

# Data Science UA. Лекция 3.

## Exploratory Data Analysis

Январь 2017

## Пару слов о домашнем задании

1. Где оно? :)
2. Не затягивайте до последнего дня.
3. Держите форму вопроса!
4. Качество данных - важно. Качество тех из них, по которым дается ответ на поставленный вопрос - критично.
5. Не выводите полные списки/DataFrame - пользуйтесь `.head()` и лимитами в списках.
6. Отделяйте исследование данных от непосредственно домашней работы.
7. Пользуйтесь функционалом языка.
8. Приводите финальный результат.

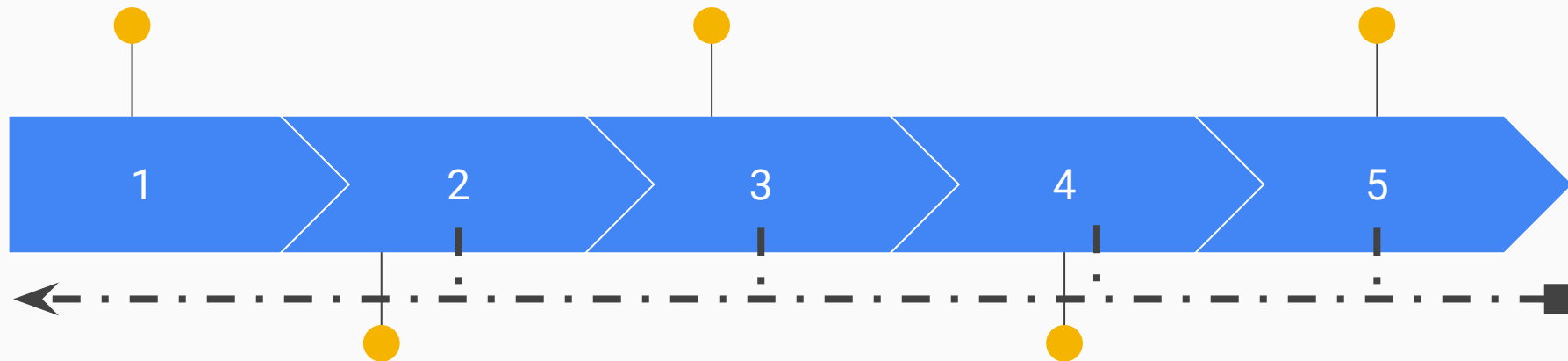
Не всегда понятно что мы ищем... пока мы это не найдем.



Формулировка  
задачи

Анализ данных

Презентация,  
использование и  
контроль модели



Сбор и очистка  
данных

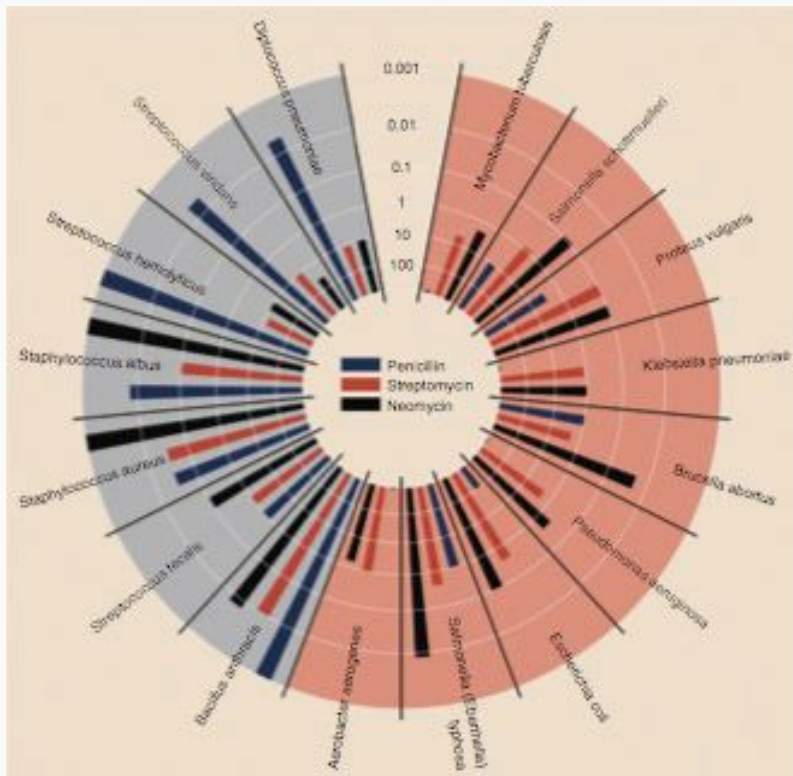
Построение модели  
данных

# Антибиотики

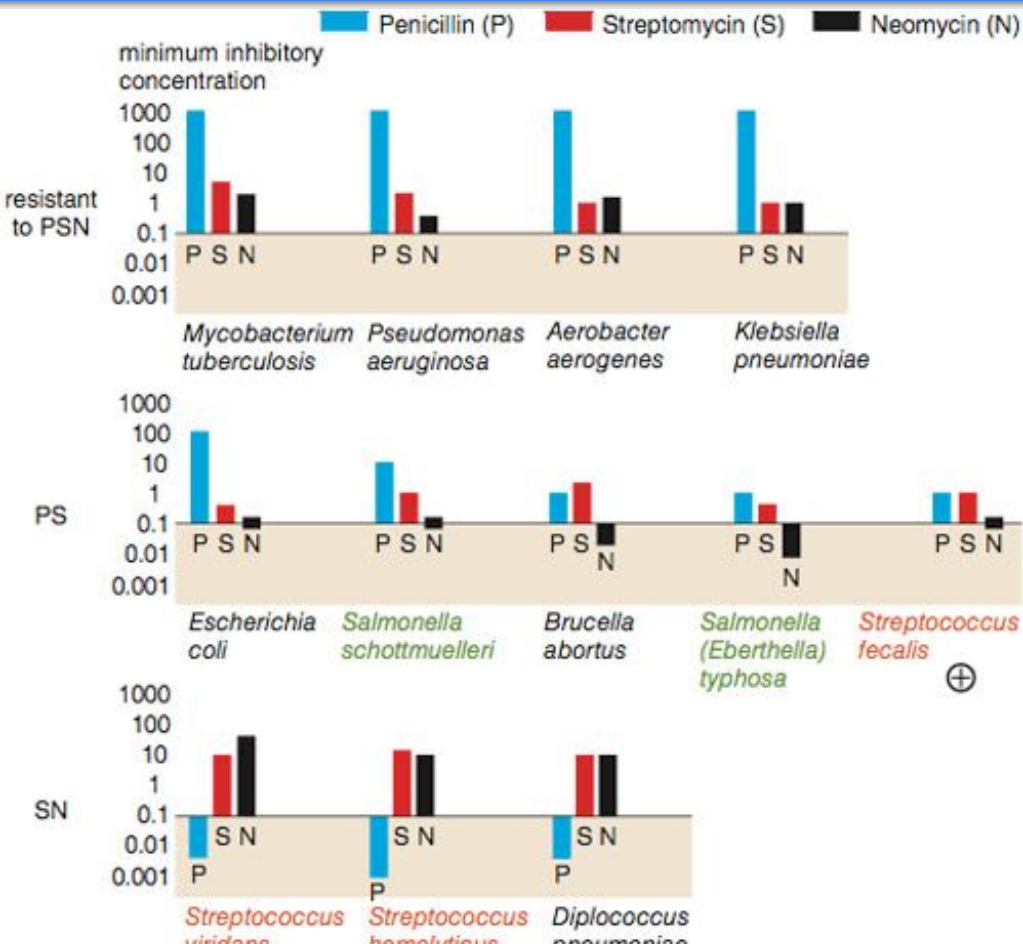
Уилл Буртин, исследование 1951 года, незадаанные вопросы

<http://www.americanscientist.org/issues/pub/thats-funny>

# Эффективность антибиотиков



Bacteria	Penicillin	Antibiotic Streptomycin	Neomycin	Gram stain
<i>Aerobacter aerogenes</i>	870	1	1.6	—
<i>Brucella abortus</i>	1	2	0.02	—
<i>Bacillus anthracis</i>	0.001	0.01	0.007	+
<i>Diplococcus pneumoniae</i>	0.005	11	10	+
<i>Escherichia coli</i>	100	0.4	0.1	—
<i>Klebsiella pneumoniae</i>	850	1.2	1	—
<i>Mycobacterium tuberculosis</i>	800	5	2	—
<i>Proteus vulgaris</i>	3	0.1	0.1	—
<i>Pseudomonas aeruginosa</i>	850	2	0.4	—
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	—
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	—
<i>Staphylococcus albus</i>	0.007	0.1	0.001	+
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	+
<i>Streptococcus faecalis</i>	1	1	0.1	+
<i>Streptococcus hemolyticus</i>	0.001	14	10	+
<i>Streptococcus viridans</i>	0.005	10	40	+



Величайшая ценность изображений в том, что они заставляют видеть то, чего вы не собирались увидеть.

Джон Таки

## Цели визуализации

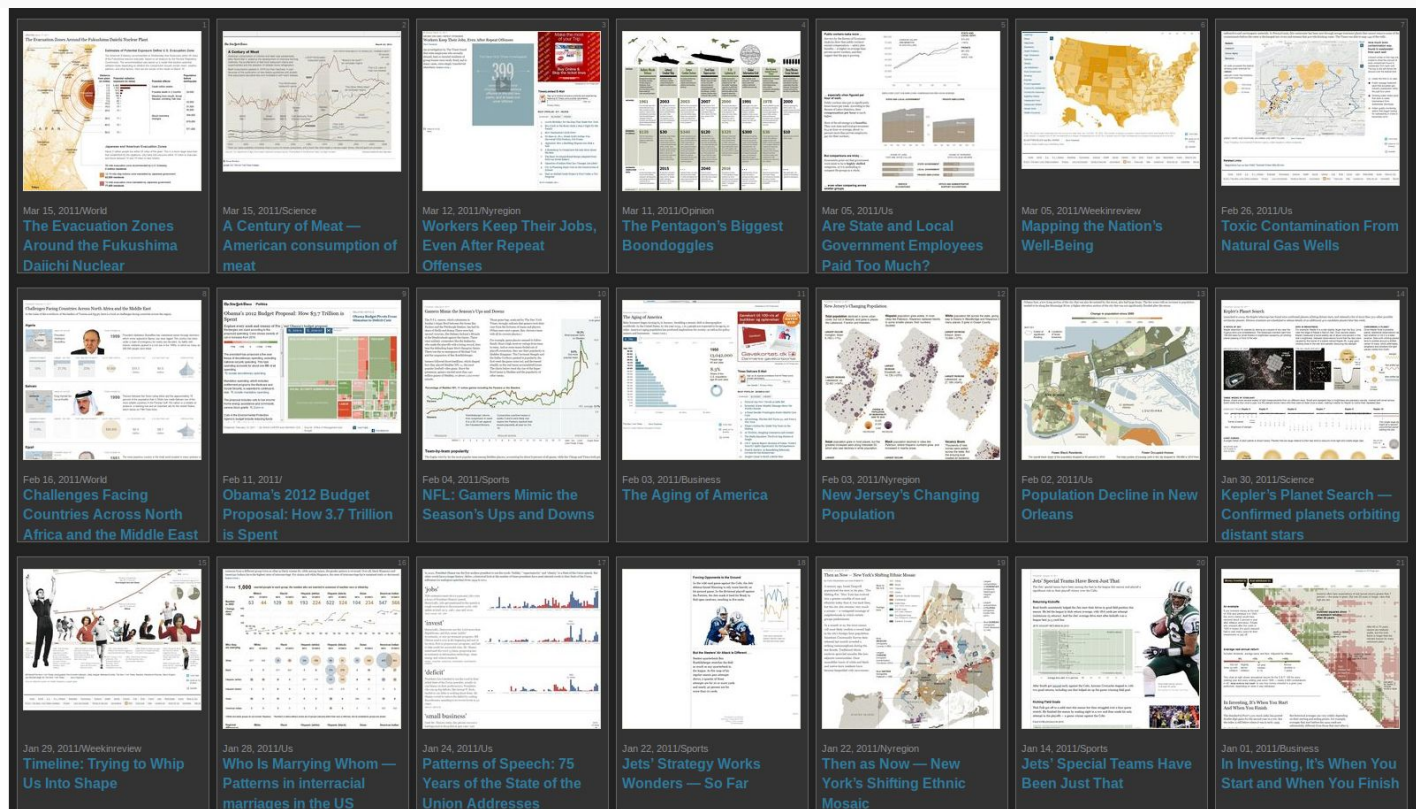
### Коммуникативная (explanatory)

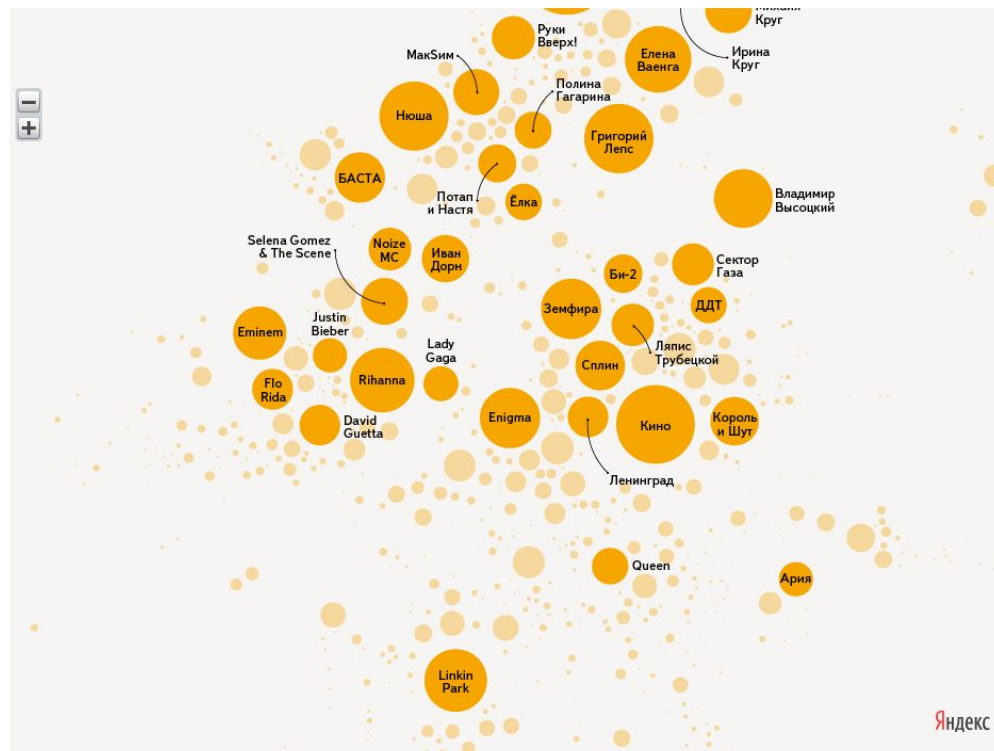
1. Представить данные и идеи
2. Донести и проинформировать
3. Поддержать и аргументировать
4. Повлиять и убедить

### Исследовательская (exploratory)

1. Исследовать данные
2. Проанализировать ситуацию
3. Определить следующие шаги
4. Вынести решение по вопросу

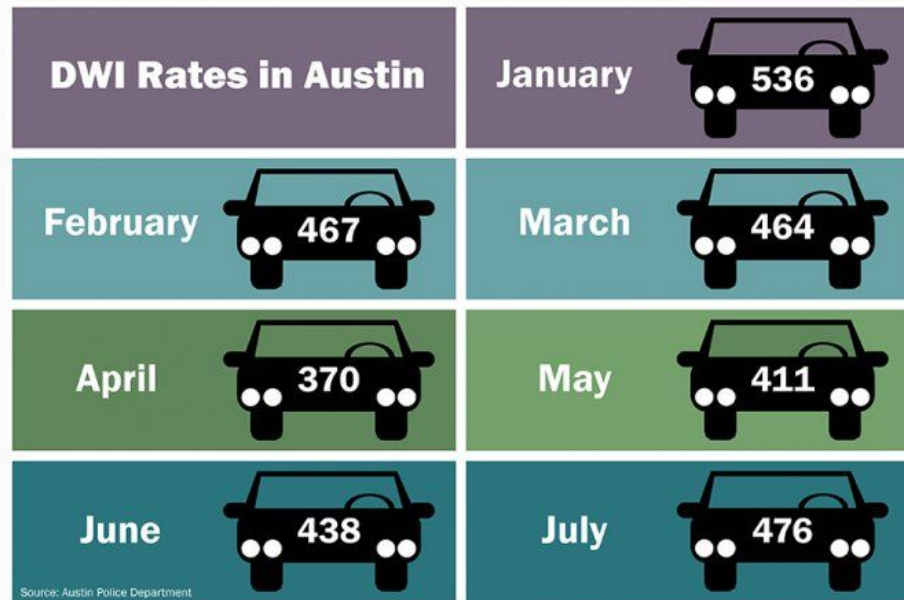
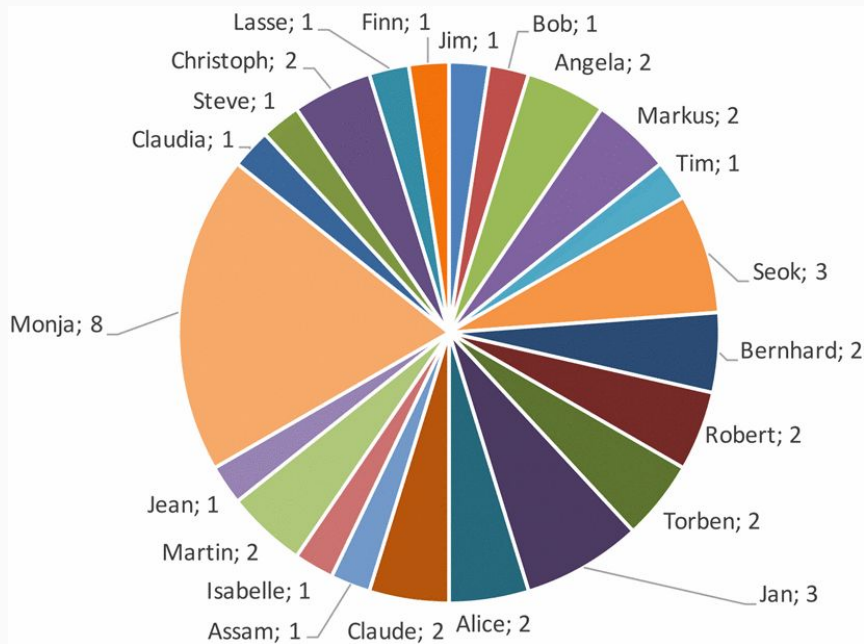






Посмотреть на карте в полный экран

# Визуализируйте (не) правильно!



## Проблемы плохой визуализации

1. Неправильный выбор типа визуализации.
2. Перегруженность материалом.
3. 3D эффекты.
4. Просто ненужные эффекты.
5. Неправильные/неэффективные цвета.
6. Слишком мало информации.
7. Информация, которую предполагается сравнивать - разнесена в пространстве.
8. Тысячи их!

Если после просмотра графика/диаграммы у вас возникает желание посмотреть исходные данные - визуализация провалилась в своей задаче.

# Визуализируйте правильно!

Ключевые моменты:

1. Графическая целостность.
2. Простота.
3. Правильная форма.
4. Правильное использование цвета.
5. Целостность повествования.

# Графическая целостность

## Weight Over Time

POUNDS

180 —

175 —

170 —

165 —

164 —

*This isn't weight. The length of a bar represents how much greater than a measurement is than the minimum, 164 pounds.*

Day 1

20

40

60

80

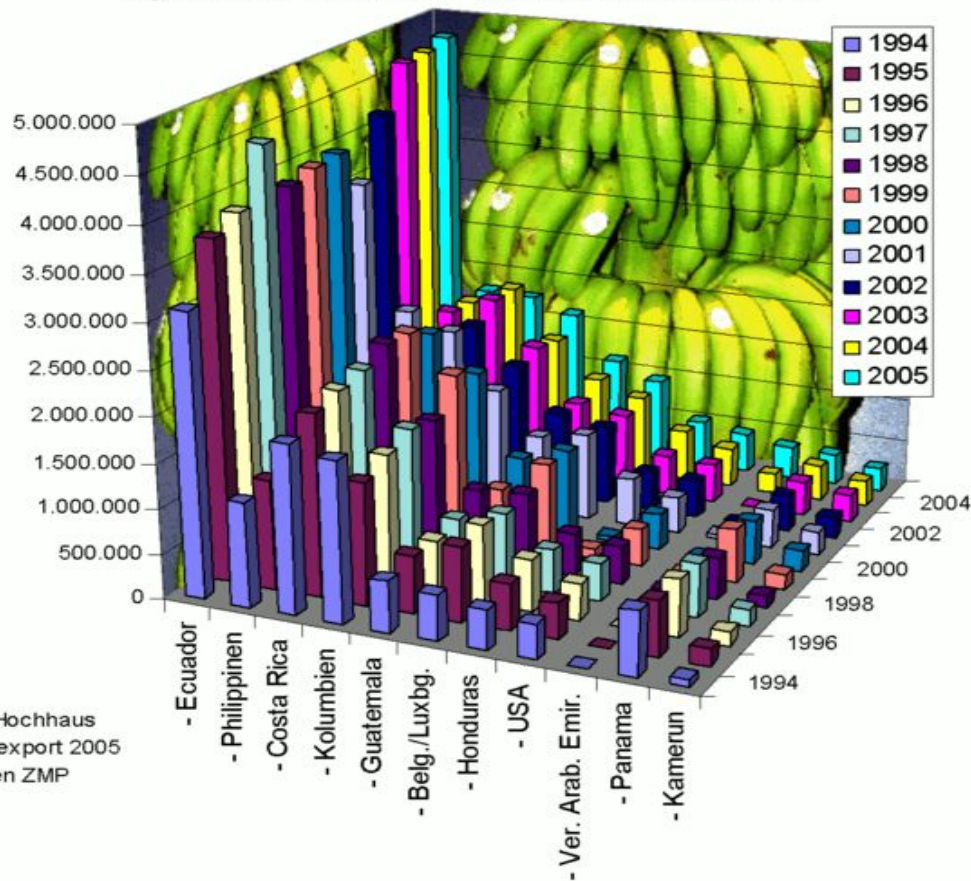
100

120

Шкалы - честная координатная сетка, честное представление данных. Равные расстояния на визуализации должны означать равную разницу в данных.

<https://flowingdata.com/2015/08/31/bar-chart-baselines-start-at-zero/>

Export von Bananen in Tonnen von 1994-2005



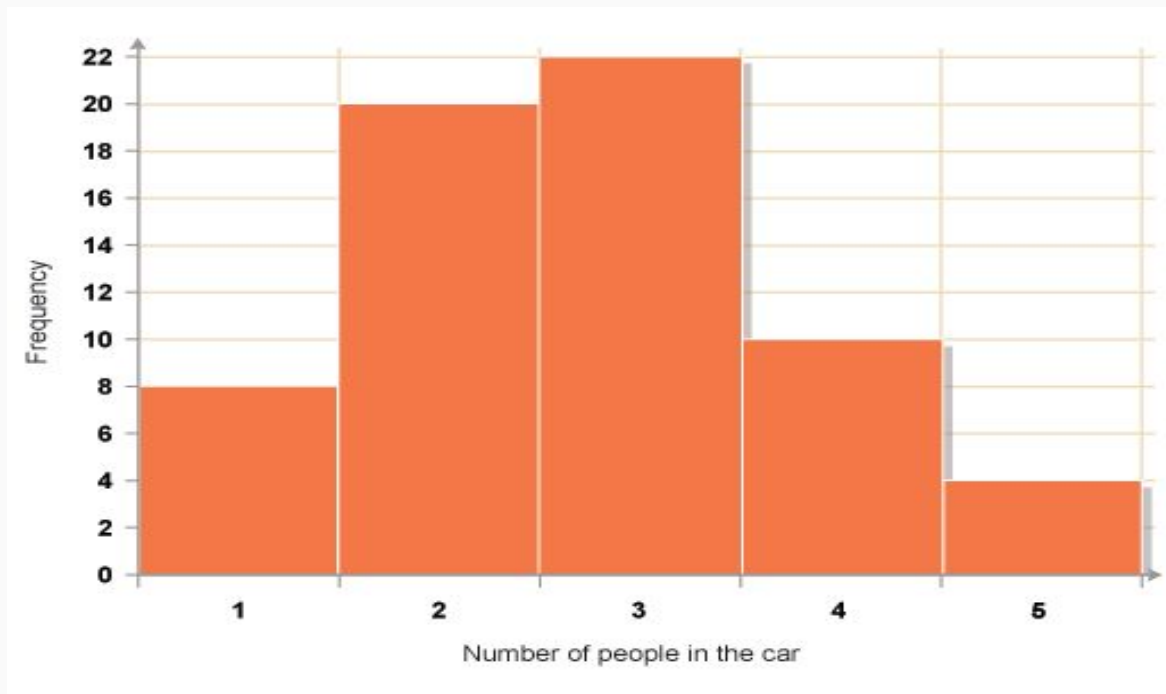
Отношение данные / чернила  
должно быть как можно  
больше. (с) Эдвард Тафти

<https://en.wikipedia.org/wiki/Chartjunk>



# Правильная форма

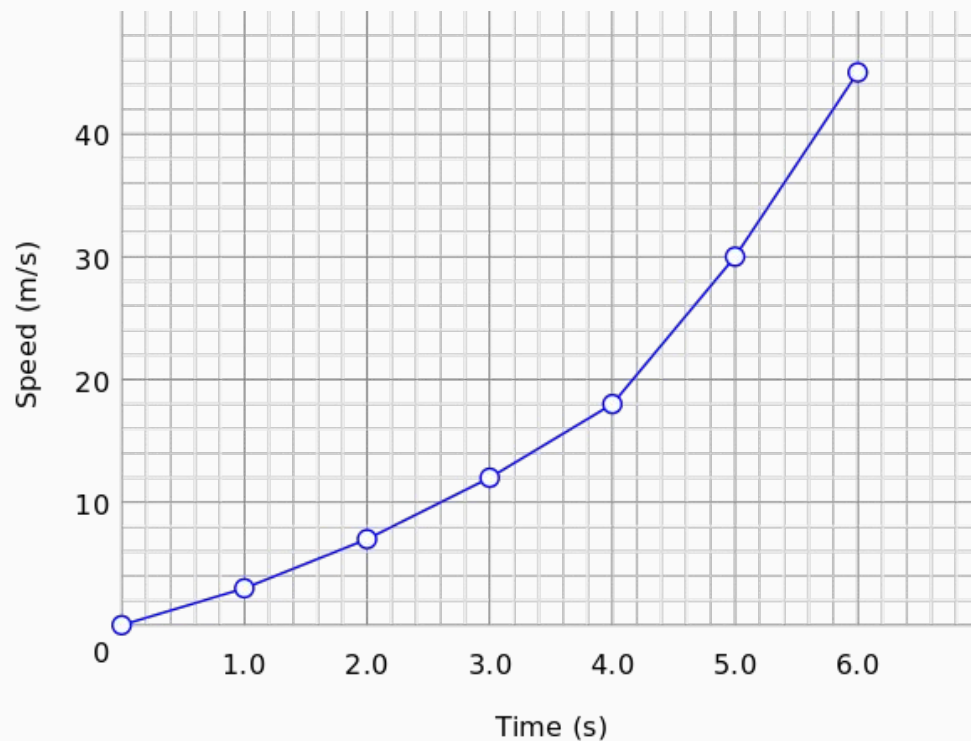
Сравнение - (столбчатая диаграмма) bar charts. Линейные графики подразумевают непрерывность





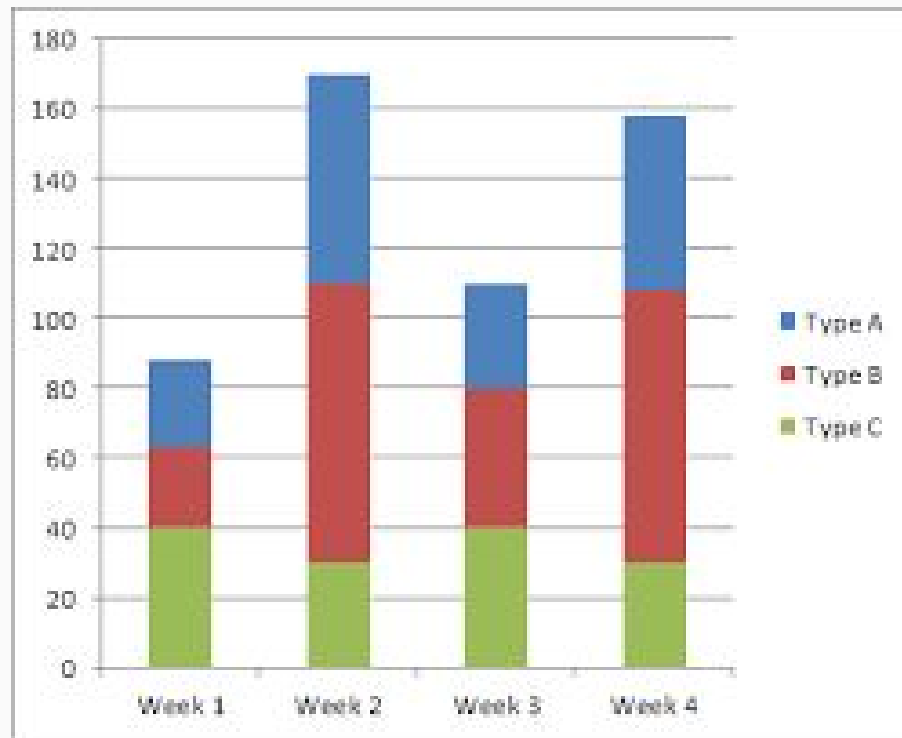
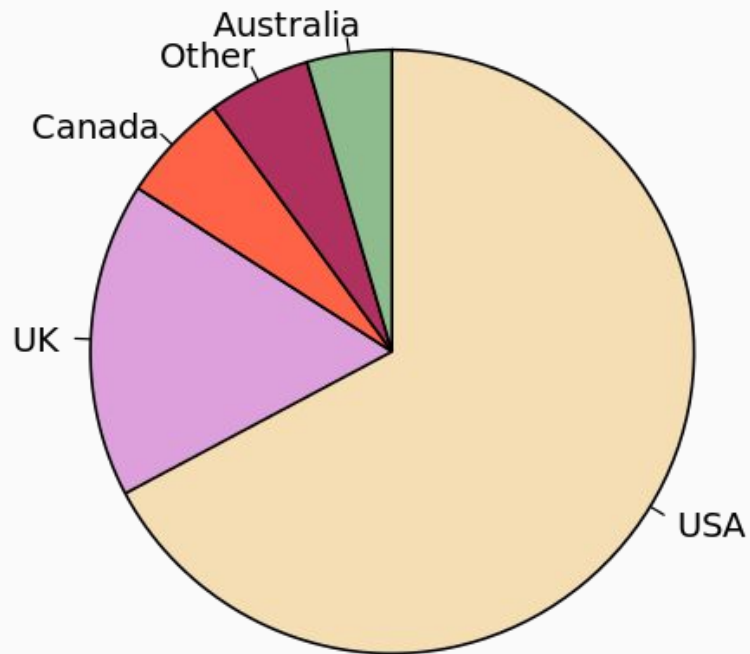
# Правильная форма

Тренды во времени, непрерывные данные - линейные графики.

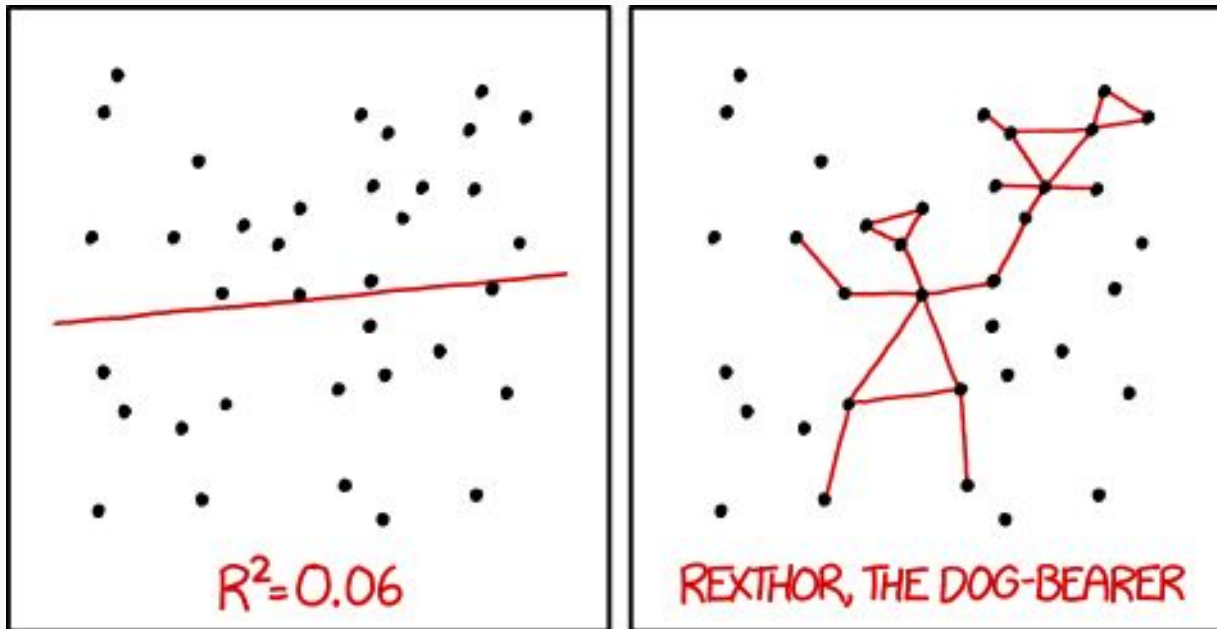


# Правильная форма

Пропорции - пироги :). Круговая диаграмма. Колонки как вариант.



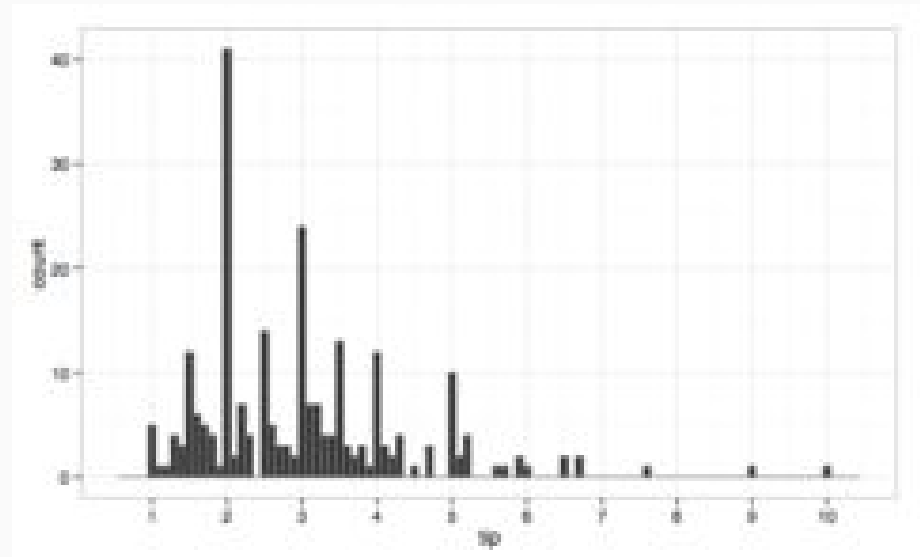
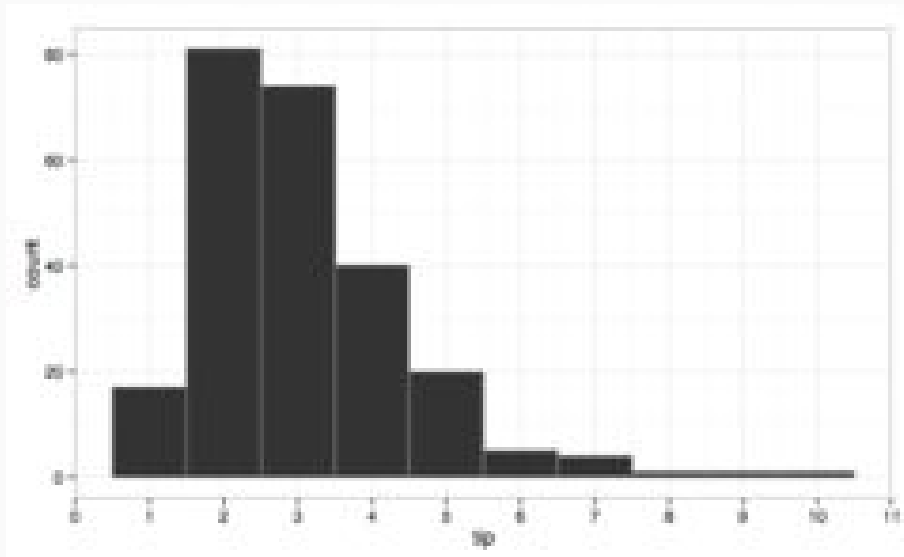
Корреляции - scatter plot (диаграмма рассеяния).



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER  
TO GUESS THE DIRECTION OF THE CORRELATION FROM THE  
SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

# Правильная форма

Распределение - гистограмма.

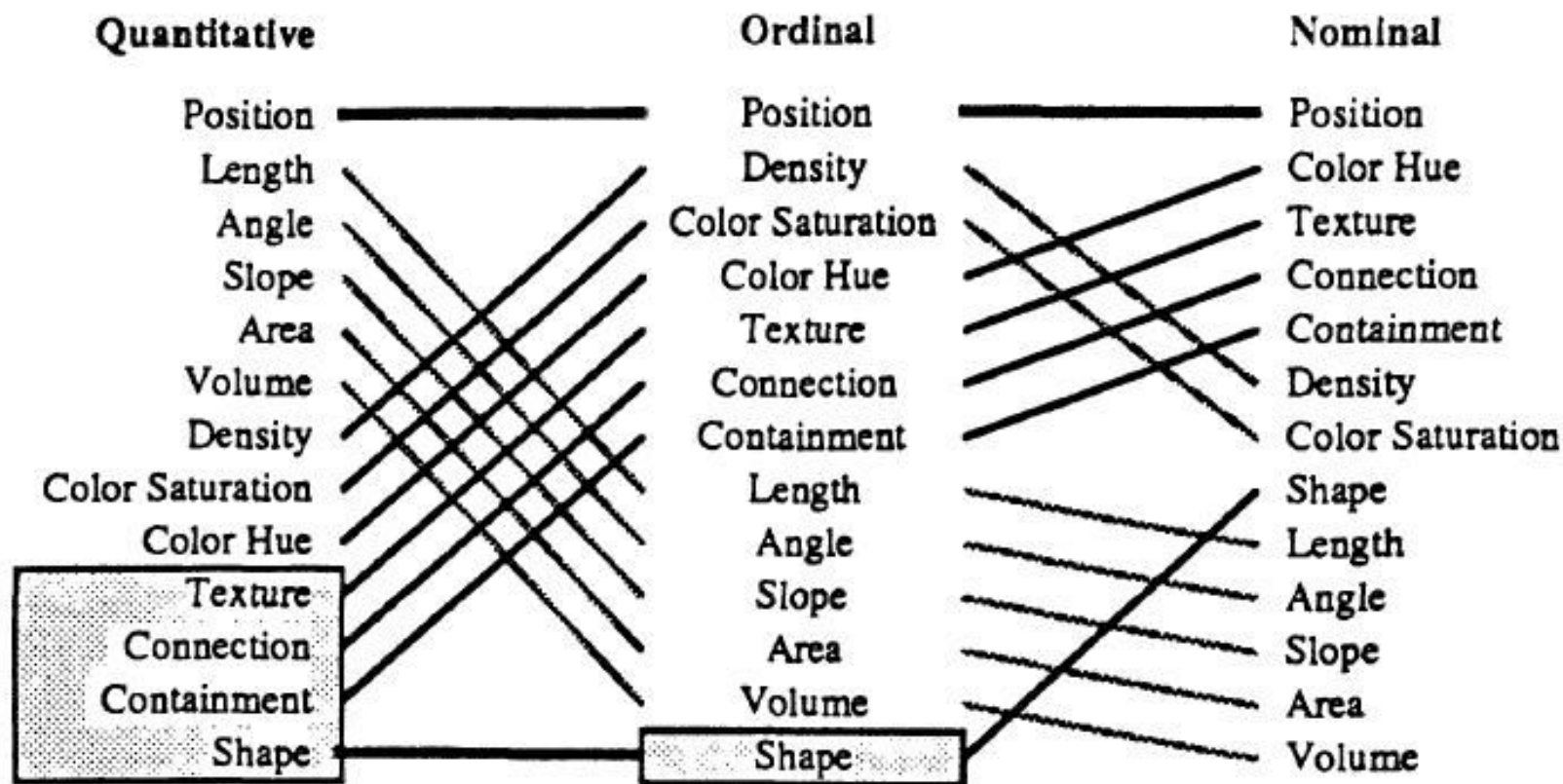


	1.before	2.after	3.in_1_week
0	10	15	10
1	20	10	15
2	30	35	40
3	40	40	35

# Эффективность восприятия

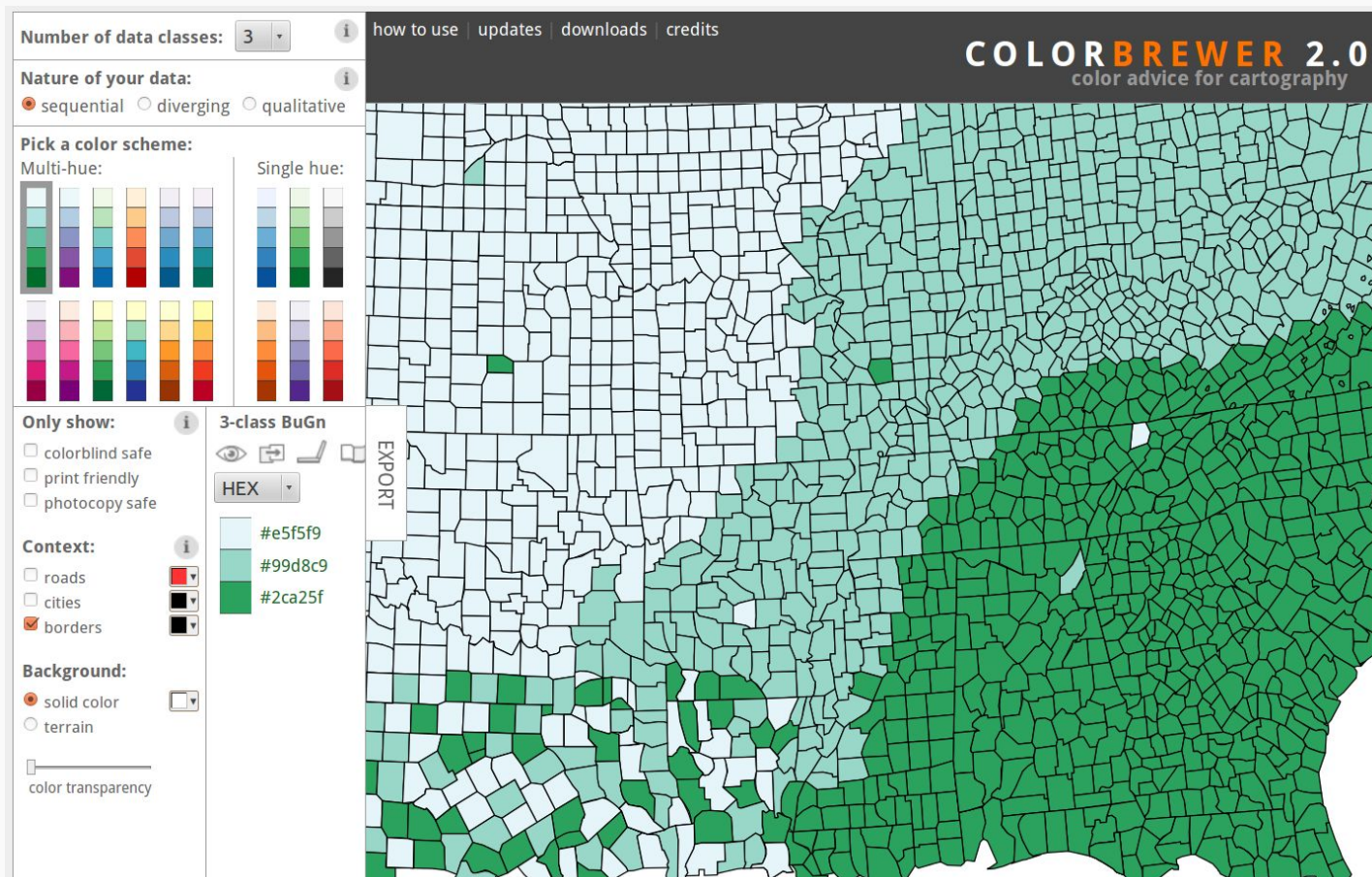
Bertin's Original Visual Variables	
<b>Position</b> changes in the x, y location	
<b>Size</b> change in length, area or repetition	
<b>Shape</b> infinite number of shapes	
<b>Value</b> changes from light to dark	
<b>Colour</b> changes in hue at a given value	
<b>Orientation</b> changes in alignment	
<b>Texture</b> variation in 'grain'	

# Эффективность восприятия





# Правильное использование цвета





# Рекомендуемые ресурсы

## 1. Сайты

- a. <https://flowingdata.com/>
- b. <http://colorbrewer2.org/#type=sequential&scheme=BuGn&n=3>
- c. <http://www.datavis.ca/milestones/index.php?group=1950%2B>
- d. [http://www.infovis-wiki.net/index.php?title=Visual\\_Variables](http://www.infovis-wiki.net/index.php?title=Visual_Variables)

## 2. Книги

- a. Эдвард Тафти
  - i. Visual Display of quantitative information
  - ii. Envisioning information
  - iii. Visual explanations
- b. Стивен Фью
  - i. Show me the numbers
  - ii. Now you see it

Вопросы?

## Домашнее задание

Seaborn - мощная библиотека для визуализации данных. Изучите ее базовый функционал и выведите график линейной регрессии по данным собранным в домашней работе №2. Подсказку можно найти в проверочном скрипте окружения.

Подумайте над следующими вопросами:

- О чем говорит наклон графика линейной модели?
- Как он соотносится с результатами полученными вами в предыдущем домашнем задании?
- Хорошо ли линейная аппроксимация подходит для ваших данных или видны систематические отклонения?

При оформлении слайдов использованы изображения с следующих страниц.

- <https://flowingdata.com/>
- <http://colorbrewer2.org/#type=sequential&scheme=BuGn&n=3>
- <http://www.datavis.ca/milestones/index.php?group=1950%2B>
- [http://www.infovis-wiki.net/index.php?title=Visual\\_Variables](http://www.infovis-wiki.net/index.php?title=Visual_Variables)
- [https://www.explainxkcd.com/wiki/index.php/1725:\\_Linear\\_Regression](https://www.explainxkcd.com/wiki/index.php/1725:_Linear_Regression)
- [https://en.wikipedia.org/wiki/Pie\\_chart](https://en.wikipedia.org/wiki/Pie_chart)
- <https://en.wikipedia.org/wiki/Histogram>
- [https://en.wikipedia.org/wiki/Line\\_chart](https://en.wikipedia.org/wiki/Line_chart)
- [http://www.bbc.co.uk/bitesize/ks3/maths/handling\\_data/representing\\_data/review/2/](http://www.bbc.co.uk/bitesize/ks3/maths/handling_data/representing_data/review/2/)
- <http://www.excelcharts.com/blog/change-bad-charts-in-the-wikipedia/>
- <http://viz.wtf/>
- [https://www.123rf.com/photo\\_25442337\\_hand-lens-that-magnifies-a-needle-in-a-haystack.html](https://www.123rf.com/photo_25442337_hand-lens-that-magnifies-a-needle-in-a-haystack.html)