

Проклятие размерности

In [1]:

```
1 import numpy as np
2 import scipy.stats as sps
3
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 sns.set(font_scale=1.3)
7
8 red = '#FF3300'
9 blue = '#0099CC'
10 green = '#00CC66'
```

Сгенерируем 1000 случайных векторов в единичном кубе размерности d . Для каждой размерности d от 1 до 3000 посчитаем среднюю норму вектора и стандартное отклонение.

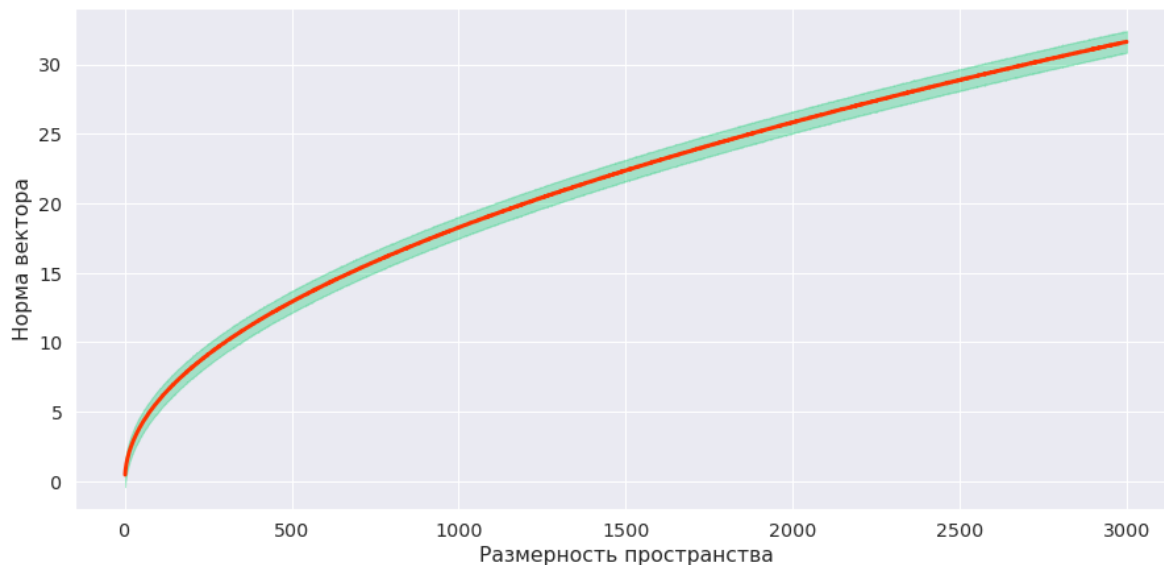
In [2]:

```
1 N_dim = 3000
2 dimentions = np.arange(1, N_dim+1)
3 sample_size = 1000
4
5 means = np.zeros(N_dim)
6 stds = np.zeros(N_dim)
7
8 for d in dimentions:
9     sample = sps.uniform.rvs(size=(sample_size, d))
10    norm = np.sqrt((sample**2).sum(axis=1))
11    means[d-1] = norm.mean()
12    stds[d-1] = norm.std()
```

Визуализируем. Видим, что средняя норма увеличивается с ростом размерности пространства, но разброс остается примерно постоянным.

In [3]:

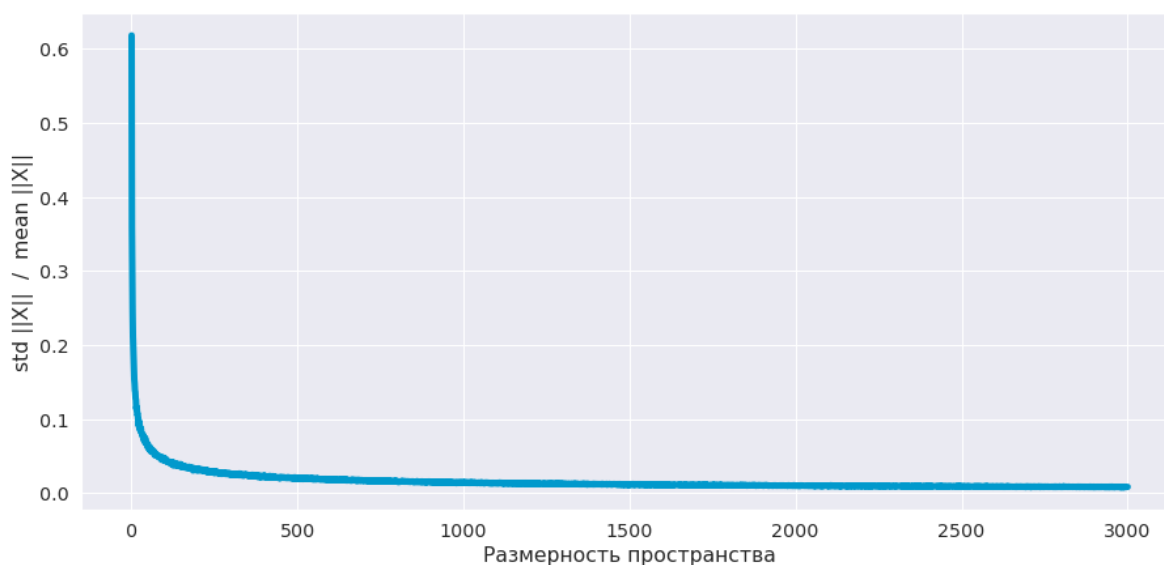
```
1 plt.figure(figsize=(15, 7))
2 plt.plot(dimensions, means, color=red, lw=3)
3 plt.fill_between(dimensions, means - 3*stds, means + 3*stds,
4                 color=green, alpha=0.3)
5 plt.xlabel('Размерность пространства')
6 plt.ylabel('Норма вектора');
```



Посмотрим на график величины разброса нормы по отношению к ее среднему значению. Как и следовало ожидать, эта величина сходится к нулю с ростом размерности пространства.

In [4]:

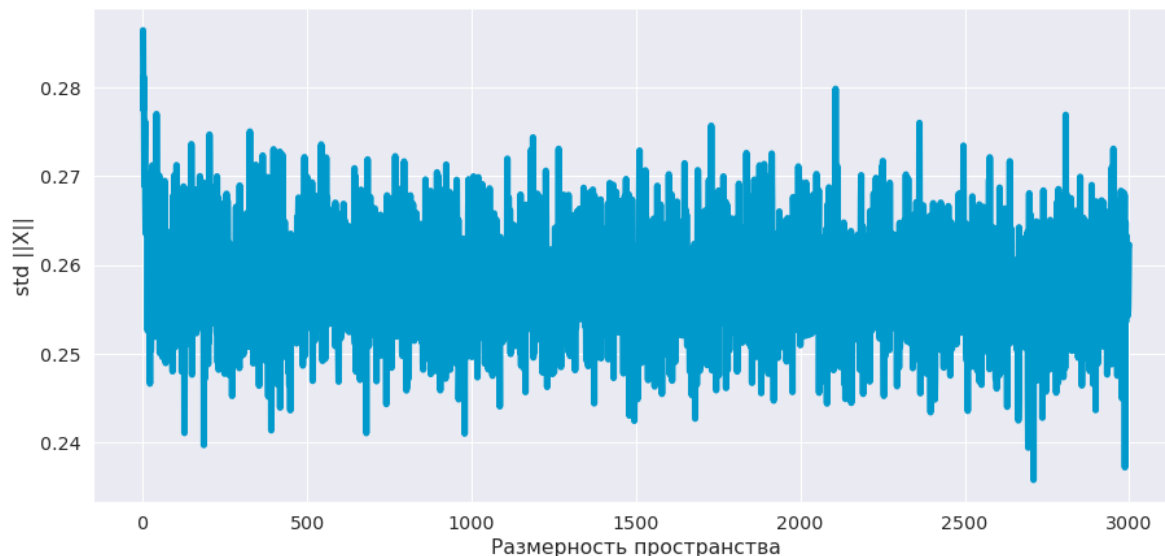
```
1 plt.figure(figsize=(15, 7))
2 plt.plot(dimensions, stds/means, lw=5, color=blue)
3 plt.xlabel('Размерность пространства')
4 plt.ylabel('std ||X|| / mean ||X||');
```



Убедимся, что разброс не меняется с ростом размерности пространства. Шум обусловлен конечным размером выборки.

In [6]:

```
1 plt.figure(figsize=(15, 7))
2 plt.plot(dimensions, stds, lw=5, color=blue)
3 plt.xlabel('Размерность пространства')
4 plt.ylabel('std ||X||');
```



Полученные свойства визуализируют одно из неприятных особенностей проклятия размерности --- неинформативность расстояний. Это проявляется тем, что при анализе данных расстояния между любыми парами точек сосредоточены вокруг некоторого значения, а изменения могут быть сопоставимы с шумом.

По этой причине на практике при анализе данных большой размерности нельзя использовать методы, явным образом использующие расстояния между объектами. Например, метод ближайшего соседа, t-SNE. Однако, некоторые методы могут модифицировать метрику, например, как это делает UMAP. Такие методы можно использовать для анализа данных большой размерности.