

In [1]:

```
1 import numpy as np
2 from sklearn.datasets import make_blobs
3 from sklearn.manifold import TSNE
4 from sklearn.decomposition import PCA
5
6 import seaborn as sns
7 import matplotlib.pyplot as plt
8 %matplotlib inline
9
10 sns.set()
```

## t-SNE

[Страница на сайте автора \(http://lvdmaaten.github.io/tsne/\)](http://lvdmaaten.github.io/tsne/).

[Реализация в R \(https://cran.r-project.org/web/packages/tsne/tsne.pdf\)](https://cran.r-project.org/web/packages/tsne/tsne.pdf).

Реализация в sklearn: `sklearn.manifold.TSNE` (<http://scikit-learn.org/stable/modules/generated/sklearn.manifold.TSNE.html>) (`n_components=2`, `perplexity=30`, ...)

- `n_components` --- число компонент;
- `perplexity` --- перплексия (сглаженный показатель эффективного числа соседей);

Методы:

- `fit_transform(X)` --- обучиться на данных `X` и вернуть сжатое представление `X`.

## Примеры использования метода t-SNE

Генерация выборки размера 500 из пяти "облаков" (гауссовских выборок) в 10-мерном пространстве.

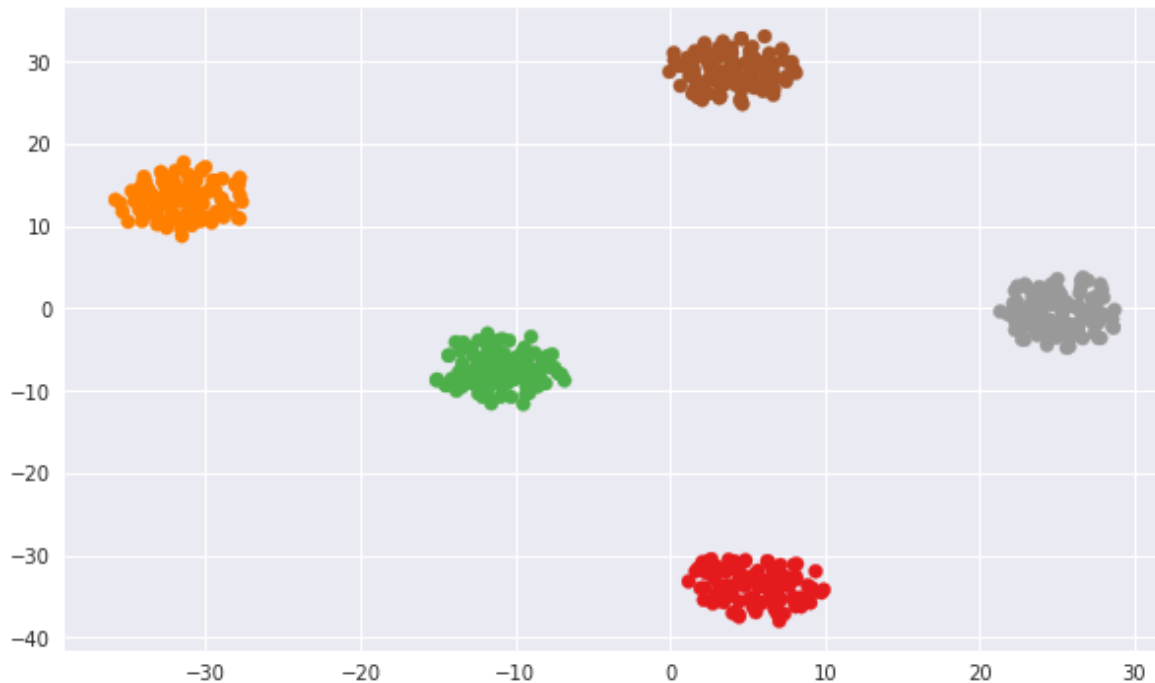
In [2]:

```
1 X, y = make_blobs(n_samples=500, n_features=10, centers=5)
```

Простое применение t-SNE для сжатия в пространство размерности 2.

In [3]:

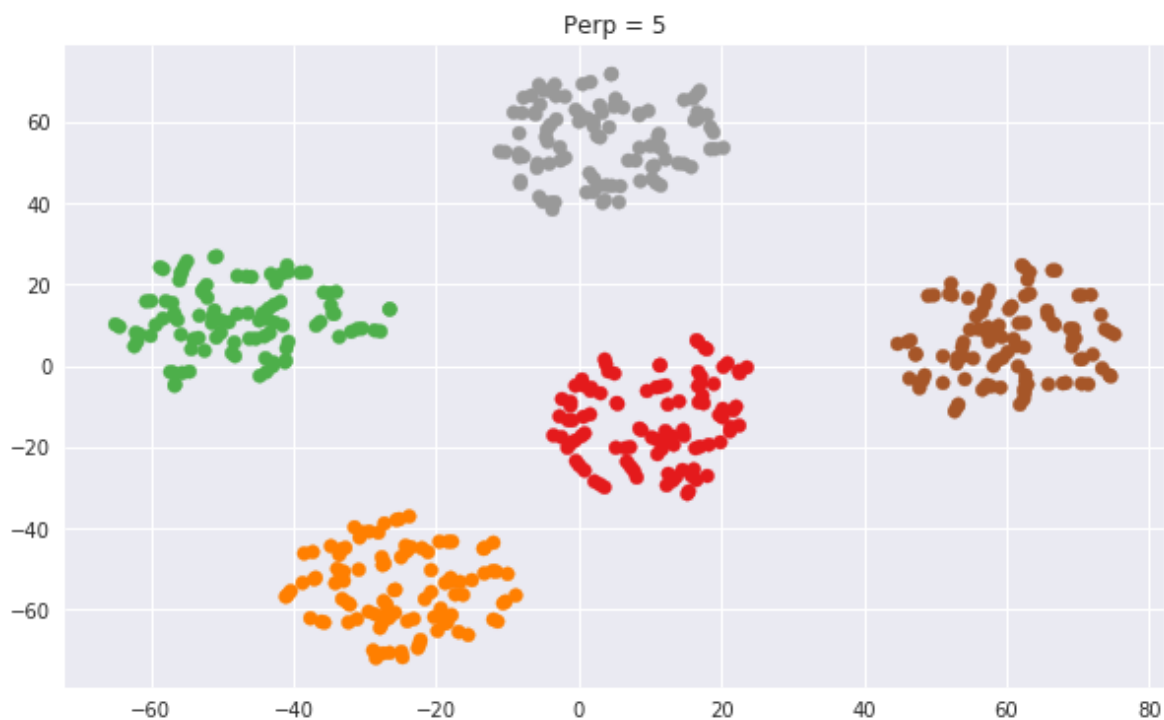
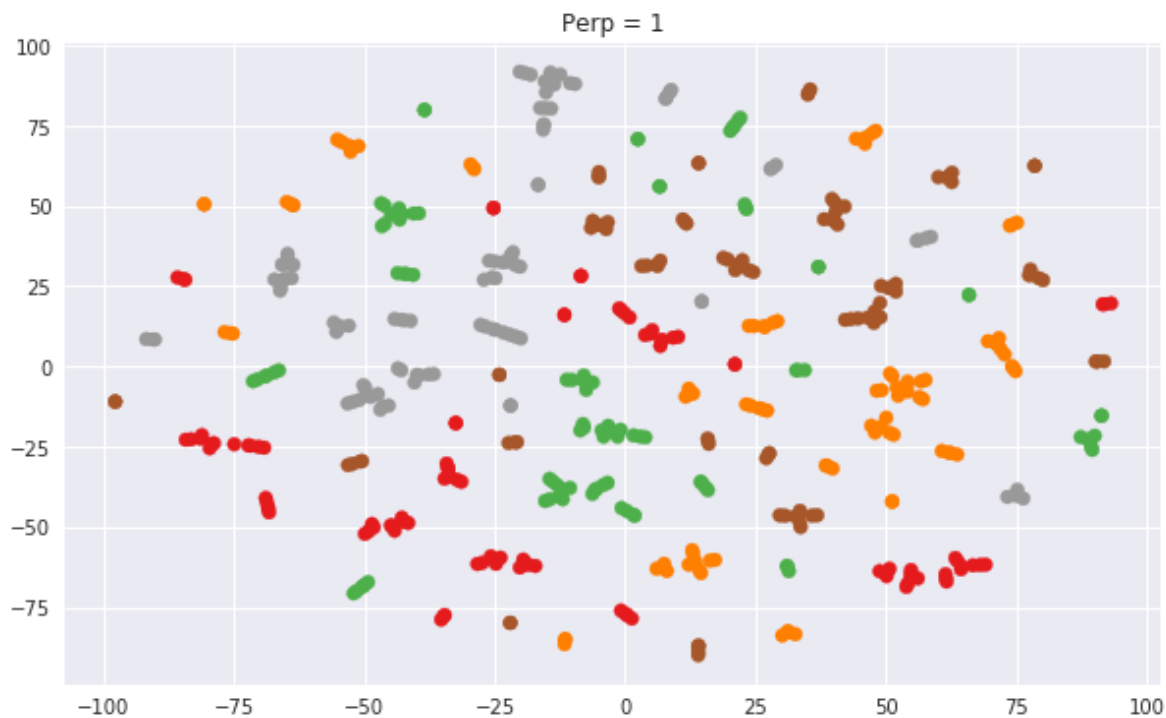
```
1 tsne = TSNE(n_components=2)
2 X_hat = tsne.fit_transform(X)
3
4 plt.figure(figsize=(10, 6))
5 plt.scatter(X_hat[:, 0], X_hat[:, 1], c=y, cmap='Set1')
6 plt.show()
```



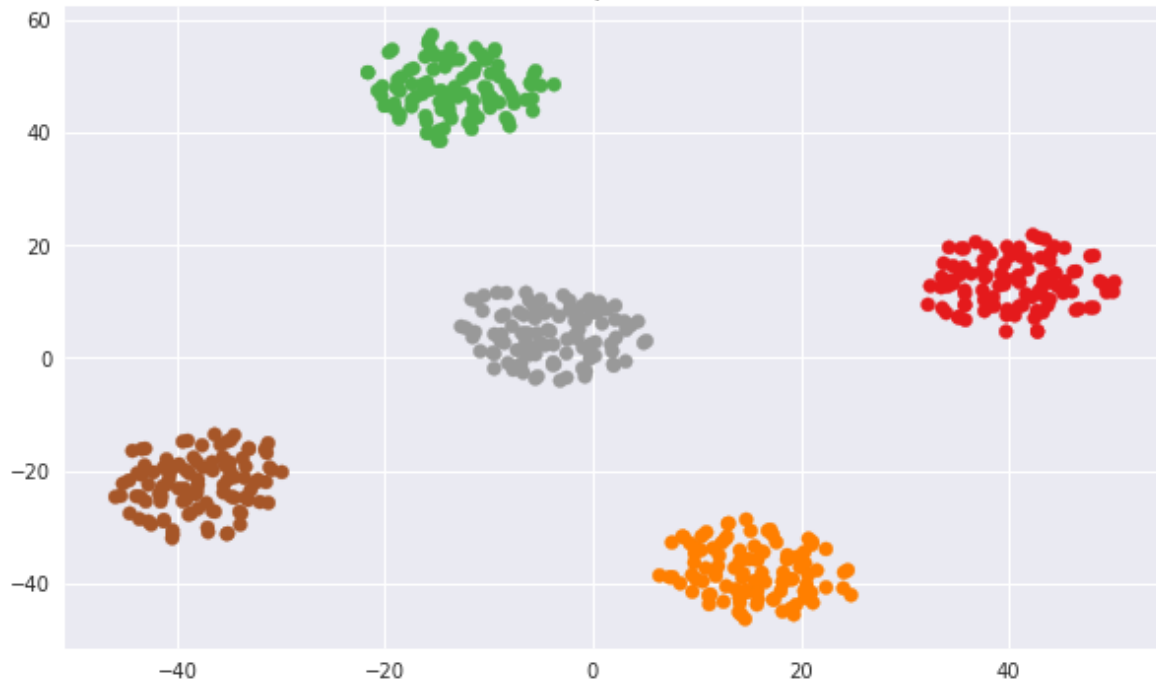
Посмотрим, как меняется расположение точек в двумерном пространстве в зависимости от значения перплексии (сглаженный показатель эффективного числа соседей).

In [4]:

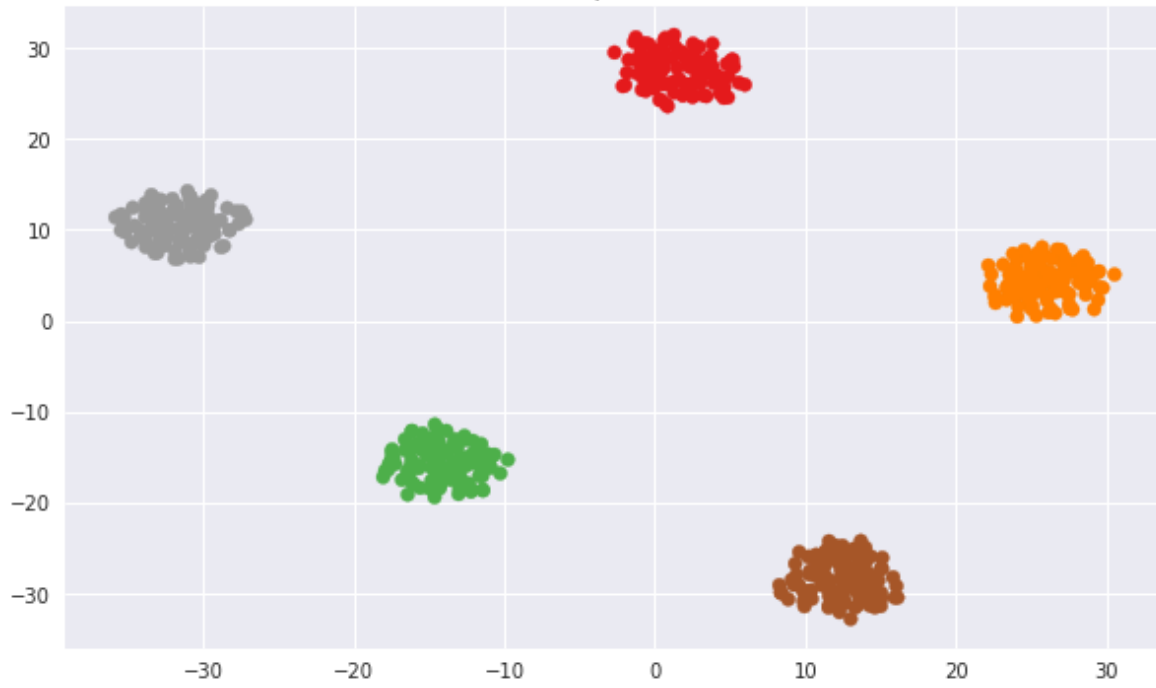
```
1 for perplexity in [1, 5, 15, 30, 50, 75, 100]:
2
3     tsne = TSNE(n_components=2, perplexity=perplexity)
4     X_hat = tsne.fit_transform(X)
5
6     plt.figure(figsize=(10, 6))
7     plt.scatter(X_hat[:, 0], X_hat[:, 1], c=y, cmap='Set1')
8     plt.title('Perp = {}'.format(perplexity))
9     plt.show()
```

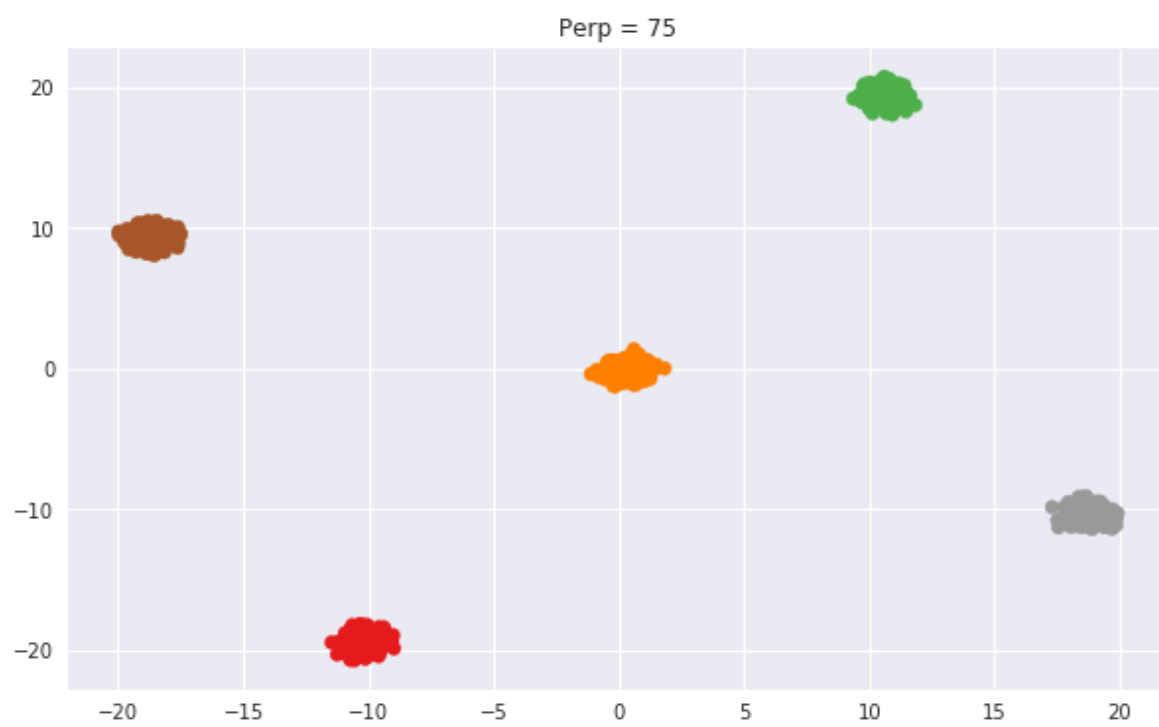
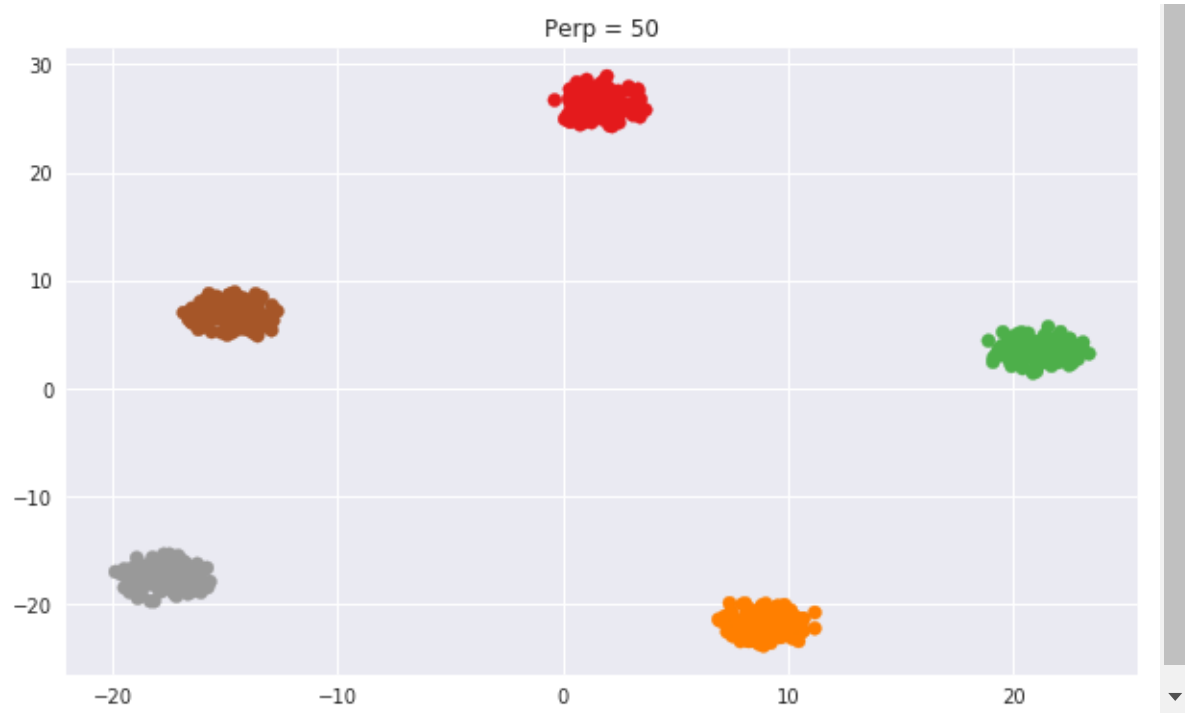


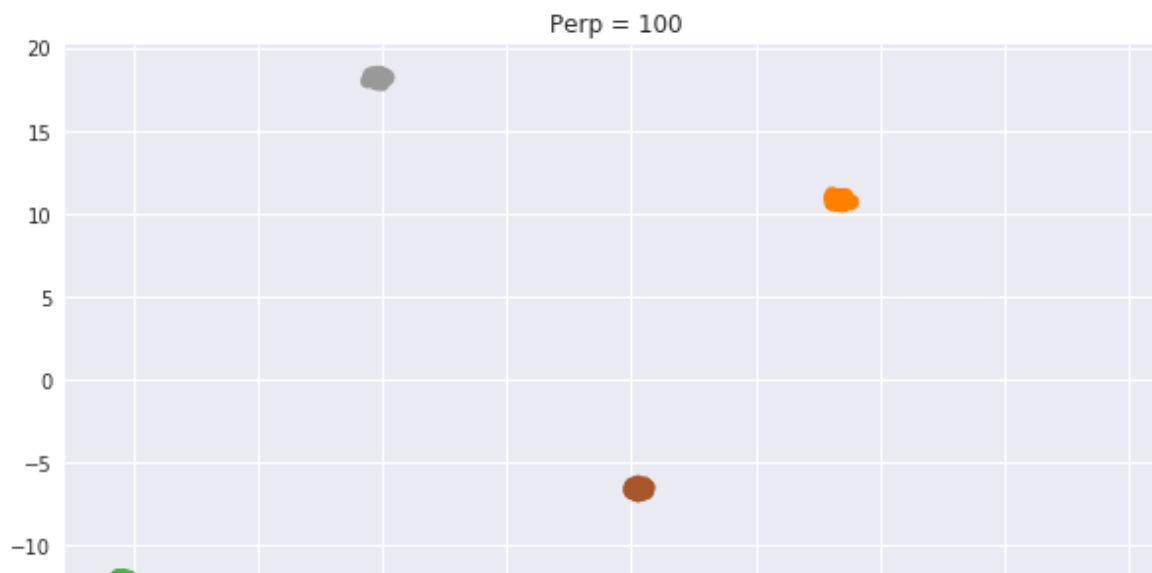
Perp = 15



Perp = 30







Как видим, чем меньше перплексия, тем более картинка получается размытой. Это происходит из-за того, что с уменьшением перплексии уменьшаются величины  $\sigma_i$ , отвечающие за дисперсию. При малом значении перплексии на расположение точки в маломерном пространстве влияют только самые ближайшие соседи.

Если же увеличить дисперсию "облаков", то картинка станет хуже при том же значении перплексии (30).

In [5]:

```
1 X, y = make_blobs(n_samples=500, n_features=10, centers=5, cluster_std=5)
2
3 tsne = TSNE(n_components=2)
4 X_hat = tsne.fit_transform(X)
5
6 plt.figure(figsize=(10, 6))
7 plt.scatter(X_hat[:, 0], X_hat[:, 1], c=y, cmap='Set1')
8 plt.show()
```



**MNIST**

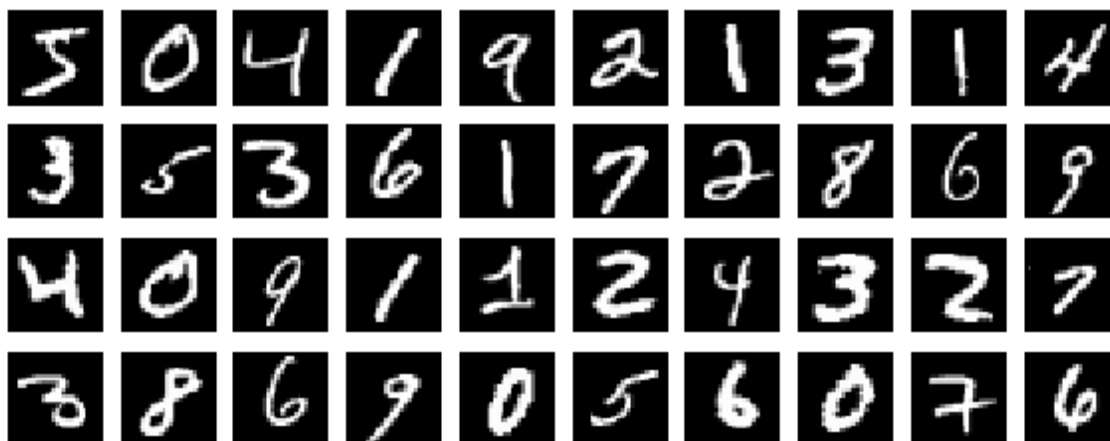
Загрузим часть обычного MNIST'a

In [ ]:

```
1 X = np.loadtxt('train.txt')
2 labels = np.loadtxt('train_labels.txt')
```

In [7]:

```
1 plt.figure(figsize=(10, 4))
2 for i in range(40):
3     plt.subplot(4, 10, i + 1)
4     plt.imshow(X[i].reshape((28, 28)), cmap='gray')
5     plt.axis('off')
```



Применим к нему сначала PCA для сжатия в пространство размерности 30, а затем t-SNE для сжатия в двумерное пространство.

In [8]:

```
1 %%time
2 pca = PCA(n_components=30)
3 Y = pca.fit_transform(X)
```

CPU times: user 11.8 s, sys: 3.2 s, total: 15 s  
Wall time: 5.11 s

Доля дисперсии данных, которую они могут объяснить

In [9]:

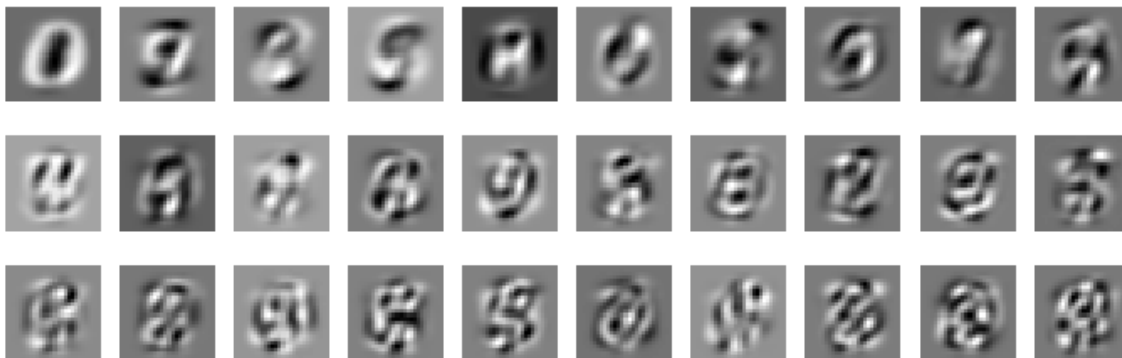
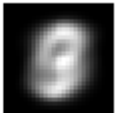
```
1 pca.explained_variance_ratio_.sum()
```

Out[9]:

0.730539514069253

In [10]:

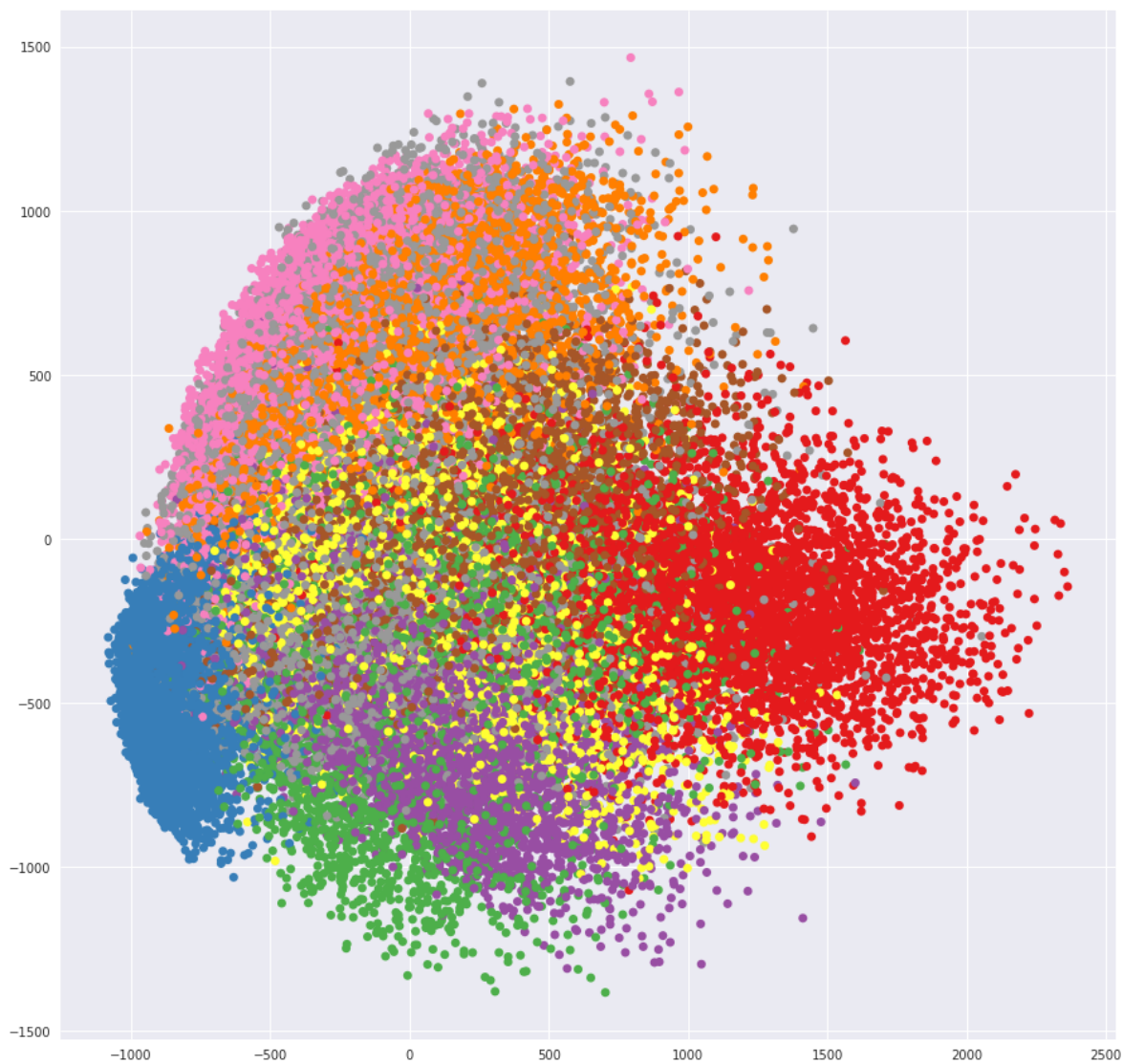
```
1 plt.figure(figsize=(1, 1))
2 plt.imshow(pca.mean_.reshape((28, 28)), cmap='gray')
3 plt.axis('off')
4 plt.show()
5
6 plt.figure(figsize=(12, 4))
7 for i in range(30):
8     plt.subplot(3, 10, i + 1)
9     plt.imshow(pca.components_[i].reshape((28, 28)), cmap='gray')
10    plt.axis('off')
```





In [11]:

```
1 plt.figure(figsize=(15, 15))
2 plt.scatter(Y[:, 0], Y[:, 1], c=labels, linewidths=0, cmap='Set1')
3 plt.show()
```



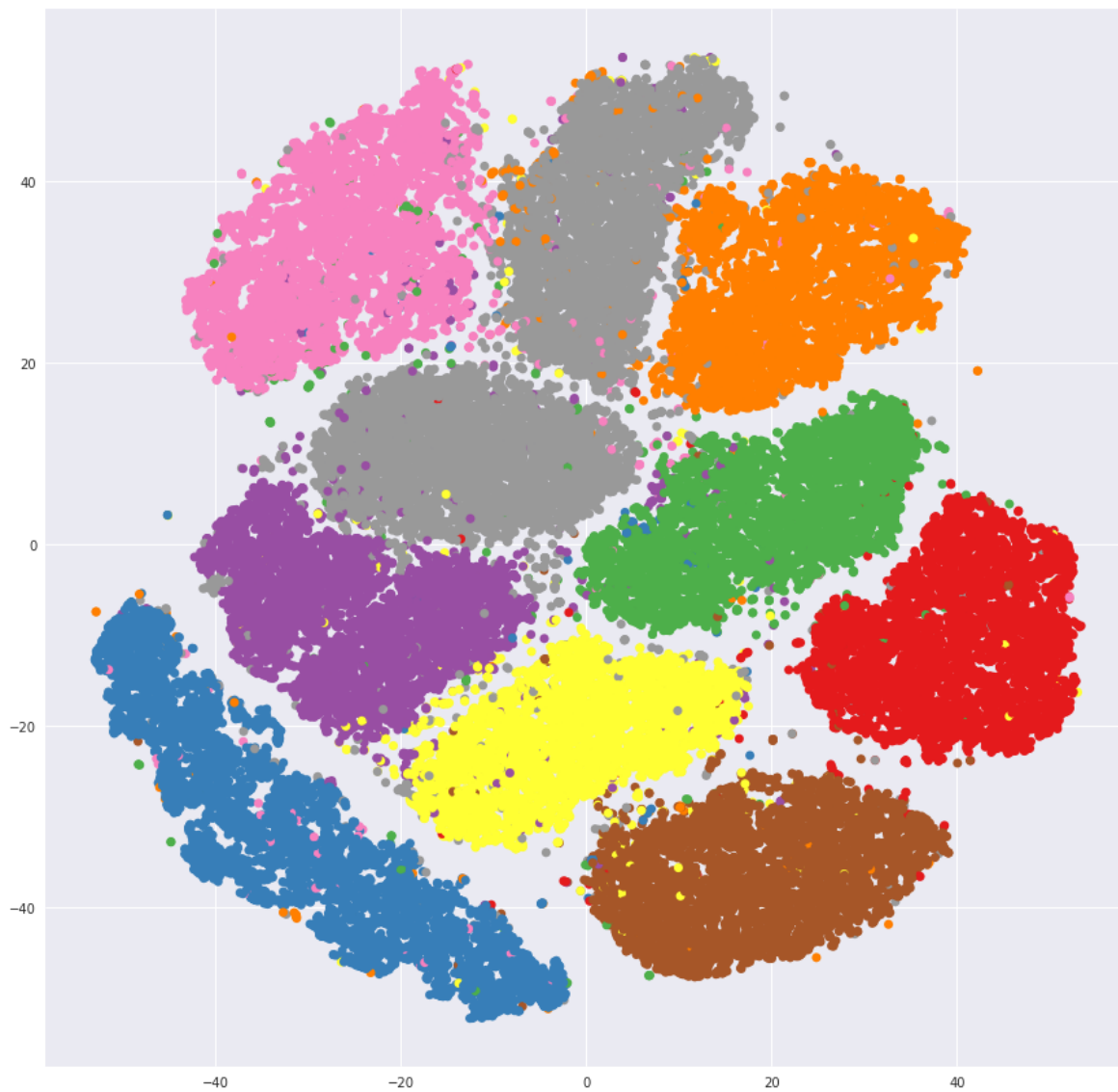
In [6]:

```
1 %%time
2 tsne = TSNE(n_components=2, perplexity=40)
3 Z = tsne.fit_transform(Y)
```

CPU times: user 27min 6s, sys: 1min 47s, total: 28min 54s  
Wall time: 28min 58s

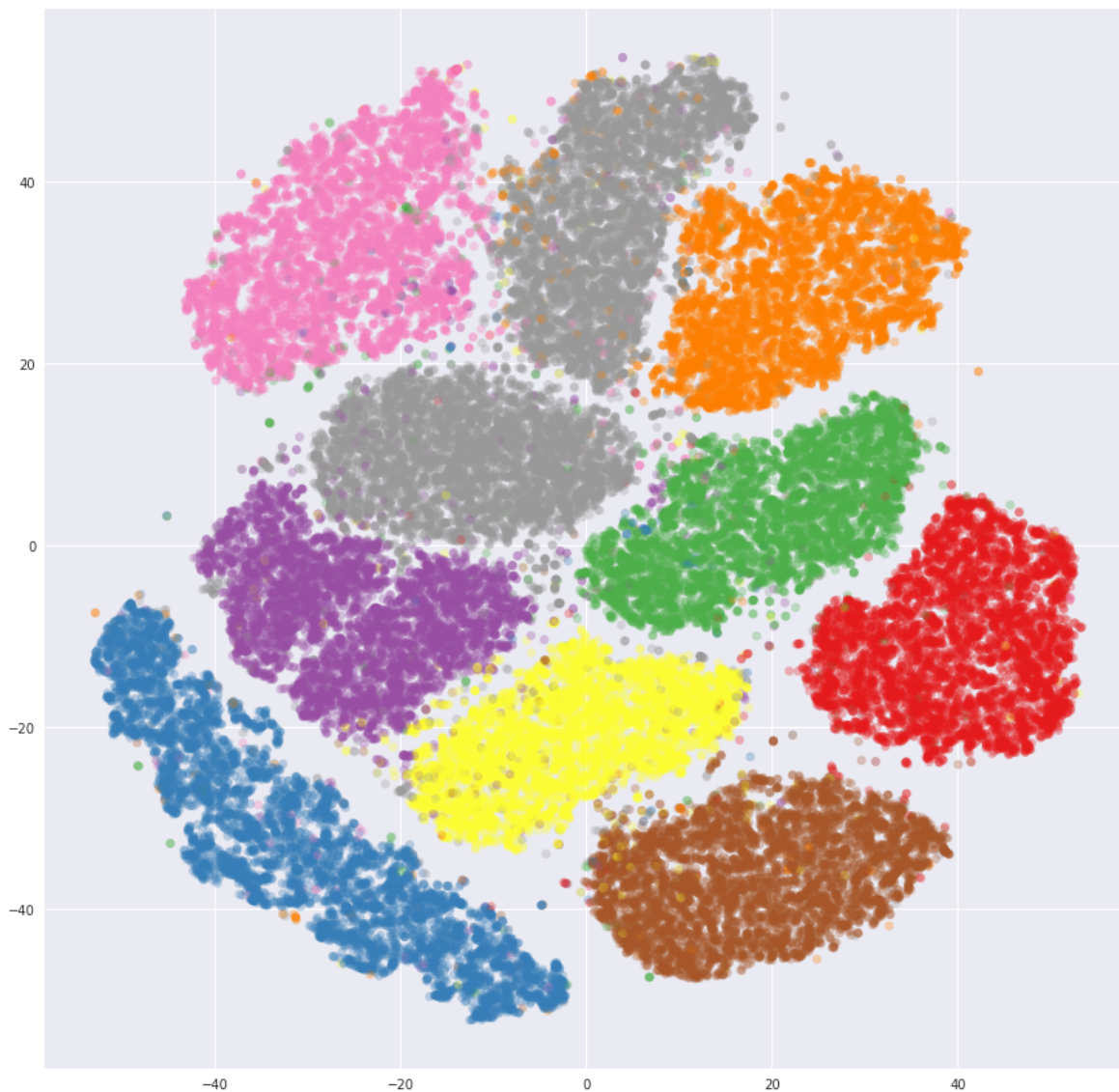
In [7]:

```
1 plt.figure(figsize=(15, 15))
2 plt.scatter(Z[:, 0], Z[:, 1], c=labels, linewidths=0, cmap='Set1')
3 plt.show()
```



In [10]:

```
1 plt.figure(figsize=(15, 15))
2 plt.scatter(Z[:, 0], Z[:, 1], c=labels, linewidths=0, cmap='Set1', alpha=0.3)
3 plt.show()
```



Оригинальная картинка [http://lvdmaaten.github.io/tsne/examples/mnist\\_tsne.jpg](http://lvdmaaten.github.io/tsne/examples/mnist_tsne.jpg)  
([http://lvdmaaten.github.io/tsne/examples/mnist\\_tsne.jpg](http://lvdmaaten.github.io/tsne/examples/mnist_tsne.jpg))

---

Прикладная статистика и анализ данных, 2019

Никита Волков

<https://mipt-stats.gitlab.io/> (<https://mipt-stats.gitlab.io/>)