

# Kernel Density Estimates and Mean-Shift Clustering

Jonas Spinner\* – 1927895  
Analytics and Statistics  
KIT – Karlsruhe Institute of Technologie

January 20, 2019

---

\*jonas.spinner@student.kit.edu

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Kernel Density Estimation</b>	<b>5</b>
2.1	Other methods . . . . .	6
2.2	Popular Kernels . . . . .	6
2.3	Bandwidth Selection . . . . .	8
<b>3</b>	<b>Mean Shift Clustering</b>	<b>8</b>
3.1	A general perspective on clustering algorithms . . . . .	8
3.2	Density gradient estimation . . . . .	9
3.3	The algorithm . . . . .	12
3.4	Discussion . . . . .	13
<b>4</b>	<b>Application</b>	<b>14</b>
4.1	Experiments . . . . .	14
4.2	Evaluation and Comparison . . . . .	14
<b>5</b>	<b>Summary</b>	<b>14</b>
<b>A</b>	<b>Code</b>	<b>15</b>
<b>B</b>	<b>Notation</b>	<b>15</b>

## List of Figures

1	Samples, histogram and kernel density estimates . . . . .	5
2	Popular kernels . . . . .	6

## List of Tables

1	Popular kernels . . . . .	8
---	---------------------------	---

## List of Algorithms

1	Mean-shift algorithm in matrix form. . . . .	13
2	Mean-shift algorithm in iterative form. . . . .	13

## List of Code Listings

1	<code>mean_shift</code> . . . . .	16
2	<code>connected_component</code> . . . . .	17
3	<code>estimate_bandwidth</code> . . . . .	18

## Abstract

In this seminar paper we are describing and evaluating the mean shift algorithm. It is a nonparametric clustering method which does not depend on prior knowledge on the number and shape of the clusters.

The clusters are defined by the modes of the data density and a iteration procedure is used to let a point converge to its cluster centroid. We introduce the kernel density estimate for estimating the data density and then derive the mean shift algorithm from it.

After presenting the methods we also present different applications of the algorithm, mainly its use in image segmentation.

## 1 Introduction

In this seminar paper I am going to discuss the method of kernel density estimation (KDE) and its use in the mean-shift clustering algorithm (MSC). Clustering is one of the main tasks in Machine-Learning and MSC has many applications, mainly in image segmentation. The mean-shift algorithm is an non-parametric approach which allows a wide range of data distributions without imposing prior knowledge.

The content of this seminar paper is mainly based on Comaniciu & Meer (2002).

### History

The mean shift algorithm was introduced in Fukunaga & Hostetler (1975) where also the term “mean shift” was established. Comaniciu & Meer (2002) and Comaniciu et al. (2003) are responsible for a gain in interest for the algorithm. These publication popularized the use of the algorithm in image segmentation and tracking.

### Notation

In this paper I’m going to use the following conventions. Bold symbols  $\mathbf{x}$  denotes a vector.  $\{\mathbf{x}_i\}_{i=1}^n$  denotes  $n$  samples  $\mathbf{x}_i$ , indexed by  $i$ .  $d$  is the dimensionality of the samples.  $K(\cdot)$  is a multivariate kernel function  $\mathbb{R}^d \rightarrow \mathbb{R}$  and  $k(\cdot)$  denotes a single variable kernel function  $\mathbb{R} \rightarrow \mathbb{R}$ .

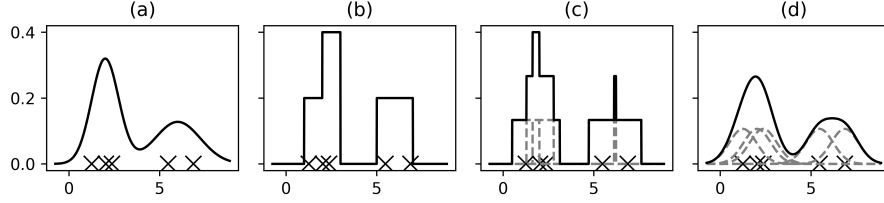


Figure 1: Samples, histogram and kernel density estimates. Author's illustration. (a) The probability density distribution  $p(x) = \frac{3}{5}N(x; 2, 0.75^2) + \frac{2}{5}N(x; 6, 1.25^2)$  and five samples drawn from it. (b) A histogram with bins  $[i, i + 1)$ ,  $i \in \mathbb{Z}$ . (c) A kernel density estimate with a uniform kernel and bandwidth  $h = 0.75$ ,  $k(x) = \frac{1}{2h}$  for  $|x| \leq h$  and 0 else. (d) A kernel density estimate with a gaussian kernel and bandwidth  $h = 0.75$ ,  $k(x) = \frac{1}{h} \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}|\frac{x}{h}|^2)$ .

## 2 Kernel Density Estimation

In this section we are going to introduce the kernel density estimate. It will be used in the derivation of the mean shift algorithm.

We assume that we have  $n$  data points  $\{\mathbf{x}_i\}_{i=1}^n \subset \mathbb{R}^d$ , which are independently, identically distributed samples from a unknown probability distribution  $p(\mathbf{x})$ ,  $\mathbf{x}_i \stackrel{i.i.d.}{\sim} p(\cdot)$ . We call the corresponding density function  $f$ .

The goal of density estimation is estimating the probability density of  $p$ , given samples  $\{\mathbf{x}_i\}_{i=1}^n$  in  $\mathbb{R}^d$ ,  $\mathbf{x} \sim p(\mathbf{x})$ .

Histograms estimate the probability distribution by bucketing the samples and counting the proportion of samples which fall in each bucket.

The **kernel density estimate** is defined as

$$\hat{f}(\mathbf{x}) = \frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right) \quad \mathbf{x} \in \mathbb{R}^d, \quad (1)$$

where  $n$  is the sample size,  $d$  is the dimensionality of the data,  $K(\mathbf{x})$  is a kernel function and  $h$  is a bandwidth parameter. The estimate is dependent on the choice of the kernel function and the bandwidth.

A **kernel function**  $K(\mathbf{x})$  is a function that satisfies  $K(\mathbf{x}) \geq 0$  and  $\int_{\mathbb{R}^d} K(\mathbf{x}) d\mathbf{x} = 1$ . These restrictions on the kernel ensure that  $\hat{f}(\mathbf{x})$  is a valid density. The **bandwidth**  $h$  controls the smoothness of the kernel density estimate. We are going

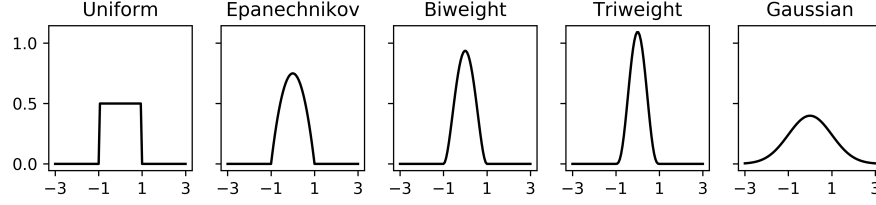


Figure 2: Popular kernels. Author’s illustration. The first four kernels all have support  $[-1, 1]$ . The gaussian kernel has the support  $(-\infty, \infty)$ .

to discuss the choices for kernel functions and the bandwidth in the next section.

The kernel density estimate can be interpreted as the sum of “bumps” centered on each data point.

For non-negative kernels, one can see that  $\hat{f}(\mathbf{x})$  is a valid probability density.

## 2.1 Other methods

Another method for estimating a probability density function are histograms. The samples  $\{\mathbf{x}_i\}_{i=1}^n$  are put into distinct buckets  $b_j \subseteq \mathbb{R}^d$  and the probability is estimated by the sample frequency  $p(j) = |\{\mathbf{x}_i \mid \mathbf{x}_i \in b_j\}|/n$  for each bucket.

## 2.2 Popular Kernels

A general Kernel is a function  $K : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfying the following properties. See (Comaniciu & Meer 2002) and (Wand & Jones 1995, p. 95).

$$\begin{aligned} \int_{\mathbb{R}^d} K(\mathbf{x}) d\mathbf{x} &= 1 & \int_{\mathbb{R}^d} \mathbf{x} K(\mathbf{x}) d\mathbf{x} &= \mathbf{0} \\ \lim_{\|\mathbf{x}\| \rightarrow \infty} \|\mathbf{x}\|^d K(\mathbf{x}) &= 0 & \int_{\mathbb{R}^d} \mathbf{x} \mathbf{x}^T K(\mathbf{x}) d\mathbf{x} &= c_K \mathbf{I} \end{aligned}$$

$$K^P(\mathbf{x}) = \prod_{j=1}^d k(x_j) \quad K^S(\mathbf{x}) = a_{k,d} k(\|\mathbf{x}\|)$$

Although there are many kernel functions, we are going to concentrate us on one class of functions: **radially symmetric kernels**. They are the kernels which can be formulated as

$$K(\mathbf{x}) = c_{k,d} k(\|\mathbf{x}\|^2), \quad (2)$$

where  $k : [0, \infty) \rightarrow [0, \infty)$  is called the **profile** of  $K$ , and  $c_{k,d}$  is a constant which ensures, that  $K(\mathbf{x})$  integrate to 1.

The derivation of the mean shift algorithm is going to be easier and nearly all popular kernels already belong to this class. The most popular kernels for mean shift clustering are the uniform, Epanechnikov and the Gaussian kernel. We are going to shortly introduce them and their profile.

The **uniform kernel** is defined as being constant in the unit hypersphere and can be written as

$$K_U(\mathbf{x}) = \begin{cases} \text{vol}(S_d)^{-1} & \text{for } \|\mathbf{x}\| \leq 1 \\ 0 & \text{else} \end{cases}, \quad (3)$$

where  $\text{vol}(S_d) = \pi^{d/2} \Gamma((d+2)/2)^{-1}$  is the volume of the d-dimensional hypersphere with  $\Gamma(\cdot)$  being the gamma function. If  $d = 1$  then  $\text{vol}(S_d) = 2$ . The corresponding profile  $k_U$  is

$$k_U(u) = \begin{cases} 1 & \text{for } 0 \leq u \leq 1 \\ 0 & \text{for } 1 < u. \end{cases} \quad (4)$$

The **Epanechnikov kernel** is defined as

$$K_E(\mathbf{x}) = \begin{cases} \frac{2+d}{2\text{vol}(S_d)} (1 - \|\mathbf{x}\|^2) & \text{for } \|\mathbf{x}\| \leq 1 \\ 0 & \text{else} \end{cases} \quad (5)$$

and its profile  $k_E$  is

$$k_E(u) = \begin{cases} 1 - u & \text{for } 0 \leq u \leq 1 \\ 0 & \text{for } 1 < u. \end{cases} \quad (6)$$

The **gaussian kernel** is defined as

$$K_G(\mathbf{x}) = (2\pi)^{-d/2} \exp\left(-\frac{1}{2} \|\mathbf{x}\|^2\right) \quad (7)$$

and its profile  $k_G$  is

$$k_G(u) = \exp\left(-\frac{1}{2}u\right). \quad (8)$$

Name	Formula	Support
Uniform	$k(x) = \frac{1}{2}$	$[-1, 1]$
Epanechnikov	$k(x) = \frac{3}{4}(1 - x^2)$	$[-1, 1]$
Biweight	$k(x) = \frac{15}{16}(1 - x^2)^2$	$[-1, 1]$
Triweight	$k(x) = \frac{35}{32}(1 - x^2)^3$	$[-1, 1]$
Gaussian	$k(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{x^2}{2})$	$(-\infty, \infty)$

Table 1: Popular kernels

## 2.3 Bandwidth Selection

One of the main challenges in using the kernel density estimator in practice is the choice of the bandwidth.

## 3 Mean Shift Clustering

Clustering is one of the main machine learning tasks, concerned with grouping objects  $\{x_i\}_{i=1}^n$  into clusters  $C_j$  resulting in a clustering  $\{C_j\}_{j=1}^K$ .

### 3.1 A general perspective on clustering algorithms

#### *introduction*

What does define a good clustering? One can try to minimize a global criterion like  $J = \sum_{j=1}^K \sum_{x_i \in C_j} \|x_i - m_j\|^2$ , like the  $K$ -Means algorithm, or minimize the intra-cluster min, max or average distances between the data points, like in hierarchical clustering. The notion of a “good” clustering brings rise to many different clustering algorithms.

Another criteria to differentiate clustering algorithms, are the assumptions that it makes on the data.  $K$ -Means assumes the data is separable into  $K$  clusters.

There are many different notions of what defines a good clustering and what shapes a cluster can take. These different notions result in different clustering algorithms. One class of clustering algorithms is centroid based clustering. Each cluster is represented by a representative called centroid. For example the representative of a  $K$ -Means clustering is the mean of the cluster  $m_j := \frac{1}{|C_j|} \sum_{x_i \in C_j} x_i$ .



*centroid based*

A common concept for several algorithms are **centroids**. A centroid is a point in the data space, which is a representative for a cluster.  $K$ -Means defines its centroids as the cluster mean  $\mathbf{m}_j = \frac{1}{|C_j|} \sum_{\mathbf{x}_i \in C_j} \mathbf{x}_i$ . A shortcoming of this definition is that the cluster mean is not guaranteed to “look like” a typical data object.

The mean shift algorithm defines its cluster representatives by the modes of the data density. This definition solves several of shortcomings of other algorithms. The algorithm does not make prior assum

Another way clustering algorithms can be differentiated is whether or not the amount of clusters is given as an input to the algorithm or not. For example the  $K$ -Means algorithm searches for exactly  $K$  clusters. But there are also other notions of a cluster center or representative. One is the added restriction, that the center is itself a data object or atleast “looks like” one. The latter is covered by using representatives which are likely also a data object. That’s the core of the mean shift clustering algorithm. It estimates the underlying density of the data objects and searches for points which high relative probability. Or in other words it identifies the modes of the estimated density.

That leaves open on how to estimate the density.

### 3.2 Density gradient estimation

Before looking at the algorithm itself, we are going to derive the name giving mean shift vector by estimating the density gradient with the gradient of the density estimate. This subsection is based on Comaniciu & Meer (2002) and follows its notation.

As we have discussed in the previous subsection, one way to define a cluster is by the modes of the the density. This is notion that is the main point of this paper. The task is to find the points  $\mathbf{x}$  where the density gradient vanishes,  $\nabla f(\mathbf{x}) = 0$ .

The density is estimated with kernel density estimation. Recall that the density estimate is formulated as

$$\hat{f}(\mathbf{x}) = \frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right). \quad (9)$$

A natural way to estimate the density gradient is to use the gradient of the density estimate. Formally we write

$$\hat{\nabla} f(\mathbf{x}) \equiv \nabla \hat{f}(\mathbf{x}) = \frac{1}{nh^d} \sum_{i=1}^n \nabla K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right). \quad (10)$$

To simplify the coming derivation we are going to concentrate on a class of kernels, **radially symmetric kernels**. These can be written as

$$K(\mathbf{x}) = c_{k,d} k(\|\mathbf{x}\|^2) \quad (11)$$

where  $k : [0, \infty] \rightarrow [0, \infty)$  is called the **profile** of  $K$ , and  $c_{k,d}$  is a constant which ensures that  $K(\mathbf{x})$  integrate to 1. The most common kernels belong to this class of functions.

Because we are going to use different kernel density estimates, we explicit state the used bandwidth  $h$  and kernel  $K$  as a subscript on the kernel density estimate  $\hat{f}_{h,K}(\mathbf{x})$ . The kernel density estimate can now be formulated as

$$\hat{f}_{h,K}(\mathbf{x}) = \frac{c_{k,d}}{nh^d} \sum_{i=1}^n k\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{h}\right\|^2\right), \quad (12)$$

and its gradient is

$$\nabla \hat{f}_{h,K}(\mathbf{x}) = \frac{2c_{k,d}}{nh^{d+2}} \sum_{i=1}^n (\mathbf{x} - \mathbf{x}_i) k'\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{h}\right\|^2\right). \quad (13)$$

With the substitution  $g(u) = -k'(u)$  and reordering we can reformulate the gradient to

$$\nabla \hat{f}_{h,K}(\mathbf{x}) = \frac{2c_{k,d}}{nh^{d+2}} \left[ \sum_{i=1}^n g\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{h}\right\|^2\right) \right] \left[ \frac{\sum_{i=1}^n \mathbf{x}_i g\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{h}\right\|^2\right)}{\sum_{i=1}^n g\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{h}\right\|^2\right)} - \mathbf{x} \right]. \quad (14)$$

The expression can be decomposed into several terms. The last term is called mean shift vector,

$$\mathbf{m}(\mathbf{x}) = \frac{\sum_{i=1}^n \mathbf{x}_i g\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{h}\right\|^2\right)}{\sum_{i=1}^n g\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{h}\right\|^2\right)} - \mathbf{x}. \quad (15)$$

The first two factors can be decomposed into a scaling factor and a kernel density estimate with the kernel  $G(\mathbf{x}) = c_{g,d}g(\|\mathbf{x}\|^2)$

$$\frac{2c_{k,d}}{nh^{d+2}} \left[ \sum_{i=1}^n g \left( \left\| \frac{\mathbf{x} - \mathbf{x}_i}{h} \right\|^2 \right) \right] = \frac{2c_{k,d}}{h^2 c_{g,d}} \left[ \frac{c_{g,d}}{nh^d} \sum_{i=1}^n g \left( \left\| \frac{\mathbf{x} - \mathbf{x}_i}{h} \right\|^2 \right) \right] \quad (16)$$

$$= \frac{2c_{k,d}}{h^2 c_{g,d}} \hat{f}_{h,G}(\mathbf{x}). \quad (17)$$

With these equations we can reformulate the density gradient estimate  $\nabla \hat{f}_{h,K}(\mathbf{x})$  into

$$\nabla \hat{f}_{h,K}(\mathbf{x}) = \frac{2c_{k,d}}{h^2 c_{g,d}} \hat{f}_{h,G}(\mathbf{x}) \mathbf{m}(\mathbf{x}). \quad (18)$$

That allows us to write the mean shift vector as a normalized gradient of the kernel density estimate  $\nabla \hat{f}_{h,K}(\mathbf{x})$

$$\mathbf{m}(x) = \frac{h^2 c_{g,d}}{2c_{k,d}} \frac{\nabla \hat{f}_{h,K}(\mathbf{x})}{\hat{f}_{h,G}(\mathbf{x})} \quad (19)$$

$$\begin{aligned} \nabla \hat{f}_{h,K}(\mathbf{x}) &= \frac{2c_{k,d}}{nh^{d+2}} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{x}) g \left( \left\| \frac{\mathbf{x} - \mathbf{x}_i}{h} \right\|^2 \right) \\ &= \frac{2c_{k,d}}{nh^{d+2}} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{x}) w_i \\ &= \frac{2c_{k,d}}{nh^{d+2}} \sum_{i=1}^n \mathbf{x}_i w_i - \mathbf{x} \sum_{i=1}^n w_i \\ &= \frac{2c_{k,d}}{nh^{d+2}} \left[ \sum_{i=1}^n w_i \right] \left[ \frac{\sum_{i=1}^n \mathbf{x}_i w_i}{\sum_{i=1}^n w_i} - \mathbf{x} \right] \\ &= \frac{2c_{k,d}}{nh^{d+2}} \left[ \sum_{i=1}^n g \left( \left\| \frac{\mathbf{x} - \mathbf{x}_i}{h} \right\|^2 \right) \right] \left[ \frac{\sum_{i=1}^n \mathbf{x}_i g \left( \left\| \frac{\mathbf{x} - \mathbf{x}_i}{h} \right\|^2 \right)}{\sum_{i=1}^n g \left( \left\| \frac{\mathbf{x} - \mathbf{x}_i}{h} \right\|^2 \right)} - \mathbf{x} \right] \\ &= \frac{2c_{k,d}}{h^2 c_{g,d}} \left[ \frac{c_{g,d}}{nh^d} \sum_{i=1}^n g \left( \left\| \frac{\mathbf{x} - \mathbf{x}_i}{h} \right\|^2 \right) \right] \mathbf{m}(\mathbf{x}) \\ &= \frac{2c_{k,d}}{h^2 c_{g,d}} \hat{f}_{h,G}(\mathbf{x}) \mathbf{m}(\mathbf{x}) \end{aligned}$$

$$\mathbf{m}(x) = \frac{h^2 c_{g,d}}{2c_{k,d}} \frac{\nabla \hat{f}_{h,K}(\mathbf{x})}{\hat{f}_{h,G}(\mathbf{x})} \quad (20)$$

Name	Support	$k(u)$	$-k'(u)$	$K(\mathbf{x})$
Epanechnikov	$[0, 1]$	$1 - u$	1	$\frac{1}{2}c_d^{-1}(d+2)(1 - \ \mathbf{x}\ ^2)$
Gaussian	$(-\infty, \infty)$	$\exp(-\frac{1}{2}u)$	$\frac{1}{2}\exp(-\frac{1}{2}u)$	$(2\pi)^{-d/2}\exp(-\frac{1}{2}\ \mathbf{x}\ ^2)$

$$k_N(u) = \exp\left(-\frac{1}{2}u\right) \quad (21)$$

$$-k'_N(u) = \frac{1}{2}\exp\left(-\frac{1}{2}u\right) \propto k(u) \quad (22)$$

$$K_N(\mathbf{x}) = (2\pi)^{-d/2}\exp\left(-\frac{1}{2}\|\mathbf{x}\|^2\right) \quad (23)$$

$$k_E(u) = \begin{cases} 1 - u & 0 \leq u \leq 1 \\ 0 & 1 < u \end{cases} \quad (24)$$

$$-k'_E(u) = \begin{cases} 1 & 0 \leq u \leq 1 \\ 0 & 1 < u \end{cases} \quad (25)$$

$$K_E(\mathbf{x}) = \begin{cases} \frac{1}{2}c_d^{-1}(d+2)(1 - \|\mathbf{x}\|^2) & \|\mathbf{x}\| \leq 1 \\ 0 & \text{else} \end{cases} \quad (26)$$

### 3.3 The algorithm

The core of the mean shift algorithm is an iterative procedure which acts on a single point  $\mathbf{x}^{(0)}$ . For the given point the algorithm finds the corresponding mode of the kernel density estimate ( $\hat{\nabla}f(\mathbf{x}) = \mathbf{0}$ ). Each iteration step the point moves in the direction of the steepest ascent. The iteration stops when the point has converged to a fixed point  $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)}$ . For the clustering of a dataset, the iteration procedure is performed on each point. The samples used for the kernel density estimate are still fixed to the original dataset.

A variant of the algorithm is the blurring mean shift algorithm. There the kernel density estimate is based on the data points from the previous iteration step.

$$\mathbf{m}(\mathbf{x}) = \frac{\sum_{i=1}^n \mathbf{x}_i g\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{h}\right\|^2\right)}{\sum_{i=1}^n g\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{h}\right\|^2\right)} - \mathbf{x} \quad (27)$$

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \mathbf{m}(\mathbf{x}^{(t)}), \quad t = 1, 2, \dots \quad (28)$$

---

**Algorithm 1** Mean-shift algorithm in matrix form.

---

```

1: function MEANSHIFT( $\{\mathbf{x}_i\}_{i=1}^n \subset \mathbb{R}^d, h : \mathbb{R}^+, \varepsilon : \mathbb{R}^+$ )
2:    $\mathbf{Z} = \mathbf{X} : \mathbb{R}^{d \times n}$ 
3:   repeat
4:      $\mathbf{W} = (\exp(-\frac{1}{2} \|(\mathbf{x}_i - \mathbf{x}_j)/h\|^2))_{i,j=1..n}$ 
5:      $\mathbf{D} = \text{diag}(\sum_i \mathbf{W}_{ij})$ 
6:      $\mathbf{Q} = \mathbf{W} \mathbf{D}^{-1} : \mathbb{R}^{n \times n}$ 
7:      $\mathbf{Z} = \mathbf{X} \mathbf{Q} : \mathbb{R}^{d \times n}$ 
8:   until stop
9:   return CONNECTEDCOMPONENTS( $\{\mathbf{z}_i\}_{i=1}^n, \varepsilon$ )

```

---



---

**Algorithm 2** Mean-shift algorithm in iterative form.

---

```

1: function MEANSHIFT( $\mathbf{x}_1, \dots, \mathbf{x}_n : \mathbb{R}^d, h, \varepsilon$ )
2:   for  $i = 1..n$  do
3:      $\mathbf{x}^{(1)} = \mathbf{x}_i, t = 0$ 
4:     repeat
5:        $t = t + 1$ 
6:        $\forall n : p(i \mid \mathbf{x}^{(t)}) = \frac{\exp(-\frac{1}{2} \|(\mathbf{x}^{(t)} - \mathbf{x}_i)/h\|^2)}{\sum_{j=1}^n \exp(-\frac{1}{2} \|(\mathbf{x}^{(t)} - \mathbf{x}_j)/h\|^2)}$ 
7:        $\mathbf{x}^{(t+1)} = \sum_{i=1}^n p(i \mid \mathbf{x}^{(t)}) \mathbf{x}_i$ 
8:       until  $\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\| < tol$ 
9:        $\mathbf{z}_i = \mathbf{x}^{(t+1)}$ 
10:   return CONNECTEDCOMPONENTS( $\{\mathbf{z}_i\}_{i=1}^n, \varepsilon$ )

```

---

### 3.4 Discussion

#### *Problems*

##### *Performance*

The algorithm inhabits several problems. Some of them are dependent on the domain and data, like the choice of bandwidth and kernel, others are regarding the performance.

##### Methods for speedup

## 4 Application

### *image segmentation*

The main application for the mean shift clustering is image segmentation. Although the raw image data in the RGB-colorspace (red, green, blue) can be used, it is often not disereable. A human percieves distances of colors differently than the RGB-colorspace indicates. For that reason the image is often transformed in a more suitable colorspace, for example the L\*u\*v\*-colorspace.

Another aspect is the spatial coherence of the clusters in the image. Two pixels with the same color, but in widely different parts of the image would be put in the same cluster. That can be prevented by adding imagespace information to the pixels.

### 4.1 Experiments

#### *basic clustering task*

*image segmentation – details*

*dataset description*

Alpert et al. (2012)

### 4.2 Evaluation and Comparison

## 5 Summary

We introduced general kernel density estimates and presented popular choices for kernels. We classify the mean shift clustering algorithm by it's understanding of a cluster and derive the iteration step coming from the kernel density estimate. We present different applications of the algorithm, mainly classic clustering tasks and

## References

- Alpert, S., Galun, M., Brandt, A. & Basri, R. (2012), ‘Image segmentation by probabilistic bottom-up aggregation and cue integration’, *IEEE Transactions on Pattern Analysis and Machine Intelligence* **34**(2), 315–326.
- Comaniciu, D. & Meer, P. (2002), ‘Mean shift: a robust approach toward feature space analysis’, *IEEE Transactions on Pattern Analysis and Machine Intelligence* **24**(5), 603–619.
- Comaniciu, D., Ramesh, V. & Meer, P. (2003), ‘Kernel-based object tracking’, *IEEE Transactions on Pattern Analysis and Machine Intelligence* **25**(5), 564–577.
- Fukunaga, K. & Hostetler, L. (1975), ‘The estimation of the gradient of a density function, with applications in pattern recognition’, *IEEE Transactions on Information Theory* **21**(1), 32–40.
- Wand, M. P. & Jones, M. C. (1995), *Kernel Smoothing*, Vol. 60 of *Monographs on Statistics and Applied Probability*, Springer US, Boston, MA and s.l.  
**URL:** <http://dx.doi.org/10.1007/978-1-4899-4493-1>

## A Code

## B Notation

Data

$n$       vs.       $N$

$x$       vs.       $\mathbf{x}$

$x_i$       vs.       $x^{(i)}$

Kernels

$\hat{f}(x)$       vs.       $\hat{p}(x)$

```

1 function [A, C, T] = mean_shift3(X, kernel, h, epsilon)
2     % Performes the mean shift algorithm with the
   % specified bandwidth h.
3     %
4     % INPUTS:
5     %     X – an d by n array of data points
6     %     h – the bandwidth used
7     %     epsilon – the parameter for cluster pruning
8     % OUTPUTS:
9     %     A – an n by 1 array of cluster indices
10    %     C – an d by n array of cluster centroids
11    %     T – an n by 1 array of iterations needed
   % for each data point
12    tol = h * 1e-2;
13    max_iter = 100;
14
15    [d, n] = size(X);
16    Z = zeros(d, n);
17    T = zeros(n, 1);
18
19    for i = 1:n
20        x = X(:, i);
21        for t = 1:max_iter
22            switch kernel
23                case 'gaussian'
24                    w = exp(-0.5 * sum(((X - x) / h).^2,
   % 1));
25                case 'uniform'
26                    w = sum(((X - x) / h).^2, 1) < 1;
27            end
28
29            x_next = X * w' / sum(w);
30
31            if (norm(x_next - x) < tol)
32                break
33            end
34            x = x_next;
35        end
36        Z(:, i) = x;
37        T(i) = t;
38    end
39
40    [A, C] = connected_component(Z, epsilon);
41 end

```

Code Listing 1: Mean shift



```

1 function [A, C] = connected_component(X, epsilon)
2     % Performes data pruning. An implicit neighborhood
3     % graph is constructed in which x_i and x_j are
4     % connected if  $||x_i - x_j|| < \text{epsilon}$ . Then points
5     % are randomly chosen and all its neighbors are
6     % assigned to the same cluster.
7     %
8     % INPUTS:
9     %       X - an d by n array with data points
10    %       epsilon - the distance in which points are merged
11    % OUTPUTS:
12    %       A - an n by 1 array with cluster indices
13    %       C - an d by max(A) array with cluster
14    %       representatives
15    [d, n] = size(X);
16    A = zeros(n, 1);
17    C = zeros(d, n);
18    % the number of connected components found yet
19    c = 0;
20
21    while any(A == 0)
22        c = c + 1;
23        % find i for which x_i has not yet been assigned
24        i = find(~A, 1);
25        x = X(:, i);
26        C(:, c) = x;
27        d = sum((X - x).^2, 1).^0.5;
28        A(d < epsilon) = c;
29    end
30    C = C(:, 1:c);
31 end

```

Code Listing 2: Connected component

```

1 function h = estimate_bandwidth(X, quantile)
2     % Bandwidth estimation for kernel density estimates.
3     % The mean distance
4     % of each data point to its (quantile * n)-nearest
5     % neighbor is used. A
6     % quantile value of 0.5 means the median value of all
7     % pairwise
8     % distances is used.
9     %
10    % INPUTS:
11    %     X - an d by n array with data points
12    %     quantile - the quantile used; must be in [0, 1]
13    %     the default value is 0.3
14    % OUTPUTS:
15    %     h - the estimated bandwidth
16    %
17    % Reference:
18    % See: scikit-learn implementation: https://github.com/scikit-learn/scikit-learn/blob/7389dba/sklearn/cluster/mean\_shift\_.py#L31
19    % Example:
20    % X = [randn(3, 20) (randn(3, 40) + [2;0;6])];
21    % estimate_bandwidth(X)
22    if nargin < 2
23        quantile = 0.3;
24    end
25    [~, n] = size(X);
26    K = ceil(n * quantile);
27    [~, D] = knnsearch(X', X', 'K', K, 'Distance', 'euclidean');
28    h = mean(D(:, K));
29 end

```

Code Listing 3: Estimate bandwidth

$$\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^n K_h(x - x_i)$$

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right) \quad \text{vs.} \quad \hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$

Radial Kernel

$$\hat{f}_h(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n K_h(\|\mathbf{x} - \mathbf{x}_i\|)$$

Product Kernel

$$\hat{f}_h(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \prod_{j=1}^d K_{h_j}(\mathbf{x}_j - \mathbf{x}_{ij}) \quad \text{vs.} \quad \hat{f}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \prod_{j=1}^d K_{h_j}(\mathbf{x}_j - \mathbf{x}_j^{(i)})$$

## Statutory Declaration

I hereby declare, that I have written this seminar paper with the title **Kernel Density Estimates and Mean-Shift Clustering** by myself and have not used other sources without declaration.

Karlsruhe, January 20, 2019

---

Jonas Spinner