

Machine Learning Homework 6

Machine Learning Homework 6

Environment

Code with detailed explanations

Libraries

Visualization

Part1

Part2

Part3

Part4

Experiments settings and results & Discussion

Part1

Part2

Part3

Part4

Observations and discussion

Environment

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

Code with detailed explanations

Libraries

Used library => Corresponding library in homework description

- Eigen => numpy
- OpenCV => visualization
- Boost (for access file system only)
- 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.

```

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.

```

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.

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 $\pi_1^{(0)}, \dots, \pi_k^{(0)}$ (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 $c^*(\mathbf{a}_i) = \operatorname{argmin}_c d(\mathbf{a}_i, \mathbf{m}_c)$, resolving ties arbitrarily. Compute the updated clusters as

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

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

```

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.

```

5 namespace mlhw6
6 {
7     enum KMeansInitMethods
8     {
9         Random,
10        Kmeansplusplus,
11    };
12
13    class BaseKMeans
14    {
15    public:
16        BaseKMeans(int numberOfClusters, int maximumEpochs, int seed, KMeansInitMethods init);
17
18        virtual void fit(const Eigen::Ref<const Eigen::MatrixXd> &x) = 0;
19        virtual Eigen::VectorXi predict(const Eigen::Ref<const Eigen::MatrixXd> &x) const = 0;
20        virtual Eigen::VectorXi fitAndPredict(const Eigen::Ref<const Eigen::MatrixXd> &x);
21
22        const std::vector<Eigen::VectorXi>& getFittingHistory() const;
23
24    protected:
25        Eigen::VectorXi initializeCenters(const Eigen::Ref<const Eigen::MatrixXd> &x, KMeansInitMethods init, int seed, bool precomputed = false) const;
26
27        int numberOfClusters;
28        int maximumEpochs;
29        int seed;
30        KMeansInitMethods init;
31
32        std::vector<Eigen::VectorXi> fittingHistory;
33    };
34
35    class KMeans : public virtual BaseKMeans
36    {
37    public:
38        KMeans(int numberOfClusters, int maximumEpochs = 200, int seed = 1234, KMeansInitMethods init = KMeansInitMethods::Kmeansplusplus);
39
40        void fit(const Eigen::Ref<const Eigen::MatrixXd> &x) override;
41        Eigen::VectorXi predict(const Eigen::Ref<const Eigen::MatrixXd> &x) const override;
42    private:
43        Eigen::VectorXi assignLabels(const Eigen::Ref<const Eigen::MatrixXd> &x) const;
44        Eigen::MatrixXd centers;
45    };
46
47    class KernelKMeans : public virtual BaseKMeans
48    {
49    public:
50        KernelKMeans(int numberOfClusters, int maximumEpochs = 200, int seed = 1234, KMeansInitMethods init = KMeansInitMethods::Kmeansplusplus);
51
52        void fit(const Eigen::Ref<const Eigen::MatrixXd> &x) override;
53        Eigen::VectorXi predict(const Eigen::Ref<const Eigen::MatrixXd> &x) const override;
54    private:
55        Eigen::VectorXi assignLabels(const Eigen::Ref<const Eigen::MatrixXd> &x) const;
56    };
57

```

- KernelKMeans constructor

```

174 KernelKMeans::KernelKMeans(int numberOfClusters, int maximumEpochs, int seed, KMeansInitMethods init) : BaseKMeans(numberOfClusters, maximumEpochs, seed, init)
175 {
176 }
177
178 BaseKMeans::BaseKMeans(int numberOfClusters, int maximumEpochs, int seed, KMeansInitMethods init) : numberOfClusters(numberOfClusters), maximumEpochs(maximumEpochs), seed(seed), init(init)
179 {
180     if (numberOfClusters > 0)
181     {
182         throw std::runtime_error("The number of clusters should be larger than 0.");
183     }
184 }

```

- preprocess function: extract RGB values and coordinates

Arguments

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

```

25 std::pair<Eigen::MatrixX3d, Eigen::MatrixX2i> preprocess(const cv::Mat &image)
26 {
27     auto rows = image.rows;
28     auto columns = image.cols;
29     auto size = rows * columns;
30
31     cv::Mat rgb;
32     cv::cvtColor(image, rgb, cv::COLOR_BGR2RGB);
33
34     std::vector<int> coordinates(size * 2);
35     #pragma omp parallel for collapse(2)
36     for (unsigned int i = 0; i < rows; i++)
37     {
38         for (unsigned int j = 0; j < columns; j++)
39         {
40             auto index = (i * columns + j) * 2;
41             coordinates[index] = i;
42             coordinates[index + 1] = j;
43         }
44     }
45
46     using MatrixX3ucRowMajor = Eigen::Matrix<unsigned char, Eigen::Dynamic, 3, Eigen::RowMajor>;
47     using MatrixX2iRowMajor = Eigen::Matrix<int, Eigen::Dynamic, 2, Eigen::RowMajor>;
48
49     return std::make_pair<Eigen::MatrixX3d, Eigen::MatrixX2i>(&
50         MatrixX3ucRowMajor::Map(rgb.data, size, 3).cast<double>(), MatrixX2iRowMajor::Map(coordinates.data(), size, 2));
51 }

```

- calculateKernel function: calculate all kernel values

Arguments

pixels: (N, 3), RGB values

coordinates: (N, 2), coordinates

gamma1: γ_c scalar

gamma2: γ_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 }
```

- o `rbf` function: RBF kernel

Arguments

x1: x vector

x2: x' vector

gamma: γ 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 }
```

- o `kernelKMeans.fit` function: fit the data

Previously, we already calculated the kernel first.

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 }
```

- `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 }

```

- `randomInitialization` function: randomly pick k centers

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 }

```

- `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

```

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 numberOfInKCluster = selector.sum();
212             Eigen::MatrixXd xToKCluster = x * selector.asDiagonal();
213
214             Eigen::VectorXd secondTerm = 2 * xToKCluster.rowwise().sum() / numberOfInKCluster;
215             auto thirdTerm = (selector.asDiagonal() * xToKCluster).sum() / std::pow(numberOfInKCluster, 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

-

Part2

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

```

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;
    }
}

```

```

version : 2.0.0 ,
"tasks": [
{
    "label": "launch kernel kmeans",
    "type": "shell",
    "command": "${command:cmake.launchTargetPath}",
    "args": ["${cwd}/data", "3", "1", "0.00001", "0.00001", ">", "output.log"],
    "options": {
        "cwd": "${workspaceFolder}/build",
        "env": {
            "OMP_NUM_THREADS": "16"
        }
    }
},
{
    "label": "launch spectral clustering",
    "type": "shell",
    "command": "${command:cmake.launchTargetPath}",
    "args": ["${cwd}/data", "3", "1", "0.00001", "0.00001", ">", "output.log"],
    "options": {
        "cwd": "${workspaceFolder}/build",
        "env": {
            "OMP_NUM_THREADS": "16"
        }
    }
}
]

```

Part3

- Kernel K-Means: 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.
1. 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.
1. 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.)

```
13 Eigen::VectorXi kMeansPlusPlusInitialization(const Eigen::Ref<const Eigen::MatrixXd> &x, int numberOfClusters, int seed, bool precomputed)
14 {
15     auto rng = std::mt19937_64(seed);
16     std::vector<int> candidates;
17
18     // 1. Choose one center uniformly at random among the data points.
19     // closed interval [0, rows - 1]
20     auto sampler = std::uniform_int_distribution<int>(0, static_cast<int>(x.rows()) - 1);
21     candidates.push_back(sampler(rng));
22     numberOfClusters--;
23
24     // 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.
25     Eigen::MatrixXd distances(x.rows(), x.rows());
26     if (!precomputed)
27     {
28         #pragma omp parallel for
29         for (Eigen::Index i = 0; i < x.rows(); i++)
30         {
31             distances.row(i) = (x.rowwise() - x.row(i)).rowwise().squaredNorm().transpose();
32         }
33     }
34     else
35     {
36         distances = x;
37     }
38
39     auto probabilityDistribution = std::uniform_real_distribution();
40     while (numberOfClusters > 0)
41     {
42         Eigen::Map<Eigen::VectorXi> eigenCandidates = Eigen::VectorXi::Map(candidates.data(), candidates.size());
43
44         // 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.
45         Eigen::VectorXd weights = distances(Eigen::all, eigenCandidates).rowwise().minCoeff();
46         weights /= weights.sum();
47
48         auto probability = probabilityDistribution(rng);
49         for (int l = 0; l < weights.rows(); l++)
50         {
51             auto weight = weights[l];
52             if (probability < weight)
53             {
54                 candidates.push_back(i);
55                 break;
56             }
57             probability -= weight;
58         }
59         numberOfClusters--;
60     }
61 }
```

- Spectral Clustering

Part4

Experiments settings and results & Discussion

Part1

- Kernel K-Means

Number of clusters: 2

Initialization method: random

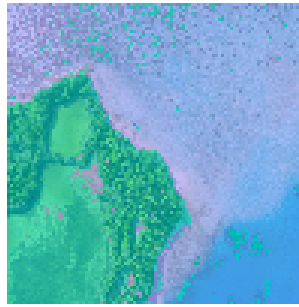
Gamma1: 0.00001

Gamma2: 0.00001

- image1

- [Video](#)
- Final

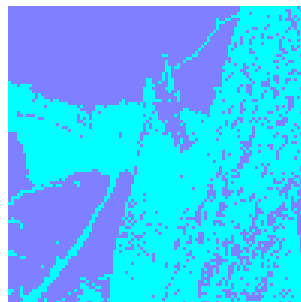
- Mask



- image2
 - [Video](#)
 - Final



- Mask



- **Spectral Clustering**

Part2

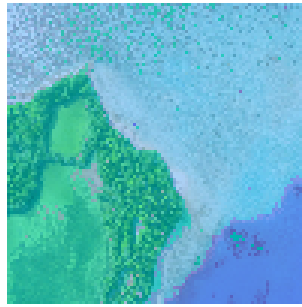
- **Kernel K-Means**

Initialization method: random

Gamma1: 0.00001

Gamma2: 0.00001

- $K = 3 \cdot \text{image1}$
 - [Video](#)
 - Final

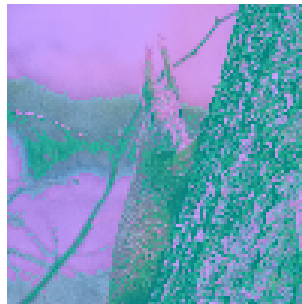


- Mask

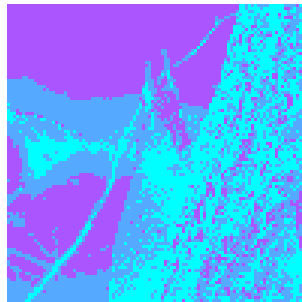


- $K = 3 \cdot \text{image2}$

- [Video](#)
- Final

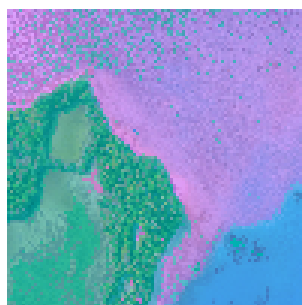


- Mask

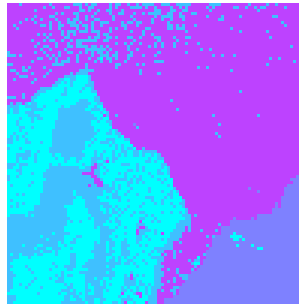


- $K = 4 \cdot \text{image1}$

- [Video](#)
- Final



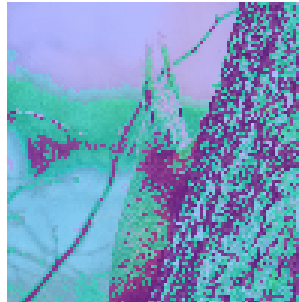
- Mask



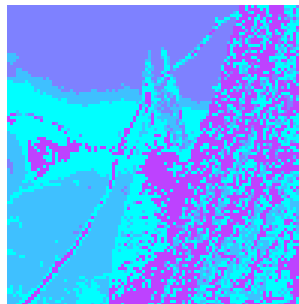
- $K = 4 \cdot \text{image2}$

- [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++

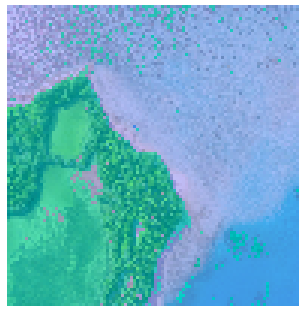
Gamma1: 0.00001

Gamma2: 0.00001

- $K = 2 \cdot \text{image1}$

- [Video](#)

- Final



- Mask

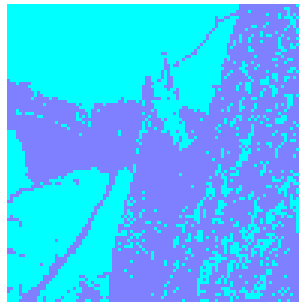


- $K = 2 \cdot \text{image2}$

- [Video](#)
- Final

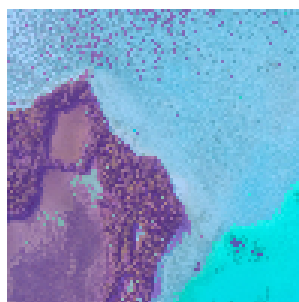


- Mask



- $K = 3 \cdot \text{image1}$

- [Video](#)
- Final



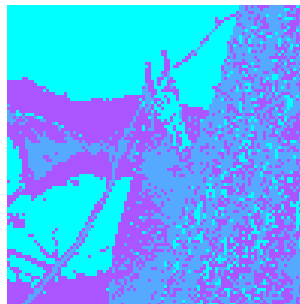
- Mask

- $K = 3 \cdot \text{image2}$

- [Video](#)
- Final

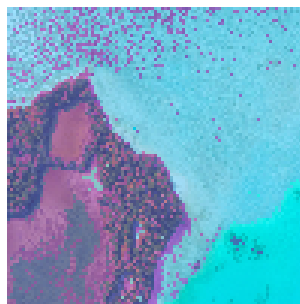


- Mask

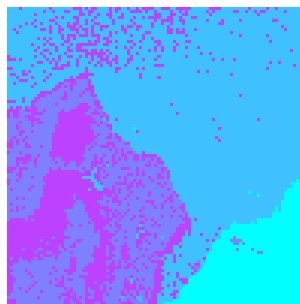


- $K = 4 \cdot \text{image1}$

- [Video](#)
- Final



- Mask



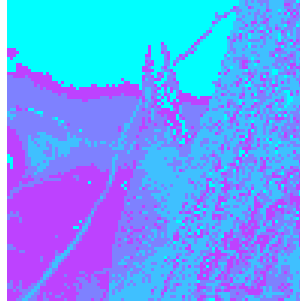
- $K = 4 \cdot \text{image2}$

- [Video](#)

- Final



- Mask



- Spectral Clustering

Part4

Observations and discussion

- Coming soon...