

Выбор метода кластеризации

Алгоритмы

Рассмотренные нами:

- К-средних
- EM-алгоритм
- Аггломеративная иерархическая кластеризация
- DBSCAN

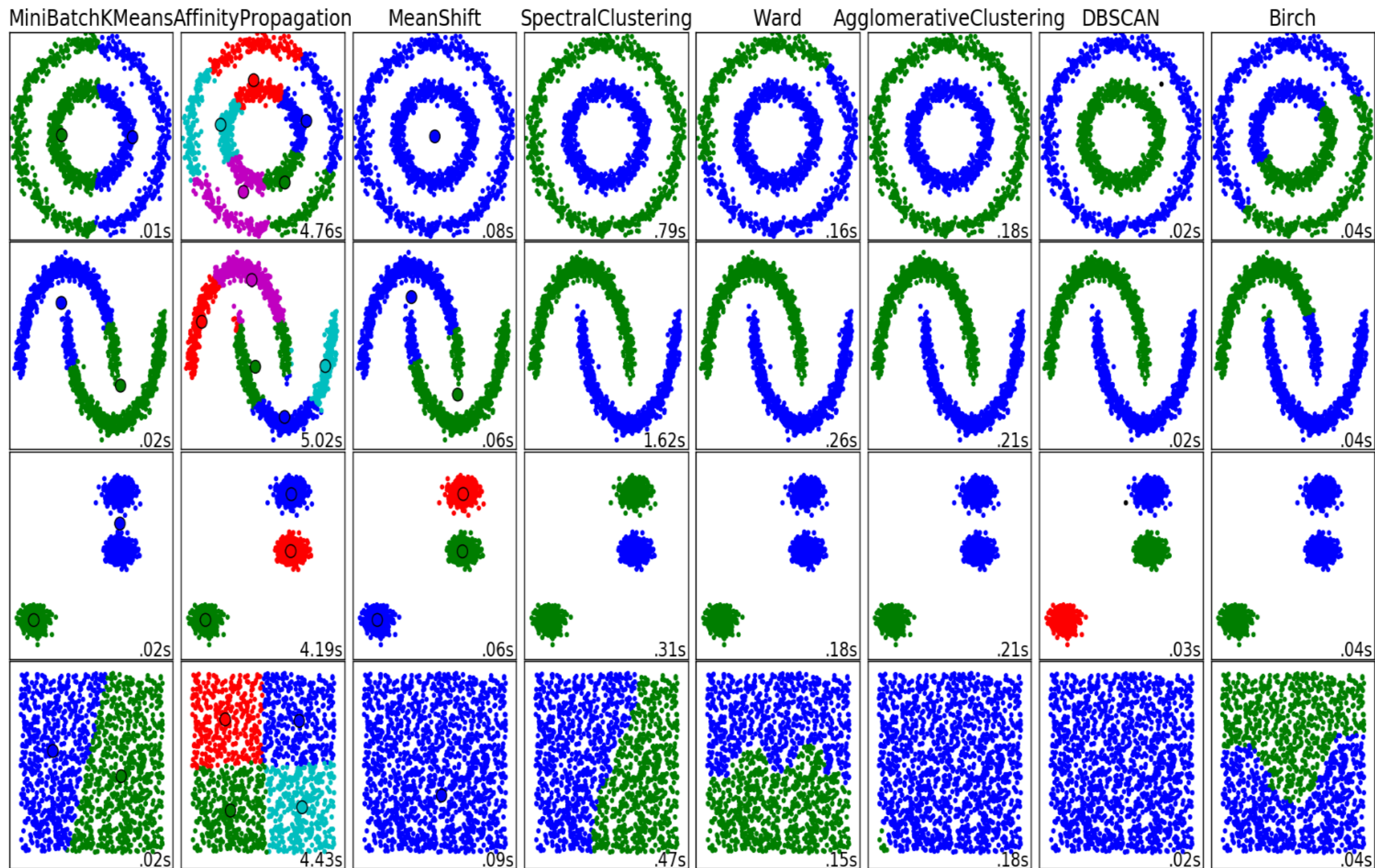
Алгоритмы

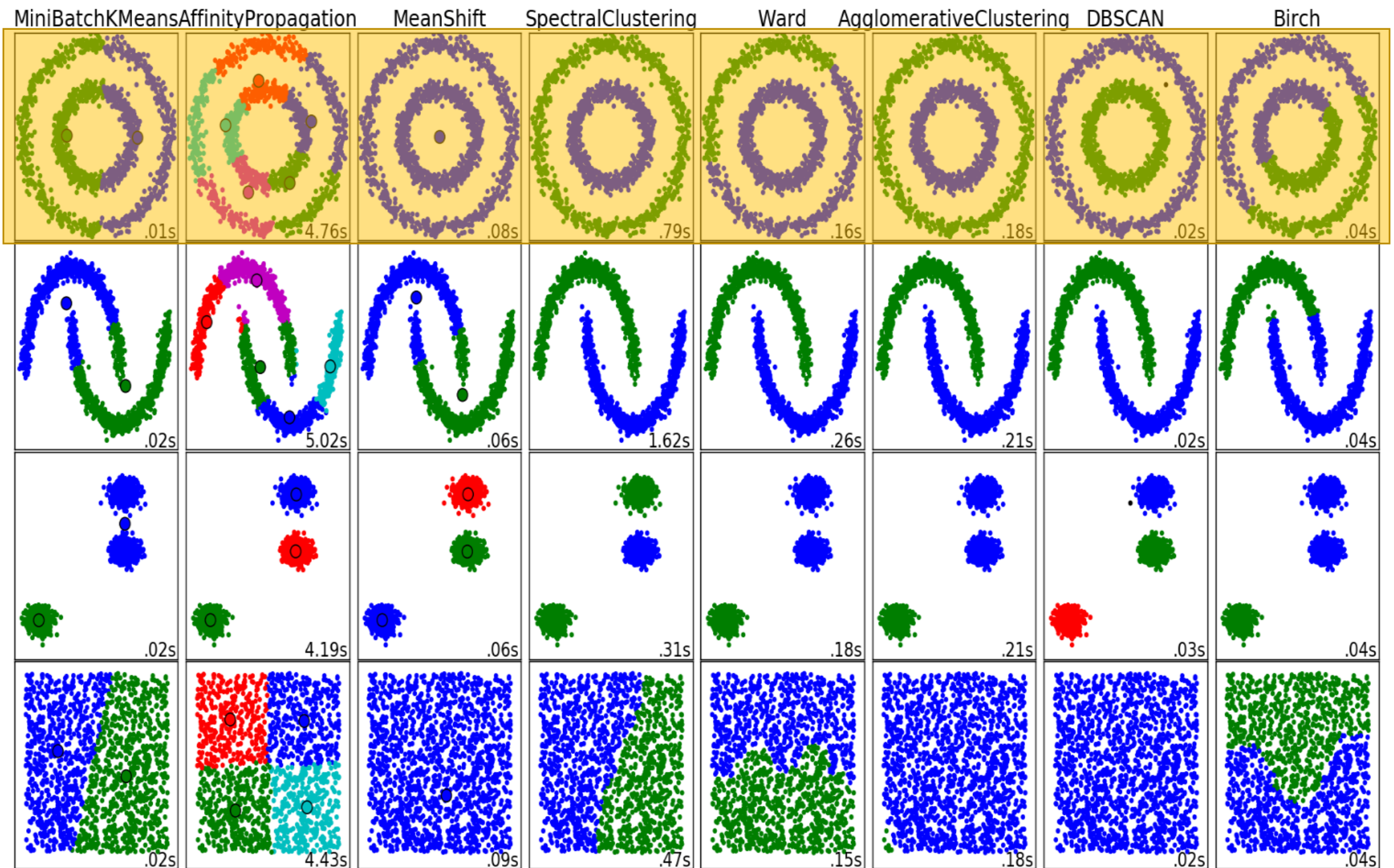
Рассмотренные нами:

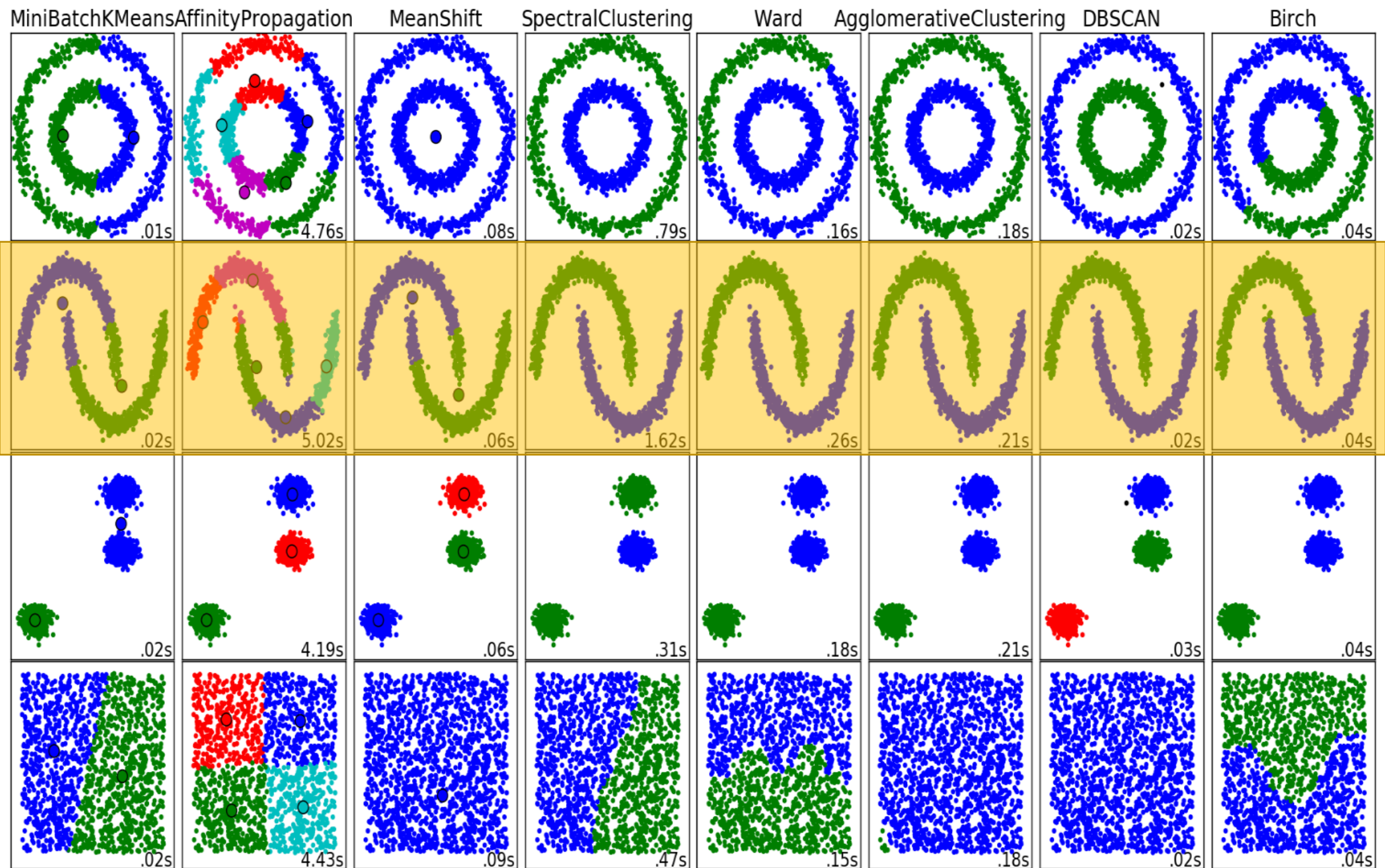
- К-средних
- EM-алгоритм
- Аггломеративная иерархическая кластеризация
- DBSCAN

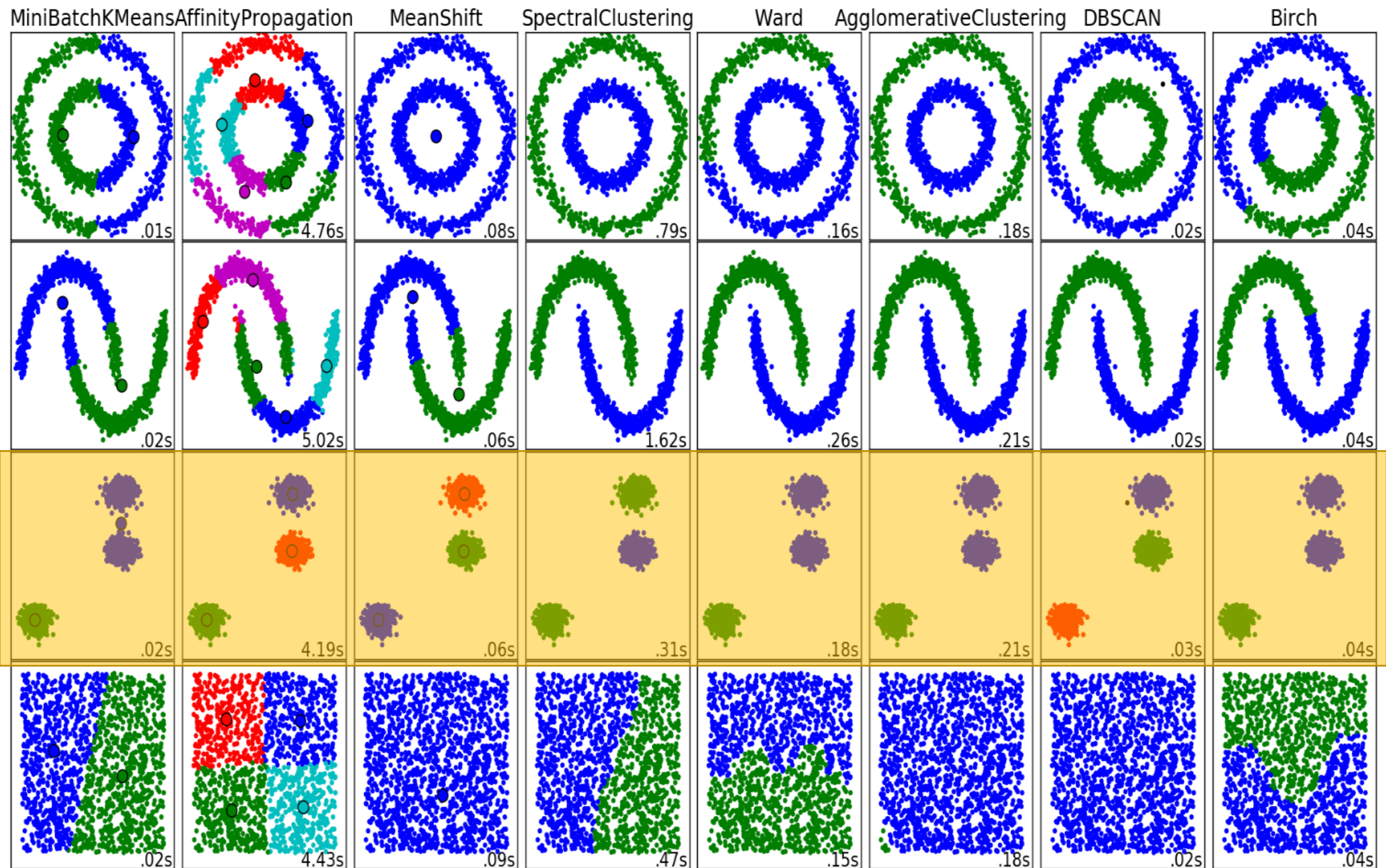
В scikit-learn:

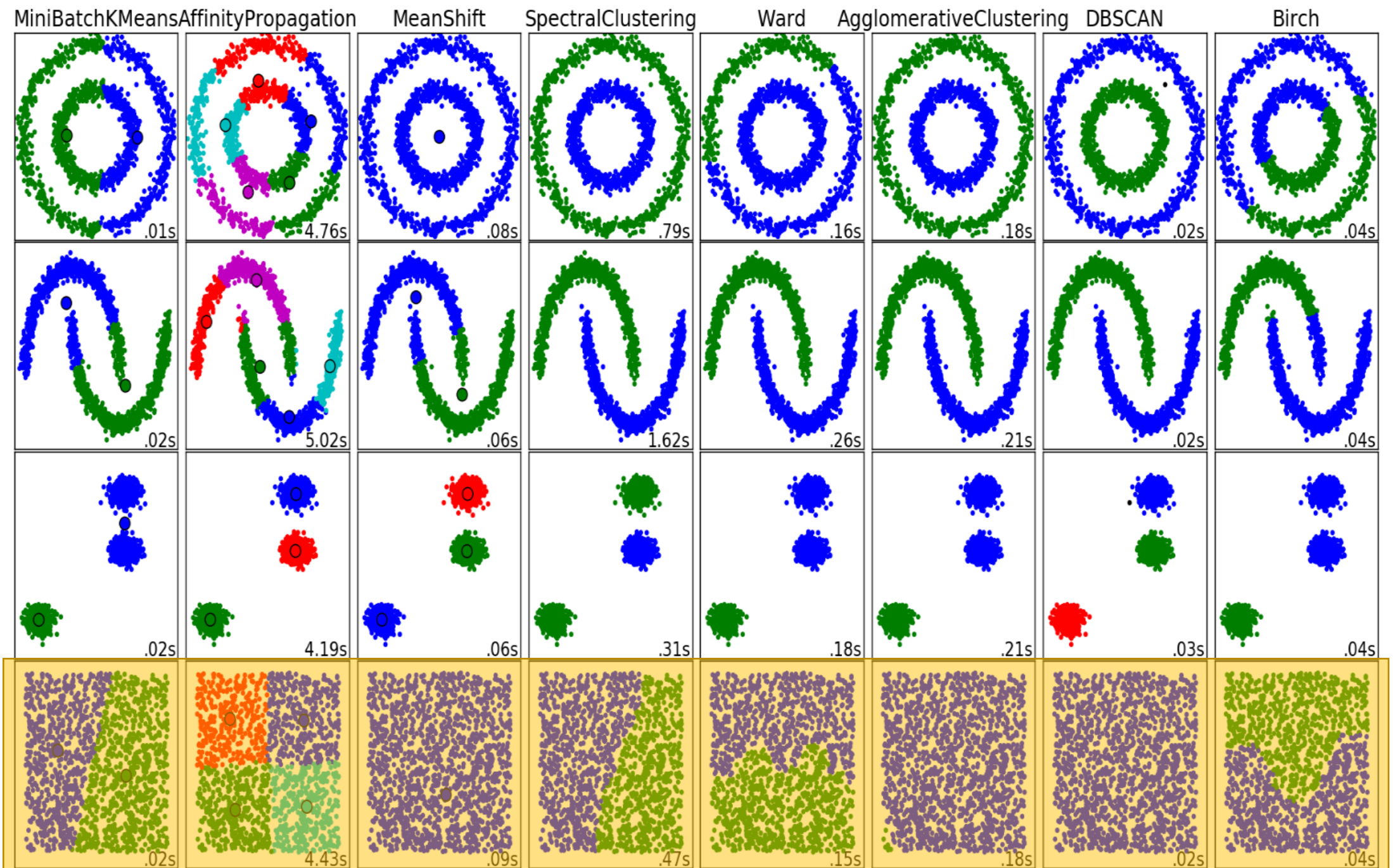
KMeans, MiniBatchKMeans, GaussianMixture,
AgglomerativeClustering, Ward, DBSCAN, MeanShift,
AffinityPropagation, SpectralClustering, Birch

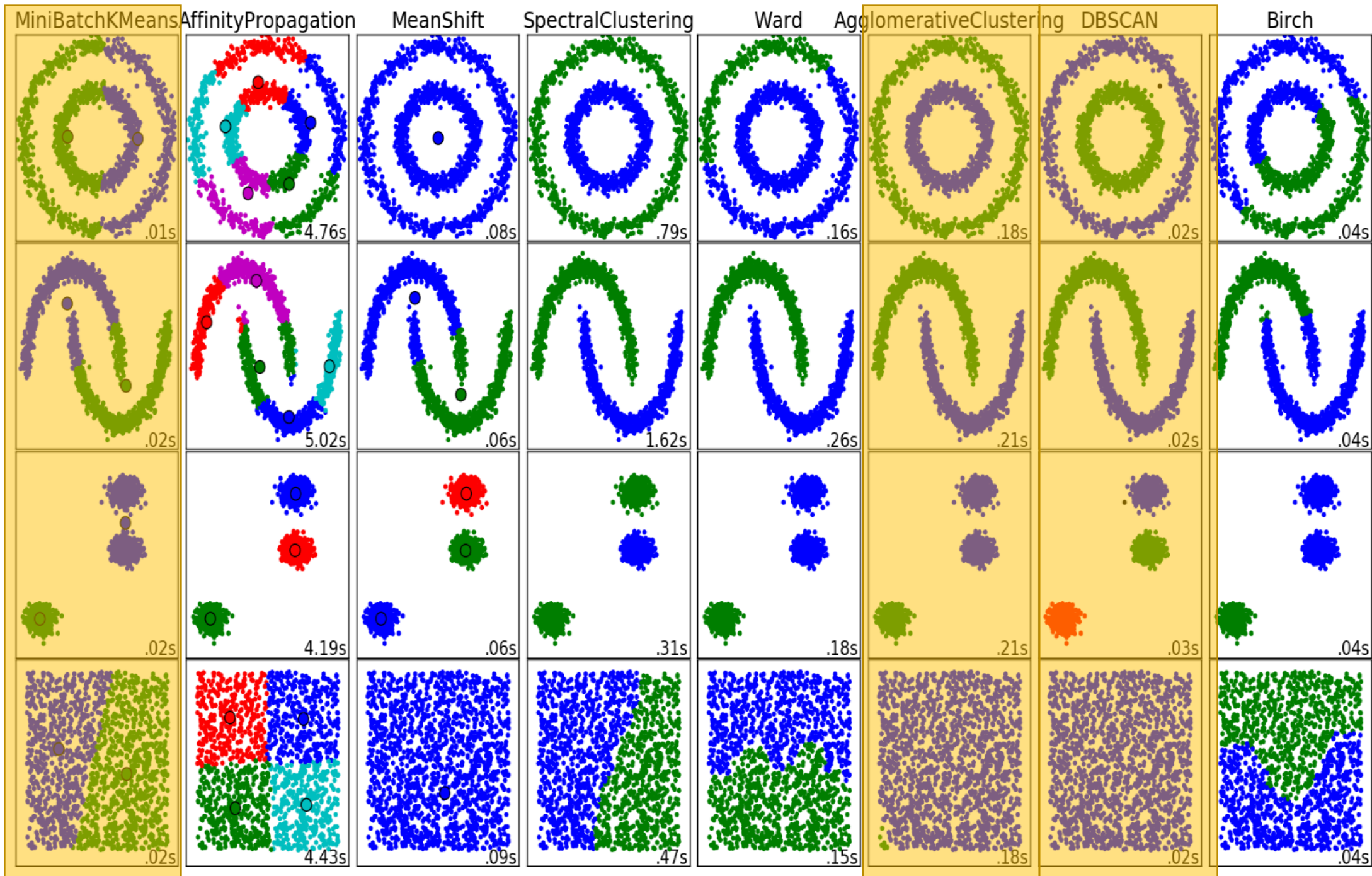












Метод	Параметры	Масштабируемость	Use-case	Геометрия
KMeans	Число кластеров	Очень много примеров (MiniBatch), среднее число кластеров	Выпуклые, примерно одинаковые кластеры	Евклидово расстояние

Метод	Параметры	Масштабируемость	Use-case	Геометрия
KMeans	Число кластеров	Очень много примеров (MiniBatch), среднее число кластеров	Выпуклые, примерно одинаковые кластеры	Евклидово расстояние
GaussianMixture	Веса, векторы средних, матрицы ковариаций	-	Восстановление плотности, выпуклые кластеры	Обобщение евклидовой метрики (с весами)

Метод	Параметры	Масштабируемость	Use-case	Геометрия
KMeans	Число кластеров	Очень много примеров (MiniBatch), среднее число кластеров	Выпуклые, примерно одинаковые кластеры	Евклидово расстояние
GaussianMixture	Веса, векторы средних, матрицы ковариаций	-	Восстановление плотности, выпуклые кластеры	Обобщение евклидовой метрики (с весами)
Agglomerative Clustering	Число кластеров, linkage, метрика	Много примеров и много кластеров	Много кластеров, нужно задавать метрику	Любая метрика/функция близости

Метод	Параметры	Масштабируемость	Use-case	Геометрия
KMeans	Число кластеров	Очень много объектов (MiniBatch), среднее число кластеров	Выпуклые, примерно одинаковые кластеры	Евклидово расстояние
GaussianMixture	Веса, векторы средних, матрицы ковариаций	-	Восстановление плотности, выпуклые кластеры	Обобщение евклидовой метрики (с весами)
Agglomerative Clustering	Число кластеров, linkage, метрика	Много объектов и много кластеров	Много кластеров, нужно задавать метрику/близость (например, косинусную)	Любая метрика/функция близости, для евклидовой - Ward
DBSCAN	Радиус окрестности, число соседей	Много объектов, среднее число кластеров	Неравные невыпуклые кластеры, выбросы,	Евклидово расстояние

Резюме

- К-средних
- EM с нормальным распределением
- Иерархическая аггломеративная кластеризация
- DBSCAN