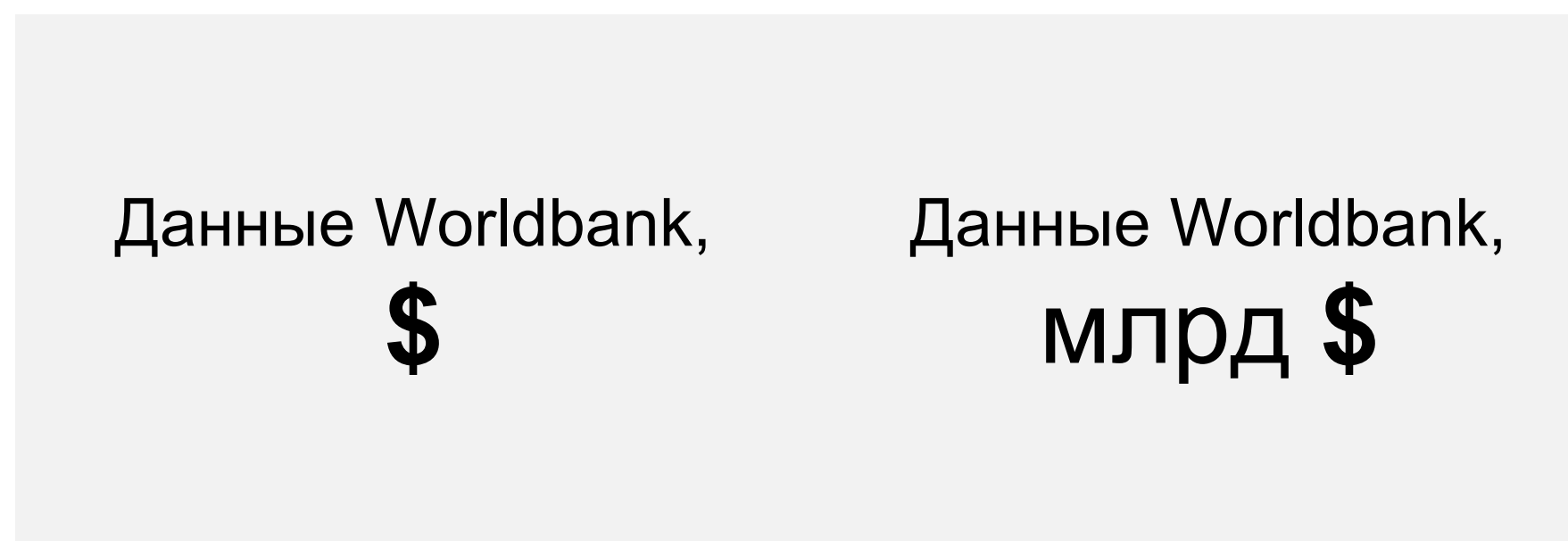
перекодировка данных



Сопоставимость данных

110 046,1

1 687,4



Удобство представления данных

1,68745E+12

1 687,4





# ЗАДАНИЕ 1

для самопроверки

Месяц	Стоимость картофеля за кг, руб
Декабрь	20
Январь	21,5
Февраль	21,5
Март	22
Апрель	35
Май	40
Июнь	40
Июль	25
Август	17
Сентябрь	18
Октябрь	18
Ноябрь	18,5

Допустим, вы исследуете стоимость овощей. Для этого приходите на рынок каждый месяц и записываете их среднюю стоимость. Собранные вами данные представлены в таблице справа.

1. Поясните, как закодировать данные на основании времени года?
2. Во сколько раз картофель стоит дороже весной, чем осенью?



# ЗАДАНИЕ 2

## для самопроверки

В Москве можно увидеть довольно много многоэтажек. Однако в последнее время также пользуется популярностью и малоэтажное строительство. Вы располагаете данными о количестве этажей в нескольких домах, а также о технологии их строительства.

1. Подумайте, сколько категорий можно придумать для группировки данных, по какому принципу?
2. Закодируйте данные по принципу технологии.
3. Сформируйте дамми-переменную, исходя из классификации:
  - а. малоэтажный- 1-4 этажа (с учетом мансарды);
  - б. многоэтажный - 5 этажей и выше

Дом	Кол-во этажей	Технология
1	12	Панельный
2	4	Монолитный
3	7	Панельный
4	3	Кирпичный
5	9	Панельный
6	5	Монолитный
7	3	Кирпичный
8	2	Монолитный



# ЗАДАНИЕ 3

## для самопроверки

Справа представлены данные прогноза погоды для США и Канады. Однако в странах применяются разные метрические системы для измерения температур: в США - градус Фаренгейта, а в Канаде – градус Цельсия.

Ответьте на вопрос, на сколько средняя температура за неделю в Цельсиях в США отличается от средней температуры за неделю в Канаде?

Дата	Температура в США, °F	Температура в Канаде, °C
26/07/2021	87	27
27/07/2021	90	26
28/07/2021	86	24
29/07/2021	87	25
30/07/2021	86	25
31/07/2021	86	25
01/08/2021	84	23





---

# ПОЛЕЗНЫЕ ССЫЛКИ

для изучения

Борис Демешев. Дамми (фиктивные) переменные. Разные зависимости для подвыборок  
<https://www.youtube.com/watch?v=fxIVaU9OMn8>



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ