# Machine Learning Homework 6

## Environment

- Language: C++
- Standard: C++20

## Code with detailed explanations

### Libraries

Used library => Corresponding library in homework description

- Eigen => numpy

- OpenCV, matplotlib++ => visualization

- Boost => filesystem, sort

- OpenMP (for parallel acceleration)

### Visualization

- **Main part**

  > Use `fittingHistory` which records the labels for every iterations to make a video.
  >
  > The duration of one iteration is 1s.
  >
  >
  > `cv::addWeighted` function is used to combine the clustering results (segmentation mask) with original image.
  >
  > `drawMask` function is used to convert the labels to segmentation mask.

```cpp
101        kernelKMeans.fit(kernel);
102        const std::vector<Eigen::VectorXi> &fittingHistory = kernelKMeans.getFittingHistory();
103
104        writer.open(imageFile + "_video.mp4", codec, fps, image.size());
105        // check if we succeeded
106        if (!writer.isOpened())
107        {
108            throw std::runtime_error("Could not open the output video file for write");
109        }
110
111        for (int i = 0; i < fps; i++)
112        {
113            writer.write(image);
114        }
115
116        cv::Mat result;
117        cv::Mat mask;
118        for (std::size_t i = 1; i < fittingHistory.size(); i++)
119        {
120            mask = drawMask(fittingHistory[i], numberOfClusters, image.cols, image.rows);
121            cv::addWeighted(image, 0.5, mask, 0.5, 0, result);
122            for (int i = 0; i < fps; i++)
123            {
124                writer.write(result);
125            }
126        }
127
128        writer.release();
129
130        cv::imwrite(imageFile + "_mask.png", mask);
131        cv::imwrite(imageFile + "_final.png", result);
```

- **Draw mask**

  Convert the labels to segmentation mask.

  `labels` contains the clustering results, from 0 to k-1.

  `cv::applyColorMap` is used to map the labels to the corresponding color to avoid duplicate colors.

```cpp
65  cv::Mat drawMask(const Eigen::VectorXi &labels, int numberOfClusters, unsigned int width, unsigned int height)
66  {
67      Eigen::VectorX<unsigned char> maskData = labels.cast<unsigned char>();
68      cv::Mat mask = cv::Mat(cv::Size(width, height), CV_8UC1, reinterpret_cast<void *>(maskData.data()));
69
70      mask *= (255 / numberOfClusters);
71
72      cv::Mat bgrMask;
73      cv::cvtColor(mask, bgrMask, cv::COLOR_GRAY2BGR);
74      cv::applyColorMap(bgrMask, bgrMask, cv::COLORMAP_COOL);
75      return bgrMask;
76  }
```

# Part1

- **Kernel K-Means**
  - Pseudo-code

    Ignore weight

ALGORITHM 1: Basic Batch Weighted Kernel $k$-means.

$\text{KERNEL\_KMEANS\_BATCH}(K, k, w, t_{max}, \{\pi_c^{(0)}\}_{c=1}^k, \{\pi_c\}_{c=1}^k)$

**Input:** $K$: kernel matrix, $k$: number of clusters, $w$: weights for each point, $t_{max}$: optional maximum number of iterations, $\{\pi_c^{(0)}\}_{c=1}^k$: optional initial clusters

**Output:** $\{\pi_c\}_{c=1}^k$: final partitioning of the points

1. If no initial clustering is given, initialize the $k$ clusters $\pi_1^{(0)}, ..., \pi_k^{(0)}$ (i.e., randomly). Set $t = 0$.

2. For each point $\mathbf{a}_i$ and every cluster $c$, compute

$$d(\mathbf{a}_i, \mathbf{m}_c) = K_{ii} - \frac{2\sum_{\mathbf{a}_j \in \pi_c^{(t)}} w_j K_{ij}}{\sum_{\mathbf{a}_j \in \pi_c^{(t)}} w_j} + \frac{\sum_{\mathbf{a}_j, \mathbf{a}_l \in \pi_c^{(t)}} w_j w_l K_{jl}}{(\sum_{\mathbf{a}_j \in \pi_c^{(t)}} w_j)^2}.$$

3. Find $c^*(\mathbf{a}_i) = \text{argmin}_c d(\mathbf{a}_i, \mathbf{m}_c)$, resolving ties arbitrarily. Compute the updated clusters as

$$\pi_c^{(t+1)} = \{\mathbf{a} \; : \; c^*(\mathbf{a}_i) = c\}.$$

4. If not converged or $t_{max} > t$, set $t = t + 1$ and go to Step 2; Otherwise, stop and output final clusters $\{\pi_c^{(t+1)}\}_{c=1}^k$.

○ Main function

> Arguments
>
> path: the image data folder
>
> numberOfClusters: the number of clusters
>
> init: the selected initialization method, random or k-means++
>
> gamma1, gamma2: the hyper-parameter of RBF kernel

```cpp
78   void run(const fs::path &path, int numberOfClusters, mlhw6::KMeansInitMethods init, double gamma1, double gamma2)
79   {
80       int fps = 30;
81       int codec = cv::VideoWriter::fourcc('m', 'p', '4', 'v');
82       cv::VideoWriter writer;
83
84       mlhw6::KernelKMeans kernelKMeans(numberOfClusters, 200, 1234, init);
85
86       for (auto imageFile : IMAGE_FILES)
87       {
88           std::cout << imageFile << std::endl;
89
90           auto fileName = (path / imageFile).generic_string();
91
92           // read image
93           auto image = cv::imread(fileName, cv::ImreadModes::IMREAD_COLOR);
94
95           // extract RGB values and coordinates
96           auto [pixels, coordinates] = preprocess(image);
97
98           // calculate kernel
99           auto kernel = calculateKernel(pixels, coordinates, gamma1, gamma2);
100
101          kernelKMeans.fit(kernel);
102          const std::vector<Eigen::VectorXi> &fittingHistory = kernelKMeans.getFittingHistory();
103
104          writer.open(imageFile + "_video.mp4", codec, fps, image.size());
```

○ `kmeans.h` header: k-means related classes

> Follow the scikit-learn logic design.

```cpp
enum KMeansInitMethods
{
    Random,
    Kmeansplusplus,
};

class BaseKMeans
{
public:
    BaseKMeans(int numberOfClusters, KMeansInitMethods init = KMeansInitMethods::Kmeansplusplus, int maximumEpochs = 200, int seed = 1234);

    virtual void fit(const Eigen::Ref<const Eigen::MatrixXd> &x) = 0;
    virtual Eigen::VectorXi predict(const Eigen::Ref<const Eigen::MatrixXd> &x) const = 0;
    virtual Eigen::VectorXi fitAndPredict(const Eigen::Ref<const Eigen::MatrixXd> &x);

    const std::vector<Eigen::VectorXi> &getFittingHistory() const;

protected:
    Eigen::VectorXi initializeCenters(const Eigen::Ref<const Eigen::MatrixXd> &x, KMeansInitMethods init, int seed, bool precomputed = false) const;

    int numberOfClusters;
    int maximumEpochs;
    int seed;
    KMeansInitMethods init;

    std::vector<Eigen::VectorXi> fittingHistory;
};

class KMeans : public virtual BaseKMeans
{
public:
    using BaseKMeans::BaseKMeans;

    void fit(const Eigen::Ref<const Eigen::MatrixXd> &x) override;
    Eigen::VectorXi predict(const Eigen::Ref<const Eigen::MatrixXd> &x) const override;

private:
    Eigen::VectorXi assignLabels(const Eigen::Ref<const Eigen::MatrixXd> &x) const;
    Eigen::MatrixXd centers;
};

class KernelKMeans : public virtual BaseKMeans
{
public:
    using BaseKMeans::BaseKMeans;

    void fit(const Eigen::Ref<const Eigen::MatrixXd> &x) override;
    Eigen::VectorXi predict(const Eigen::Ref<const Eigen::MatrixXd> &x) const override;

private:
    Eigen::VectorXi assignLabels(const Eigen::Ref<const Eigen::MatrixXd> &x) const;
};
```

- `KernelKMeans` constructor

```cpp
KernelKMeans::KernelKMeans(int numberOfClusters, int maximumEpochs, int seed, KMeansInitMethods init) : BaseKMeans(numberOfClusters, maximumEpochs, seed, init)
{
}
```

```cpp
BaseKMeans::BaseKMeans(int numberOfClusters, int maximumEpochs, int seed, KMeansInitMethods init) : numberOfClusters(numberOfClusters), maximumEpochs(maximumEpochs), seed(seed), init(init)
{
    if (!numberOfClusters > 0)
    {
        throw std::runtime_error("The number of clusters should be larger than 0.");
    }
}
```

- `preprocess` function: extract RGB values and coordinates

  > Arguments
  >     image: (H, W, 3), BGR values

```cpp
std::pair<Eigen::MatrixX3d, Eigen::MatrixX2i> preprocess(const cv::Mat &image)
{
    auto rows = image.rows;
    auto columns = image.cols;
    auto size = rows * columns;

    cv::Mat rgb;
    cv::cvtColor(image, rgb, cv::COLOR_BGR2RGB);

    std::vector<int> coordinates(size * 2);
#pragma omp parallel for collapse(2)
    for (unsigned int i = 0; i < rows; i++)
    {
        for (unsigned int j = 0; j < columns; j++)
        {
            auto index = (i * columns + j) * 2;
            coordinates[index] = i;
            coordinates[index + 1] = j;
        }
    }

    using MatrixX3ucRowMajor = Eigen::Matrix<unsigned char, Eigen::Dynamic, 3, Eigen::RowMajor>;
    using MatrixX2iRowMajor = Eigen::Matrix<int, Eigen::Dynamic, 2, Eigen::RowMajor>;

    return std::make_pair<Eigen::MatrixX3d, Eigen::MatrixX2i>(
        MatrixX3ucRowMajor::Map(rgb.data, size, 3).cast<double>(), MatrixX2iRowMajor::Map(coordinates.data(), size, 2));
}
```

- `calculateKernel` function: calculate all kernel values

  > Arguments
  >     pixels: (N, 3), RGB values
  >     coordinates: (N, 2), coordinates
  >     gamma1: $\gamma_c$ scalar
  >     gamma2: $\gamma_s$ scalar
  >
  > Formula

$$k(x, x') = e^{-\gamma_s \|S(x) - S(x')\|^2} \times e^{-\gamma_c \|C(x) - C(x')\|^2}$$

$S(x)$ is the spatial information (coordinate)

$C(x)$ is the color information (RGB)

```
55 ∨ Eigen::MatrixXd calculateKernel(const Eigen::MatrixX3d &pixels, const Eigen::MatrixX2i &coordinates, double gamma1, double gamma2)
56   {
57       auto colorKernel = mlhw6::rbf(pixels, pixels, gamma1);
58
59       const Eigen::MatrixX2d &tmp = coordinates.cast<double>();
60       auto coordinateKernel = mlhw6::rbf(tmp, tmp, gamma2);
61
62       return colorKernel.cwiseProduct(coordinateKernel);
63   }
```

○ `rbf` function: RBF kernel

Arguments

　x1: $x$ vector

　x2: $x'$ vector

　gamma: $\gamma$ scalar

Formula

$k(x, x') = e^{-\gamma \|x - x'\|^2}$

```
 9     template <typename DerivedA, typename DerivedB, typename Out = Eigen::Matrix<double, DerivedA::RowsAtCompileTime, DerivedB::RowsAtCompileTime>>
10     Out rbf(const Eigen::MatrixBase<DerivedA> &x1, const Eigen::MatrixBase<DerivedB> &x2, double gamma)
11     {
12         Out result(x1.rows(), x2.rows());
13 #pragma omp parallel for
14         for (Eigen::Index i = 0; i < x1.rows(); i++)
15         {
16             result.row(i) = (-gamma * (x2.rowwise() - x1.row(i)).rowwise().squaredNorm().transpose()).array().exp();
17         }
18         return result;
19     }
20   }
```

○ `kernelKMeans.fit` function: fit the data

Previously, we already calculated the kernel.

So, we use the precomputed kernel values directly.

Arguments

　x: (N, N), the kernel values (gram matrix, similarity matrix)

Steps

　1. Pick k centers

　2. Calculate the cost between the data and centers

　3. Assign the label which has the smallest distance to the data

　4. Keep repeating 2, 3 step until the labels are not changed

`fittingHistory` is used to store the labels for every iterations.

Note: At the line 183, we do $1 - x$ to get the distance matrix.

```
178 ∨     void KernelKMeans::fit(const Eigen::Ref<const Eigen::MatrixXd> &x)
179       {
180           this→fittingHistory = std::vector<Eigen::VectorXi>{Eigen::VectorXi::Constant(x.rows(), -1)};
181
182           // x is similarity matrix (gram matrix)
183           auto centers = this→initializeCenters(1 - x.array(), this→init, this→seed, true);
184           this→fittingHistory.back()(centers).setLinSpaced(0, this→numberOfClusters - 1);
185
186           int epoch = 0;
187           bool sameLabels = false;
188 ∨         do
189           {
190               this→fittingHistory.push_back(this→assignLabels(x));
191
192               sameLabels = (this→fittingHistory[epoch].array() == this→fittingHistory[epoch + 1].array()).all();
193               epoch++;
194           } while (!sameLabels && epoch < this→maximumEpochs);
195
196           std::cout << "Finished at " << epoch << " epoch" << std::endl;
197       }
```

- ○ `assignLabels` function: calculate the cost and assign labels

  Arguments

  x: (N, N), the kernel values

  Steps

  1. Calculate the cost between the data and centers

  $$\mathbf{k}(x_j, x_j) - \frac{2}{|C_k|}\sum_n \alpha_{kn}\mathbf{k}(x_j, x_n) + \frac{1}{|C_k|^2}\sum_p \sum_q \alpha_{kp}\alpha_{kq}\mathbf{k}(x_p, x_q)$$

  2. Assign the label which has the smallest distance to the data

```cpp
204    Eigen::VectorXi KernelKMeans::assignLabels(const Eigen::Ref<const Eigen::MatrixXd> &x) const
205    {
206        Eigen::MatrixXd distance(x.rows(), this→numberOfClusters);
207    #pragma omp parallel for
208        for (int k = 0; k < this→numberOfClusters; k++)
209        {
210            Eigen::VectorXd selector = this→fittingHistory.back().cwiseEqual(k).cast<double>();
211            auto numberOfXInKCluster = selector.sum();
212            Eigen::MatrixXd xToKCluster = x * selector.asDiagonal();
213
214            Eigen::VectorXd secondTerm = 2 * xToKCluster.rowwise().sum() / numberOfXInKCluster;
215            auto thirdTerm = (selector.asDiagonal() * xToKCluster).sum() / std::pow(numberOfXInKCluster, 2);
216            distance(Eigen::all, k) = (x.diagonal() - secondTerm).array() + thirdTerm;
217        }
218
219        std::vector<int> labels(x.rows());
220    #pragma omp parallel for
221        for (Eigen::Index i = 0; i < x.rows(); i++)
222        {
223            // find nearest neighbor
224            distance.row(i).minCoeff(&labels[i]);
225        }
226        return Eigen::VectorXi::Map(labels.data(), labels.size());
227    }
```

- **Spectral Clustering**

  - ○ Pseudo-code

    - ■ Spectral Clustering (ratio cut)

      > **Unnormalized spectral clustering**
      >
      > Input: Similarity matrix $S \in \mathbb{R}^{n \times n}$, number $k$ of clusters to construct.
      > - Construct a similarity graph by one of the ways described in Section 2. Let $W$ be its weighted adjacency matrix.
      > - Compute the unnormalized Laplacian $L$.
      > - Compute the first $k$ eigenvectors $u_1, \ldots, u_k$ of $L$.
      > - Let $U \in \mathbb{R}^{n \times k}$ be the matrix containing the vectors $u_1, \ldots, u_k$ as columns.
      > - For $i = 1, \ldots, n$, let $y_i \in \mathbb{R}^k$ be the vector corresponding to the $i$-th row of $U$.
      > - Cluster the points $(y_i)_{i=1,\ldots,n}$ in $\mathbb{R}^k$ with the $k$-means algorithm into clusters $C_1, \ldots, C_k$.
      > Output: Clusters $A_1, \ldots, A_k$ with $A_i = \{j| y_j \in C_i\}$.

    - ■ Normalized Spectral Clustering (normalized cut)

      > **Normalized spectral clustering according to Shi and Malik (2000)**
      >
      > Input: Similarity matrix $S \in \mathbb{R}^{n \times n}$, number $k$ of clusters to construct.
      > - Construct a similarity graph by one of the ways described in Section 2. Let $W$ be its weighted adjacency matrix.
      > - Compute the unnormalized Laplacian $L$.
      > - Compute the first $k$ generalized eigenvectors $u_1, \ldots, u_k$ of the generalized eigenproblem $Lu = \lambda Du$.
      > - Let $U \in \mathbb{R}^{n \times k}$ be the matrix containing the vectors $u_1, \ldots, u_k$ as columns.
      > - For $i = 1, \ldots, n$, let $y_i \in \mathbb{R}^k$ be the vector corresponding to the $i$-th row of $U$.
      > - Cluster the points $(y_i)_{i=1,\ldots,n}$ in $\mathbb{R}^k$ with the $k$-means algorithm into clusters $C_1, \ldots, C_k$.
      > Output: Clusters $A_1, \ldots, A_k$ with $A_i = \{j| y_j \in C_i\}$.

  - ○ Main function

```cpp
switch (type)
{
case SpectralClusteringType::Unnormalized:
{
    auto model = mlhw6::SpectralClustering(numberOfClusters, init);
    run(path, model, numberOfClusters, gamma1, gamma2);
    break;
}
case SpectralClusteringType::Normalized:
{
    auto model = mlhw6::NormalizedSpectralClustering(numberOfClusters, init);
    run(path, model, numberOfClusters, gamma1, gamma2);
    break;
}
default:
    std::cerr << "Unknown spectral clustering type." << std::endl;
    return 1;
}
```

```cpp
85   void run(const fs::path &path, mlhw6::BaseSpectralClustering &model, int numberOfClusters, double gamma1, double gamma2)
86   {
87       int fps = 30;
88       int codec = cv::VideoWriter::fourcc('m', 'p', '4', 'v');
89       cv::VideoWriter writer;
90
91       for (auto imageFile : IMAGE_FILES)
92       {
93           std::cout << imageFile << std::endl;
94
95           auto fileName = (path / imageFile).generic_string();
96
97           // read image
98           auto image = cv::imread(fileName, cv::ImreadModes::IMREAD_COLOR);
99
100          // extract RGB values and coordinates
101          auto [pixels, coordinates] = preprocess(image);
102
103          // calculate kernel
104          auto kernel = calculateKernel(pixels, coordinates, gamma1, gamma2);
105
106          model.fit(kernel);
107          const std::vector<Eigen::VectorXi> &fittingHistory = model.getFittingHistory();
108
109          writer.open(imageFile + "_video.mp4", codec, fps, image.size());
```

○ `spectral.h` header: spectral clustering related classes

> Follow the scikit-learn logic design.

```cpp
9    class BaseSpectralClustering
10   {
11   public:
12       BaseSpectralClustering(int numberOfClusters, KMeansInitMethods init = KMeansInitMethods::Kmeansplusplus, int maximumEpochs = 200, int seed = 1234);
13
14       virtual void fit(const Eigen::Ref<const Eigen::MatrixXd> &x) = 0;
15       virtual Eigen::VectorXi fitAndPredict(const Eigen::Ref<const Eigen::MatrixXd> &x);
16
17       const std::vector<Eigen::VectorXi> &getFittingHistory() const;
18       const Eigen::MatrixXd &getEigenMatrix() const;
19
20   protected:
21       KMeans kMeans;
22       Eigen::MatrixXd eigenMatrix;
23
24       int numberOfClusters;
25       unsigned int numberOfThreads;
26   };
27
28   class SpectralClustering : public virtual BaseSpectralClustering
29   {
30   public:
31       using BaseSpectralClustering::BaseSpectralClustering;
32
33       void fit(const Eigen::Ref<const Eigen::MatrixXd> &x) override;
34   };
35
36   class NormalizedSpectralClustering : public virtual BaseSpectralClustering
37   {
38   public:
39       using BaseSpectralClustering::BaseSpectralClustering;
40
41       void fit(const Eigen::Ref<const Eigen::MatrixXd> &x) override;
42   };
```

○ `kmeans.h` header: k-means related classes

> Follow the scikit-learn logic design.

```cpp
enum KMeansInitMethods
{
    Random,
    Kmeansplusplus,
};

class BaseKMeans
{
public:
    BaseKMeans(int numberOfClusters, KMeansInitMethods init = KMeansInitMethods::Kmeansplusplus, int maximumEpochs = 200, int seed = 1234);

    virtual void fit(const Eigen::Ref<const Eigen::MatrixXd> &x) = 0;
    virtual Eigen::VectorXi predict(const Eigen::Ref<const Eigen::MatrixXd> &x) const = 0;
    virtual Eigen::VectorXi fitAndPredict(const Eigen::Ref<const Eigen::MatrixXd> &x);

    const std::vector<Eigen::VectorXi> &getFittingHistory() const;

protected:
    Eigen::VectorXi initializeCenters(const Eigen::Ref<const Eigen::MatrixXd> &x, KMeansInitMethods init, int seed, bool precomputed = false) const;

    int numberOfClusters;
    int maximumEpochs;
    int seed;
    KMeansInitMethods init;

    std::vector<Eigen::VectorXi> fittingHistory;
};

class KMeans : public virtual BaseKMeans
{
public:
    using BaseKMeans::BaseKMeans;

    void fit(const Eigen::Ref<const Eigen::MatrixXd> &x) override;
    Eigen::VectorXi predict(const Eigen::Ref<const Eigen::MatrixXd> &x) const override;

private:
    Eigen::VectorXi assignLabels(const Eigen::Ref<const Eigen::MatrixXd> &x) const;
    Eigen::MatrixXd centers;
};

class KernelKMeans : public virtual BaseKMeans
{
public:
    using BaseKMeans::BaseKMeans;

    void fit(const Eigen::Ref<const Eigen::MatrixXd> &x) override;
    Eigen::VectorXi predict(const Eigen::Ref<const Eigen::MatrixXd> &x) const override;

private:
    Eigen::VectorXi assignLabels(const Eigen::Ref<const Eigen::MatrixXd> &x) const;
};
```

- `preprocess` function: extract RGB values and coordinates

> Arguments
>
> image: (H, W, 3), BGR values

```cpp
std::pair<Eigen::MatrixX3d, Eigen::MatrixX2i> preprocess(const cv::Mat &image)
{
    auto rows = image.rows;
    auto columns = image.cols;
    auto size = rows * columns;

    cv::Mat rgb;
    cv::cvtColor(image, rgb, cv::COLOR_BGR2RGB);

    std::vector<int> coordinates(size * 2);
#pragma omp parallel for collapse(2)
    for (unsigned int i = 0; i < rows; i++)
    {
        for (unsigned int j = 0; j < columns; j++)
        {
            auto index = (i * columns + j) * 2;
            coordinates[index] = i;
            coordinates[index + 1] = j;
        }
    }

    using MatrixX3ucRowMajor = Eigen::Matrix<unsigned char, Eigen::Dynamic, 3, Eigen::RowMajor>;
    using MatrixX2iRowMajor = Eigen::Matrix<int, Eigen::Dynamic, 2, Eigen::RowMajor>;

    return std::make_pair<Eigen::MatrixX3d, Eigen::MatrixX2i>(
        MatrixX3ucRowMajor::Map(rgb.data, size, 3).cast<double>(), MatrixX2iRowMajor::Map(coordinates.data(), size, 2));
}
```

- `calculateKernel` function: calculate all kernel values

> Arguments
>
> pixels: (N, 3), RGB values
>
> coordinates: (N, 2), coordinates
>
> gamma1: $\gamma_c$ scalar
>
> gamma2: $\gamma_s$ scalar
>
> Formula

$$k(x, x') = e^{-\gamma_s \|S(x) - S(x')\|^2} \times e^{-\gamma_c \|C(x) - C(x')\|^2}$$

> $S(x)$ is the spatial information (coordinate)
>
> $C(x)$ is the color information (RGB)

```
55 v Eigen::MatrixXd calculateKernel(const Eigen::MatrixX3d &pixels, const Eigen::MatrixX2i &coordinates, double gamma1, double gamma2)
56   {
57       auto colorKernel = mlhw6::rbf(pixels, pixels, gamma1);
58
59       const Eigen::MatrixX2d &tmp = coordinates.cast<double>();
60       auto coordinateKernel = mlhw6::rbf(tmp, tmp, gamma2);
61
62       return colorKernel.cwiseProduct(coordinateKernel);
63   }
```

- `rbf` function: RBF kernel

  Arguments
    x1: $x$ vector
    x2: $x'$ vector
    gamma: $\gamma$ scalar

  Formula
  $$k(x, x') = e^{-\gamma||x-x'||^2}$$

```
 9       template <typename DerivedA, typename DerivedB, typename Out = Eigen::Matrix<double, DerivedA::RowsAtCompileTime, DerivedB::RowsAtCompileTime>>
10       Out rbf(const Eigen::MatrixBase<DerivedA> &x1, const Eigen::MatrixBase<DerivedB> &x2, double gamma)
11       {
12           Out result(x1.rows(), x2.rows());
13   #pragma omp parallel for
14           for (Eigen::Index i = 0; i < x1.rows(); i++)
15           {
16               result.row(i) = (-gamma * (x2.rowwise() - x1.row(i)).rowwise().squaredNorm().transpose()).array().exp();
17           }
18           return result;
19       }
20   }
```

- `model.fit` function: fit the data

  Previously, we already calculated the kernel.

  So, we use the precomputed kernel values directly.

  Arguments
    x: (N, N), the kernel values (gram matrix, similarity matrix)

    - Spectral Clustering (ratio cut)

      Algorithm

      ---

      **Unnormalized spectral clustering**

      Input:  Similarity matrix $S \in \mathbb{R}^{n \times n}$, number $k$ of clusters to construct.
      - Construct a similarity graph by one of the ways described in Section 2.  Let $W$ be its weighted adjacency matrix.
      - Compute the unnormalized Laplacian $L$.
      - **Compute the first $k$ eigenvectors $u_1, \ldots, u_k$ of $L$.**
      - Let $U \in \mathbb{R}^{n \times k}$ be the matrix containing the vectors $u_1, \ldots, u_k$ as columns.
      - For $i = 1, \ldots, n$, let $y_i \in \mathbb{R}^k$ be the vector corresponding to the $i$-th row of $U$.
      - Cluster the points $(y_i)_{i=1,\ldots,n}$ in $\mathbb{R}^k$ with the $k$-means algorithm into clusters $C_1, \ldots, C_k$.
      Output:  Clusters $A_1, \ldots, A_k$ with $A_i = \{j|\, y_j \in C_i\}$.

      ---

```cpp
void SpectralClustering::fit(const Eigen::Ref<const Eigen::MatrixXd> &x)
{
    Eigen::DiagonalMatrix<double, Eigen::Dynamic> degreeMatrix = x.rowwise().sum().asDiagonal();

    // L = D - W
    Eigen::MatrixXd laplacianMatrix = -x;
    laplacianMatrix.diagonal() += degreeMatrix.diagonal();

    // solve eigen decomposition
    Eigen::EigenSolver<Eigen::MatrixXd> solver(laplacianMatrix, true);
    const Eigen::VectorXcd &eigenValues = solver.eigenvalues();
    const Eigen::MatrixXcd &eigenVectors = solver.eigenvectors();

    std::cout << eigenValues.topRows(5) << std::endl;
    std::cout << eigenVectors.leftCols(5) << std::endl;

    // sort eigenvalues
    std::vector<std::pair<double, Eigen::Index>> eigenPairs(eigenValues.rows());
#pragma omp parallel for
    for (Eigen::Index i = 0; i < eigenValues.rows(); i++)
    {
        eigenPairs[i] = std::make_pair(eigenValues[i].real(), i);
    }
    boost::sort::parallel_stable_sort(eigenPairs.begin(), eigenPairs.end(), this→numberOfThreads);
    // std::partial_sort(eigenPairs.begin(), eigenPairs.begin() + this→numberOfClusters, eigenPairs.end());

    // pick k eigenvectors
    this→eigenMatrix = Eigen::MatrixXd(eigenVectors.rows(), this→numberOfClusters);
#pragma omp parallel for
    for (int i = 0; i < this→numberOfClusters; i++)
    {
        this→eigenMatrix.col(i) = eigenVectors.col(eigenPairs[i].second).real();
    }

    // perform k-means clustering on eigen space
    this→kMeans.fit(this→eigenMatrix);
}
```

■ Normalized Spectral Clustering (normalized cut)

Algorithm

**Normalized spectral clustering according to Shi and Malik (2000)**

Input:  Similarity matrix $S \in \mathbb{R}^{n \times n}$, number $k$ of clusters to construct.
- Construct a similarity graph by one of the ways described in Section 2.  Let $W$ be its weighted adjacency matrix.
- Compute the unnormalized Laplacian $L$.
- Compute the first $k$ generalized eigenvectors $u_1, \ldots, u_k$ of the generalized eigenproblem $Lu = \lambda Du$.
- Let $U \in \mathbb{R}^{n \times k}$ be the matrix containing the vectors $u_1, \ldots, u_k$ as columns.
- For $i = 1, \ldots, n$, let $y_i \in \mathbb{R}^k$ be the vector corresponding to the $i$-th row of $U$.
- Cluster the points $(y_i)_{i=1,\ldots,n}$ in $\mathbb{R}^k$ with the $k$-means algorithm into clusters $C_1, \ldots, C_k$.

Output:  Clusters $A_1, \ldots, A_k$ with $A_i = \{j | y_j \in C_i\}$.

```cpp
void NormalizedSpectralClustering::fit(const Eigen::Ref<const Eigen::MatrixXd> &x)
{
    Eigen::DiagonalMatrix<double, Eigen::Dynamic> degreeMatrix = x.rowwise().sum().asDiagonal();

    // L = D - W
    Eigen::MatrixXd laplacianMatrix = -x;
    laplacianMatrix.diagonal() += degreeMatrix.diagonal();

    // solve generalized eigen decomposition
    Eigen::GeneralizedEigenSolver<Eigen::MatrixXd> solver(laplacianMatrix, degreeMatrix, true);
    const Eigen::VectorXcd &eigenValues = solver.eigenvalues();
    const Eigen::MatrixXcd &eigenVectors = solver.eigenvectors();

    std::cout << eigenValues.topRows(5) << std::endl;
    std::cout << eigenVectors.leftCols(5) << std::endl;

    // sort eigenvalues
    std::vector<std::pair<double, Eigen::Index>> eigenPairs(eigenValues.rows());
#pragma omp parallel for
    for (Eigen::Index i = 0; i < eigenValues.rows(); i++)
    {
        eigenPairs[i] = std::make_pair(eigenValues[i].real(), i);
    }
    boost::sort::parallel_stable_sort(eigenPairs.begin(), eigenPairs.end(), this→numberOfThreads);
    // std::partial_sort(eigenPairs.begin(), eigenPairs.begin() + this→numberOfClusters, eigenPairs.end());

    // pick k eigenvectors
    this→eigenMatrix = Eigen::MatrixXd(eigenVectors.rows(), this→numberOfClusters);
#pragma omp parallel for
    for (int i = 0; i < this→numberOfClusters; i++)
    {
        this→eigenMatrix.col(i) = eigenVectors.col(eigenPairs[i].second).real();
    }

    // perform k-means clustering on eigen space
    this→kMeans.fit(this→eigenMatrix);
}
```

○ `kMeans.fit` function: fit the data

Use the points in eigen space to perform k-means clustering.

Steps

1. Pick k centers

2. E step

    1. Calculate the cost between the data and centers

    2. Assign the label which has the smallest distance to the data

3. M step

    Recalculate the centers by averaging the points which belong to the same cluster.

4. Keep repeating 2, 3 step until the labels are not changed

```cpp
123    void KMeans::fit(const Eigen::Ref<const Eigen::MatrixXd> &x)
124    {
125        this→fittingHistory = std::vector<Eigen::VectorXi>{Eigen::VectorXi::Constant(x.rows(), -1)};
126
127        auto centers = this→initializeCenters(x, this→init, this→seed);
128        this→fittingHistory.back()(centers).setLinSpaced(0, this→numberOfClusters - 1);
129        this→centers = x(centers, Eigen::all);
130
131        int epoch = 0;
132        bool sameLabels = false;
133        do
134        {
135            std::cout << "Epoch: " << epoch << std::endl;
136
137            // E step
138            this→fittingHistory.push_back(this→assignLabels(x));
139
140            // M step
141    #pragma omp parallel for
142            for (int k = 0; k < this→numberOfClusters; k++)
143            {
144                Eigen::VectorXd selector = this→fittingHistory.back().cwiseEqual(k).cast<double>();
145                // calculate center
146                this→centers.row(k) = (selector.asDiagonal() * x).colwise().sum() / selector.sum();
147            }
148
149            sameLabels = (this→fittingHistory[epoch].array() == this→fittingHistory[epoch + 1].array()).all();
150            epoch++;
151        } while (!sameLabels && epoch < this→maximumEpochs);
152    }
```

- `assignLabels` function: calculate the cost and assign labels

    Arguments

        x: (N, N), the kernel values

    Steps

        1. Calculate the cost (Euclidean distance) between the data and centers

        2. Assign the label which has the smallest distance to the data

```cpp
159    Eigen::VectorXi KMeans::assignLabels(const Eigen::Ref<const Eigen::MatrixXd> &x) const
160    {
161        std::vector<int> labels(x.rows());
162    #pragma omp parallel for
163        for (Eigen::Index i = 0; i < x.rows(); i++)
164        {
165            // find nearest neighbor
166            (this→centers.rowwise() - x.row(i)).rowwise().squaredNorm().minCoeff(&labels[i]);
167        }
168        return Eigen::VectorXi::Map(labels.data(), labels.size());
169    }
```

## Part2

Use command to control the number of clusters and other parameters.

- **Kernel K-Means**

init

0 => use random initialization

1 => use k-means++ initialization

```
int main(int argc, char *argv[])
{
    if (argc < 6)
    {
        std::cerr << "Usage: " << argv[0] << " <data path> <number of cluster> <init> <gamma1> <gamma2>" << std::endl;
        return 1;
    }
```

- **Spectral Clustering**

  spectral clustering type

  0 => use spectral clustering algorithm

  1 => use normalized clustering algorithm

  init

  0 => use random initialization

  1 => use k-means++ initialization

```
140   int main(int argc, char *argv[])
141   {
142       if (argc < 7)
143       {
144           std::cerr << "Usage: " << argv[0] << " <data path> <number of cluster> <spectral clustering type> <init> <gamma1> <gamma2>" << std::endl;
145           return 1;
146       }
147
```

# Part3

- `initializeCenters` function: pick k centers initialized by the selected method

  Arguments
  - x: (N, N) precomputed distance matrix or (N, features) data
  - init: the selected initialization method, random or k-means++
  - seed: the random seed
  - precomputed: x is the precomputed distance matrix or not.

```
97    Eigen::VectorXi BaseKMeans::initializeCenters(const Eigen::Ref<const Eigen::MatrixXd> &x, KMeansInitMethods init, int seed, bool precomputed) const
98    {
99        if (precomputed)
100       {
101           if (x.rows() != x.cols())
102           {
103               throw std::runtime_error("The precomputed x should be a squared matrix.");
104           }
105       }
106
107       switch (init)
108       {
109       case KMeansInitMethods::Random:
110           return randomInitialization(x, this->numberOfClusters, seed);
111       case KMeansInitMethods::Kmeansplusplus:
112           return kMeansPlusPlusInitialization(x, this->numberOfClusters, seed, precomputed);
113       default:
114           throw std::runtime_error("The initialization method is not supported.");
115       };
116   }
```

- **Random initialization**

  Arguments
  - x: (N, N) precomputed distance matrix or (N, features) data
  - numberOfClusters: k clusters
  - seed: the random seed

  Steps

  1. generate the sequence of indexes, 0 ~ N-1

  2. shuffle the sequence

  3. pick the top k rows as the centers

```
65    Eigen::VectorXi randomInitialization(const Eigen::Ref<const Eigen::MatrixXd> &x, int numberOfClusters, int seed)
66    {
67        auto rng = std::mt19937_64(seed);
68        std::vector<int> sequence(x.rows());
69        std::iota(sequence.begin(), sequence.end(), 0);
70        std::shuffle(sequence.begin(), sequence.end(), rng);
71
72        Eigen::Map<Eigen::VectorXi> tmp = Eigen::VectorXi::Map(sequence.data(), sequence.size());
73        return tmp.topRows(numberOfClusters);
74    }
```

- **K-Means++ initialization**

> Arguments
>
> x: (N, N) precomputed distance matrix or (N, features) data
>
> numberOfClusters: k clusters
>
> seed: the random seed
>
> Steps
>
> 1. Choose one center uniformly at random among the data points.
>
> 2. For each data point x not chosen yet, compute $D(x)^2$ (the squared Euclidean distance) or use the precomputed distance matrix, the distance between x and the nearest center that has already been chosen.
>
> 3. Choose one new data point at random as a new center, using a weighted probability distribution where a point x is chosen with probability proportional to $D(x)^2$. (The farthest point will be chosen.)

```
12    Eigen::VectorXi kMeansPlusPlusInitialization(const Eigen::Ref<const Eigen::MatrixXd> &x, int numberOfClusters, int seed, bool precomputed)
13    {
14        auto rng = std::mt19937_64(seed);
15        std::vector<int> candidates;
16
17        // 1. Choose one center uniformly at random among the data points.
18        // closed interval [0, rows - 1]
19        auto sampler = std::uniform_int_distribution(0, static_cast<int>(x.rows()) - 1);
20        candidates.push_back(sampler(rng));
21        numberOfClusters--;
22
23        // 2. For each data point x not chosen yet, compute D(x)^2, the distance between x and the nearest center that has already been chosen.
24        Eigen::MatrixXd distances(x.rows(), x.rows());
25        if (!precomputed)
26        {
27 #pragma omp parallel for
28            for (Eigen::Index i = 0; i < x.rows(); i++)
29            {
30                distances.row(i) = (x.rowwise() - x.row(i)).rowwise().squaredNorm().transpose();
31            }
32        }
33        else
34        {
35            distances = x;
36        }
37
38        auto probabilityDistribution = std::uniform_real_distribution();
39        while (numberOfClusters > 0)
40        {
41            Eigen::Map<Eigen::VectorXi> eigenCandidates = Eigen::VectorXi::Map(candidates.data(), candidates.size());
42
43            // 3. Choose one new data point at random as a new center, using a weighted probability distribution where a point x is chosen with probability proportional to D(x)^2.
44            Eigen::VectorXd weights = distances(Eigen::all, eigenCandidates).rowwise().minCoeff();
45            weights /= weights.sum();
46
47            auto probability = probabilityDistribution(rng);
48            for (int i = 0; i < weights.rows(); i++)
49            {
50                auto weight = weights[i];
51                if (probability < weight)
52                {
53                    candidates.push_back(i);
54                    break;
55                }
56                probability -= weight;
57            }
58
59            numberOfClusters--;
60        }
```

# Part4

> Visualize the eigen space.

```
85    void plotEigenSpace(const std::string &path, const Eigen::MatrixXd &eigenMatrix)
86    {
87        matplot::figure();
88
89        const Eigen::VectorXd &vecX = eigenMatrix.col(0);
90        const Eigen::VectorXd &vecY = eigenMatrix.col(1);
91
92        std::vector<double> x(vecX.begin(), vecX.end());
93        std::vector<double> y(vecY.begin(), vecY.end());
94        matplot::scatter(x, y);
95        matplot::save(path);
96    }
```

# Experiments settings and results & Discussion

## Part1

- **Kernel K-Means**

  > Number of clusters: 2
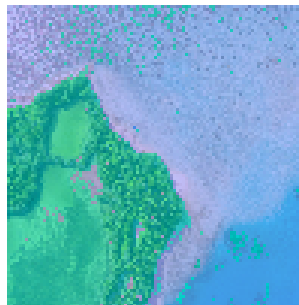  >
  > Initialization method: random
  >
  > Gamma1: 0.00001
  >
  > Gamma2: 0.00001
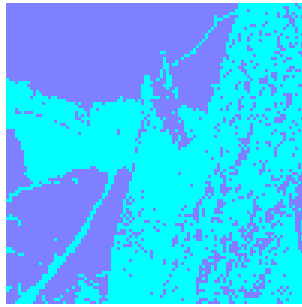
  - image1
    - [Video](#)
    - Final

      

    - Mask

      

  - image2
    - [Video](#)
    - Final

      

    - Mask

## Part1

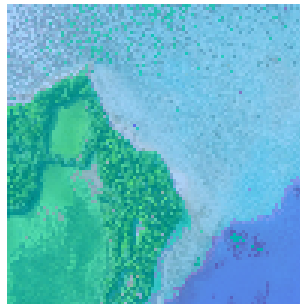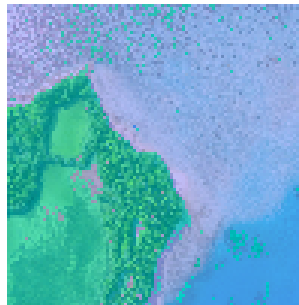- **Kernel K-Means**

- **Spectral Clustering**

# Part2

- **Kernel K-Means**

> Initialization method: random
>
> Gamma1: 0.00001
>
> Gamma2: 0.00001

- K = 3 · image1
    - [Video](#)
    - Final



    - Mask



- K = 3 · image2
    - [Video](#)
    - Final



    - Mask

- K = 4 · image1
  - [Video](Video)
  - Final



  - Mask



- K = 4 · image2
  - [Video](Video)

  - Final



  - Mask



- **Spectral Clustering**

# Part3

- **Kernel K-Means**

  I think there is no significant difference between random and k-means++.

  > Initialization method: k-means++
  >
  > Gamma1: 0.00001
  >
  > Gamma2: 0.00001

  - K = 2 · image1
    - [Video](#)
    - Final

      

    - Mask

      

  - K = 2 · image2
    - [Video](#)
    - Final

      

    - Mask

      

- K = 3 · image1
  - [Video](#)
  - Final

    

  - Mask

    

- K = 3 · image2
  - [Video](#)
  - Final

    

  - Mask

    

- K = 4 · image1
  - [Video](#)
  - Final

- Mask



- K = 4 · image2
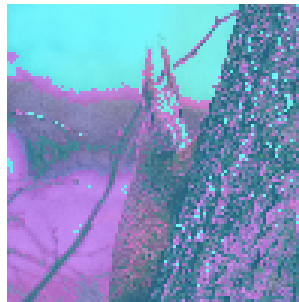  - [Video](#)
  - Final



- Mask



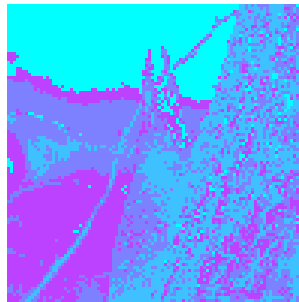- **Spectral Clustering**

## Part4

# Observations and discussion

- Coming soon...