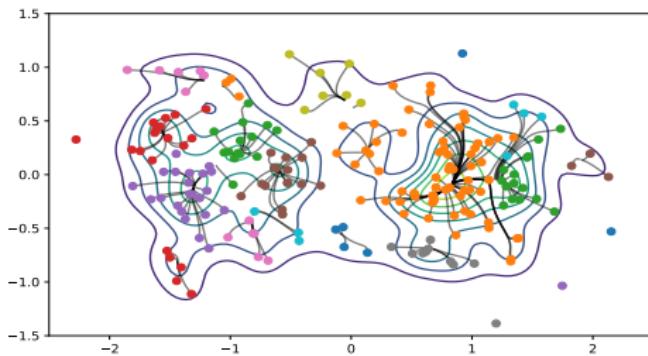
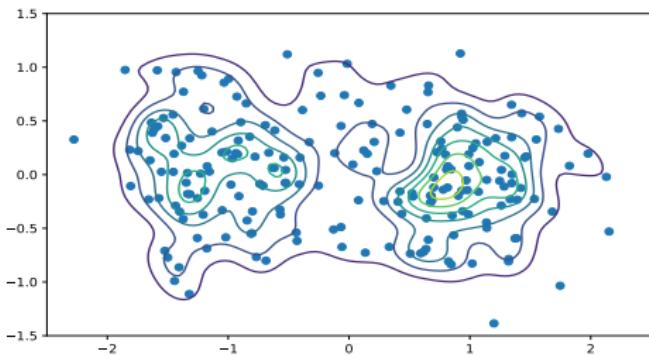


Kernel Density Estimates and Mean Shift Clustering

Jonas Spinner | February 4, 2019

ANALYTICS AND STATISTICS AT THE INSTITUTE OF OPERATIONS RESEARCH



Outline

1 Introduction

2 Kernel Density Estimates

- Definition
- Kernel functions
- Bandwidth

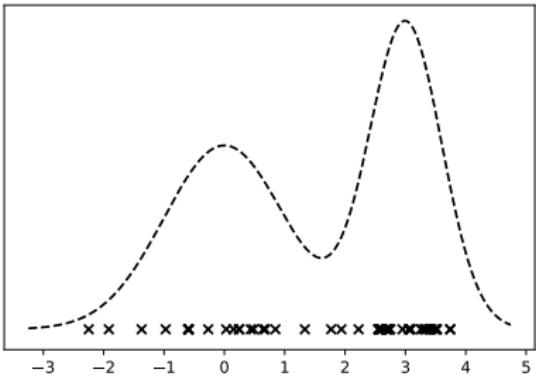
3 Mean Shift Clustering

- Algorithm
- Convergence
- Connection to KDE
- Bandwidth effects
- Speedup methods and discussion

4 Application

- Image segmentation

Density estimation



- **Task:** Given samples $x_1, \dots, x_n \in \mathbb{R}^d$, estimate the underlying density.

The kernel density estimate

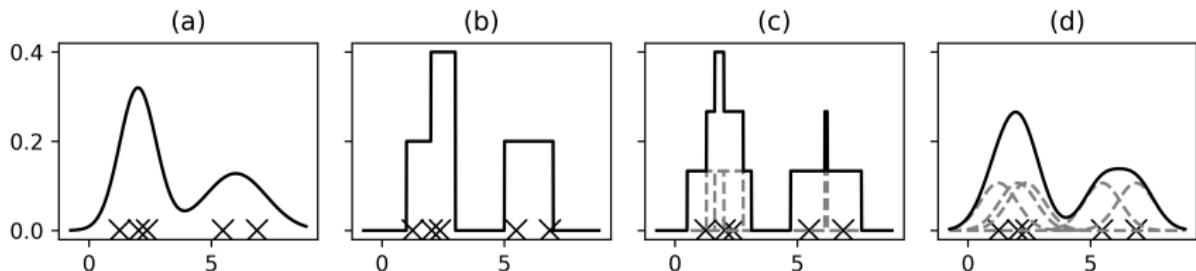
The **kernel density estimate** is defined as

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

- K is a *kernel function*.
- h is a *bandwidth parameter*.
- When $\int_{\mathbb{R}^d} K(\mathbf{x}) d\mathbf{x} = 1$ and $K(\mathbf{x}) \geq 0$ for all \mathbf{x} then $\hat{f}(\mathbf{x})$ is a valid probability density function.

The kernel density estimate

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



- (a) Underlying density.
- (b) Histogram.
- (c) KDE with uniform kernel.
- (d) KDE with gaussian kernel.

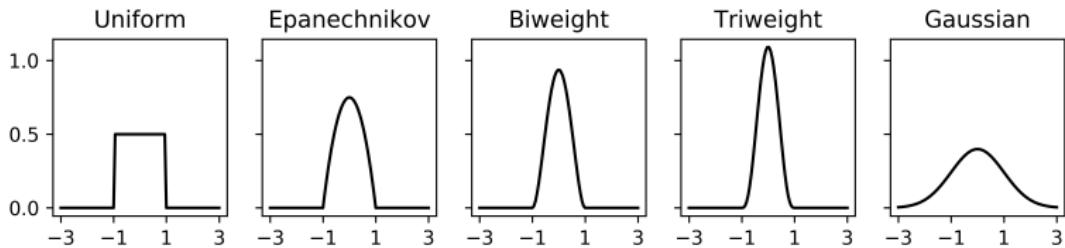
Popular kernel functions

Radially symmetric kernel functions are kernel functions which can be represented as

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

- $k(u) : [0, \infty) \rightarrow [0, \infty)$ is called the **profile** of K .
- For example the gaussian kernel has the profile $k(u) = \exp(-\frac{1}{2}u)$.
- Nearly all popular kernel belong to this class of kernels.

Popular kernel functions



| Name | Profile support | $k(u)$ | $-k'(u)$ | $K(x)$ |
|--------------|---------------------|----------------------------------|--|--|
| Uniform | $u \in [0, 1]$ | 1 | 0 | $\text{vol}(S_d)^{-1}$ |
| Epanechnikov | $u \in [0, 1]$ | $1 - u$ | 1 | $\frac{1}{2} \text{vol}(S_d)^{-1} (d+2) (1 - \ x\ ^2)$ |
| Biweight | $u \in [0, 1]$ | $(1-u)^2$ | $2(1-u)$ | $\propto (1 - \ x\ ^2)^2$ |
| Triweight | $u \in [0, 1]$ | $(1-u)^3$ | $3(1-u)^2$ | $\propto (1 - \ x\ ^2)^3$ |
| Gaussian | $u \in [0, \infty)$ | $\exp\left(-\frac{1}{2}u\right)$ | $\frac{1}{2} \exp\left(-\frac{1}{2}u\right)$ | $(2\pi)^{-d/2} \exp\left(-\frac{1}{2}\ x\ ^2\right)$ |

Bandwidth

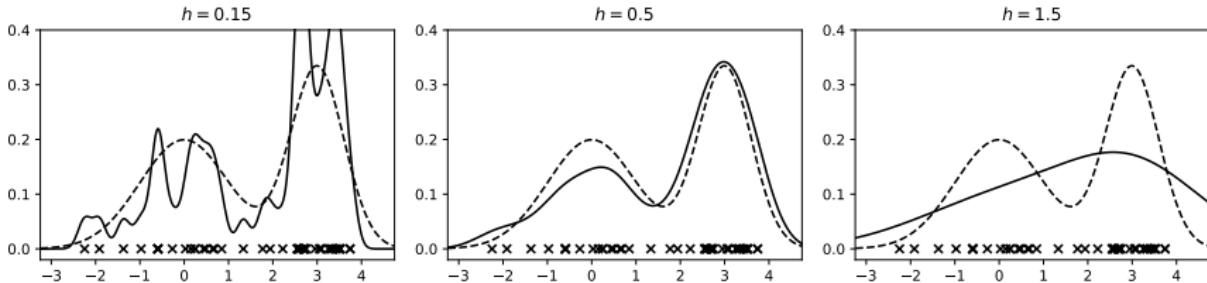


Abbildung: Gaussian kernel, $n = 50$.

- The choice of bandwidth is a bias-variance tradeoff for the estimate $\hat{f}(\mathbf{x})$.
- A small bandwidth results in high variance, a large bandwidth introduces a bias.

Bandwidth

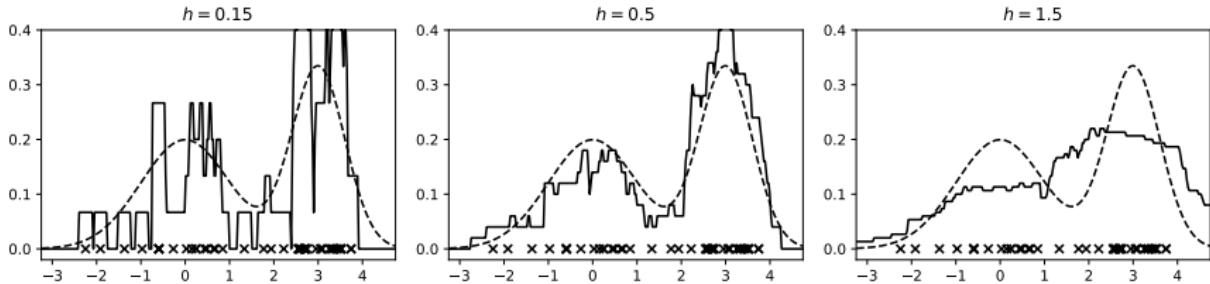


Abbildung: Uniform kernel, $n = 50$.

- The choice of bandwidth is a bias-variance tradeoff for the estimate $\hat{f}(\mathbf{x})$.
- A small bandwidth results in high variance, a large bandwidth introduces a bias.

Bandwidth

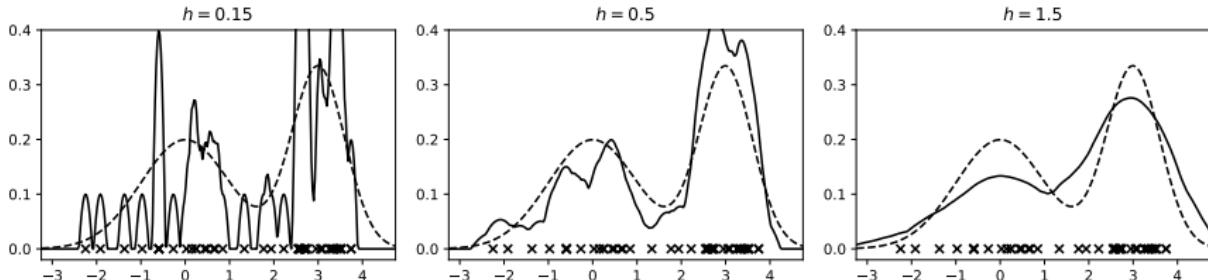
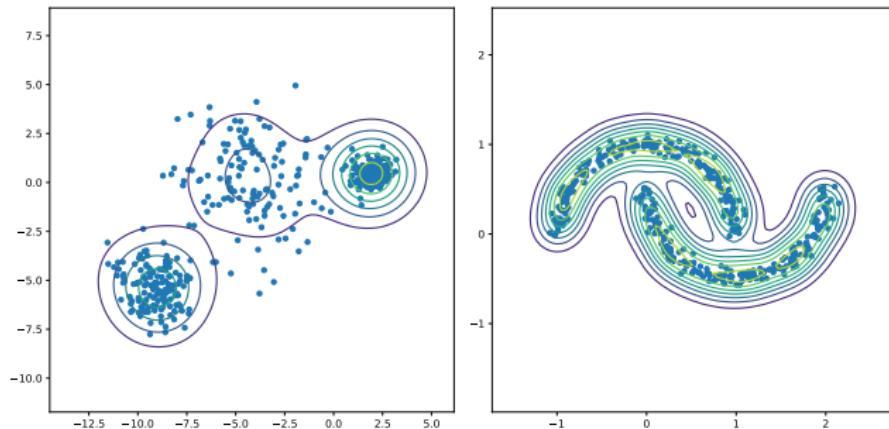


Abbildung: Epanechnikov kernel, $n = 50$.

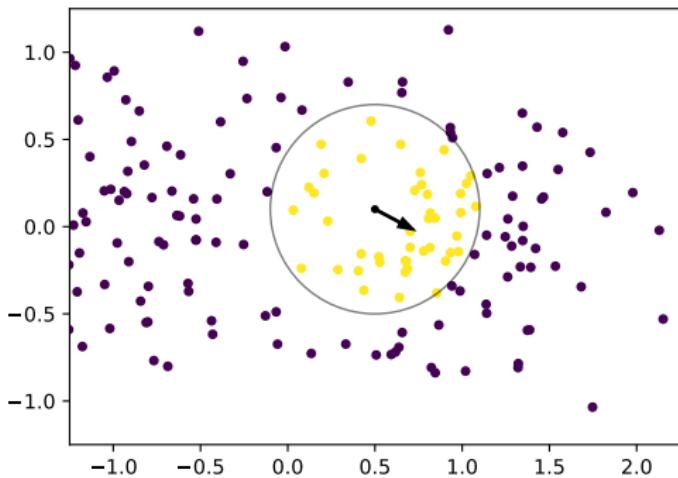
- The choice of bandwidth is a bias-variance tradeoff for the estimate $\hat{f}(x)$.
- A small bandwidth results in high variance, a large bandwidth introduces a bias.

Mean shift clustering



- Modes of the kernel density estimate can be used to identify clusters.
- Mode finding procedure results in a non-parametric clustering algorithm.

The idea



- Iteratively move point to “high-density” region.
- Direction is determined by local neighborhood.

Locally weighted mean

The weighted mean is

$$\mu^* = \frac{\sum_{i=1}^n x_i w_i}{\sum_{i=1}^n w_i}$$

The **locally weighted mean** has weights depending on the distance to a data point x

$$\mu^*(x) = \frac{\sum_{i=1}^n x_i g\left(\left\|\frac{x-x_i}{h}\right\|^2\right)}{\sum_{i=1}^n g\left(\left\|\frac{x-x_i}{h}\right\|^2\right)}$$

Mean shift procedure

The **mean shift vector** is defined as

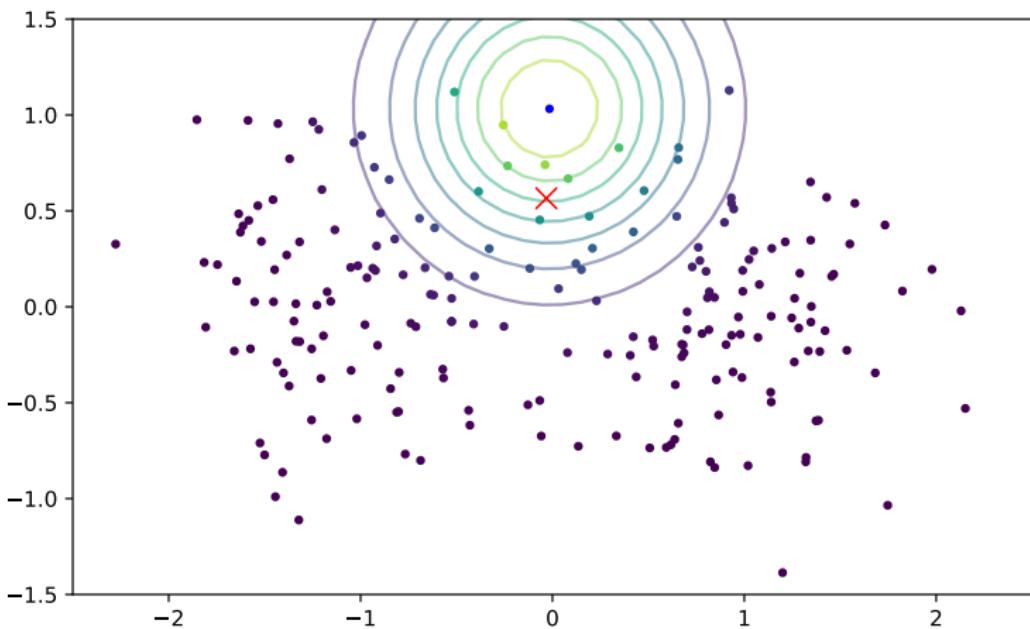
$$\mathbf{m}(\mathbf{x}) = \mu^*(\mathbf{x}) - \mathbf{x}$$

The **mean shift procedure** is

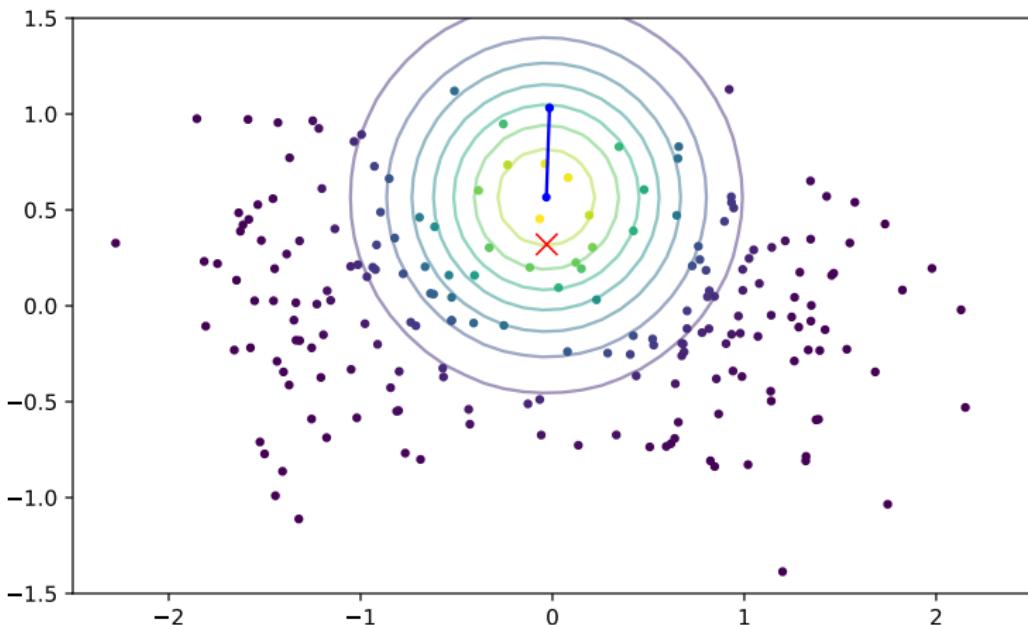
$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \mathbf{m}(\mathbf{x}^{(t)}) \quad \text{for } i = 1, 2, \dots$$

- When the process converges we assign \mathbf{x} to the mode $\mathbf{x}^{(\infty)}$.
- Points at the same mode are considered to be in the same cluster.

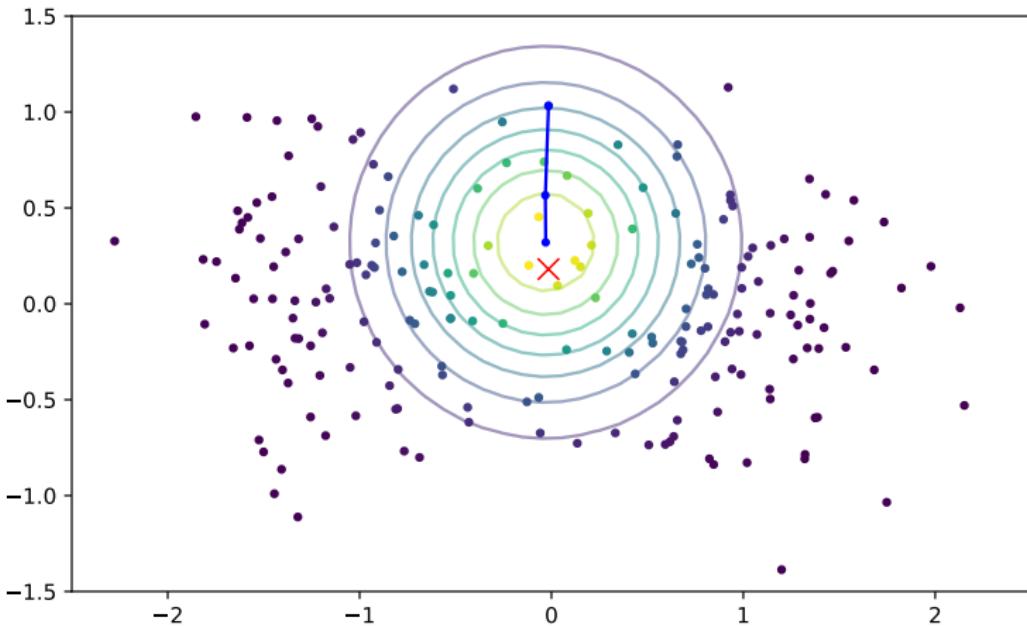
Convergence – Gaussian kernel



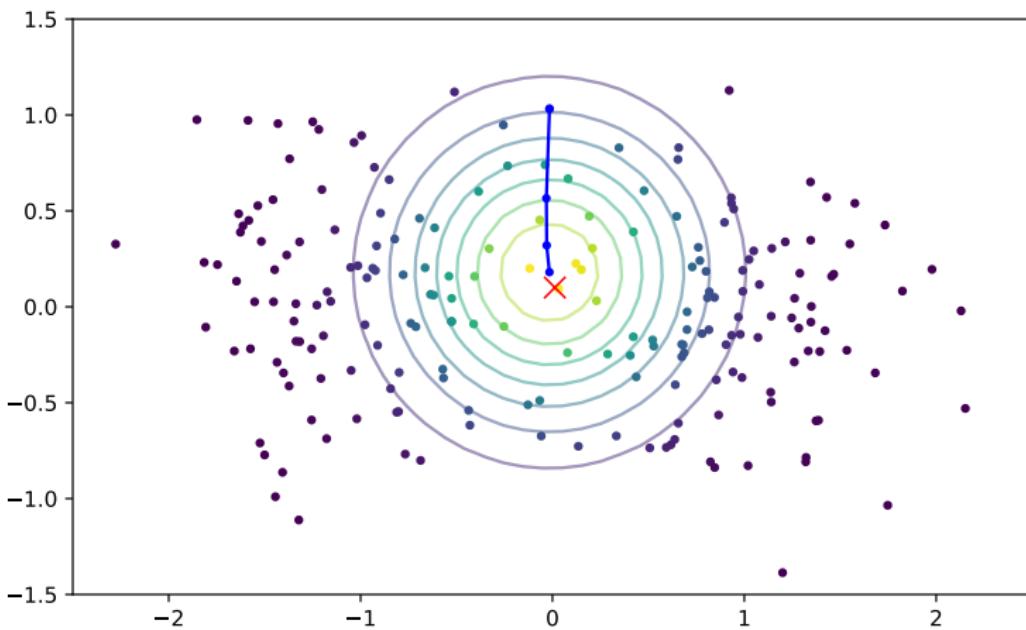
Convergence – Gaussian kernel



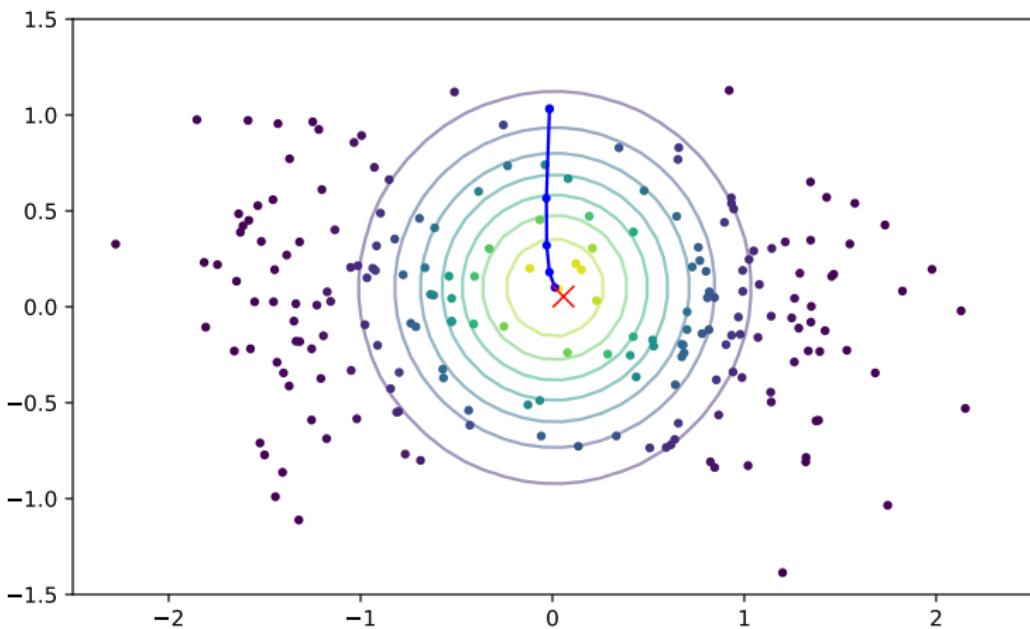
Convergence – Gaussian kernel



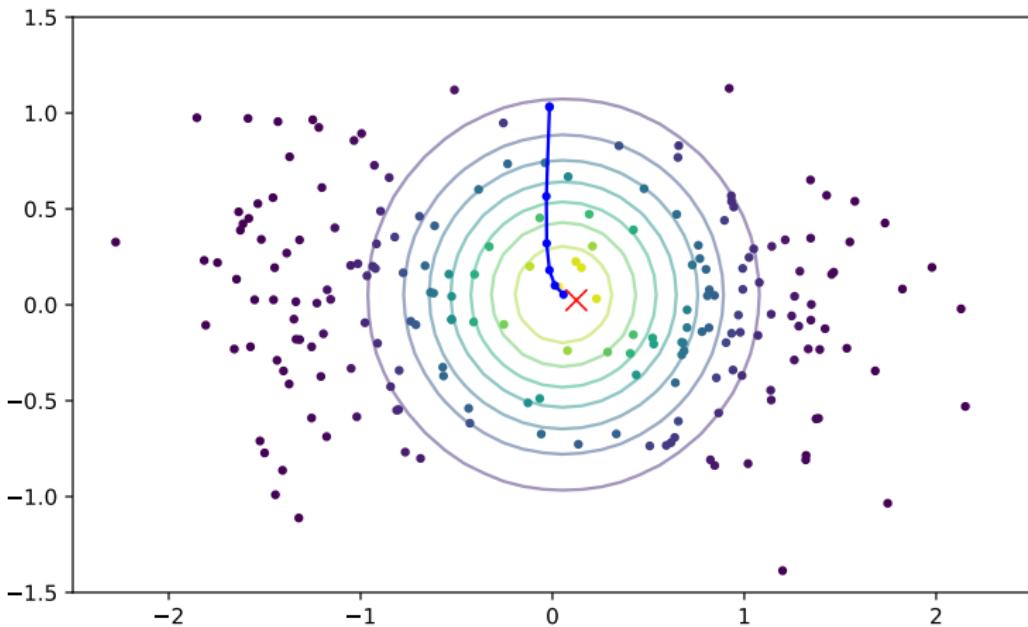
Convergence – Gaussian kernel



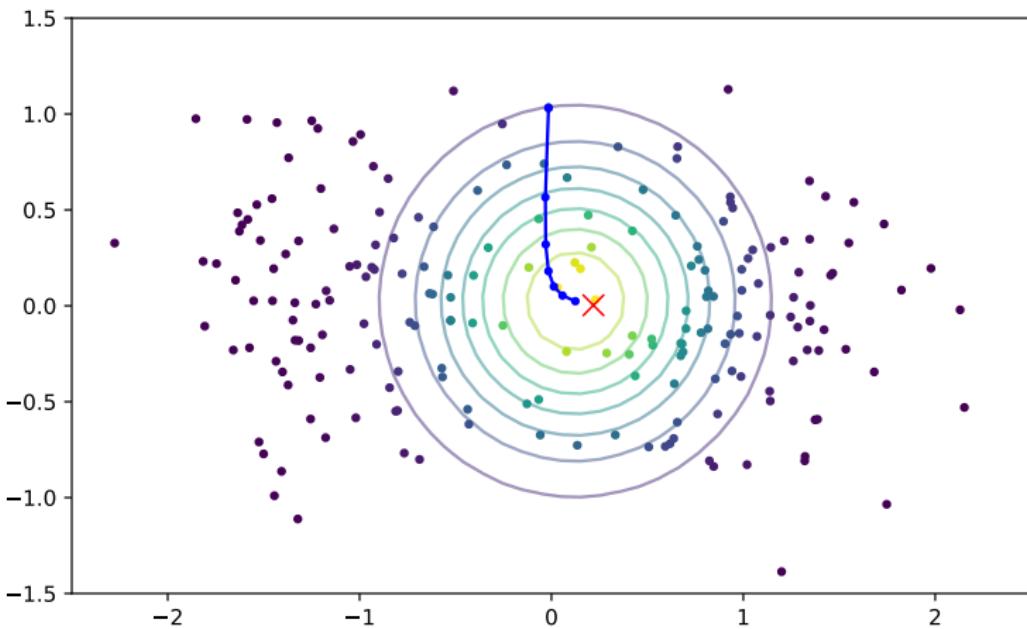
Convergence – Gaussian kernel



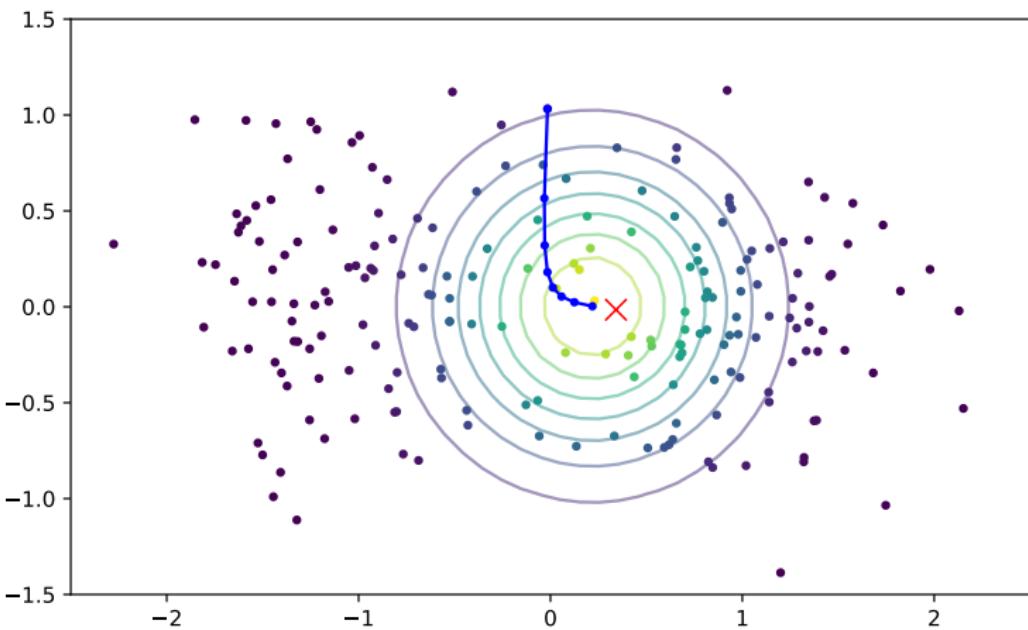
Convergence – Gaussian kernel



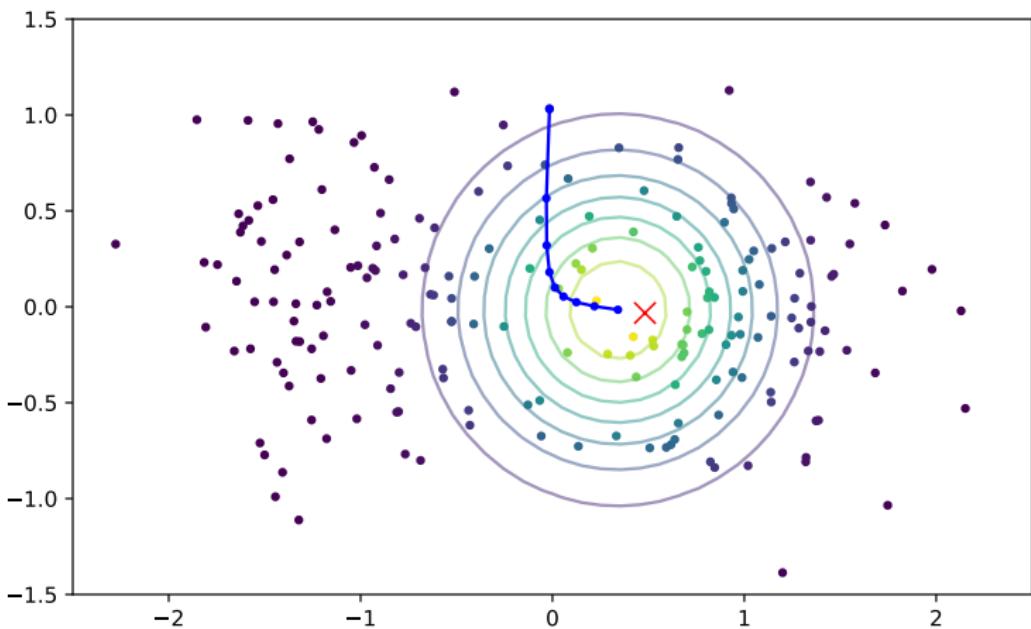
Convergence – Gaussian kernel



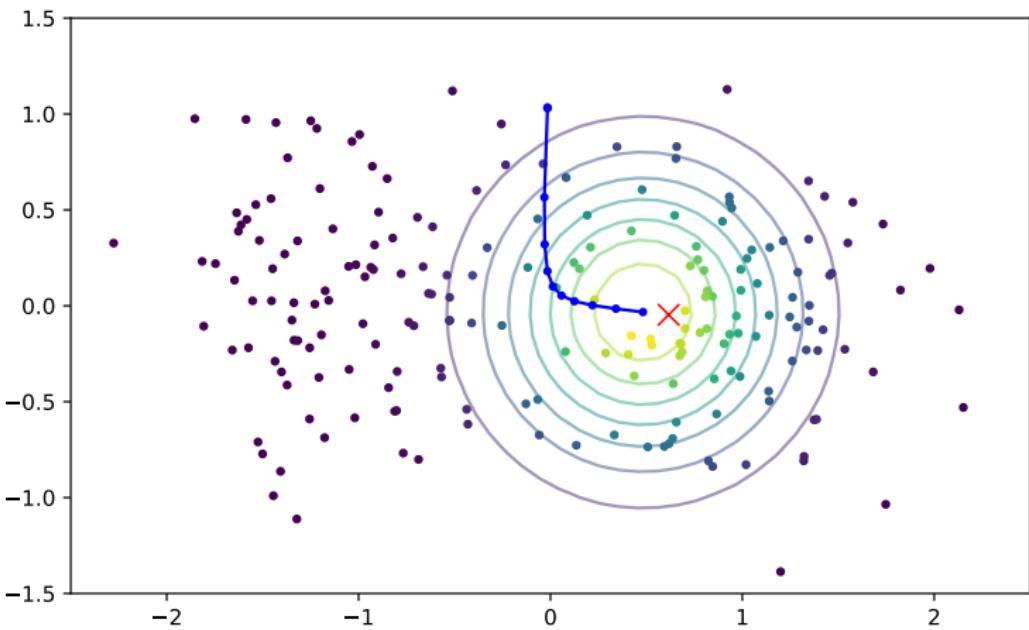
Convergence – Gaussian kernel



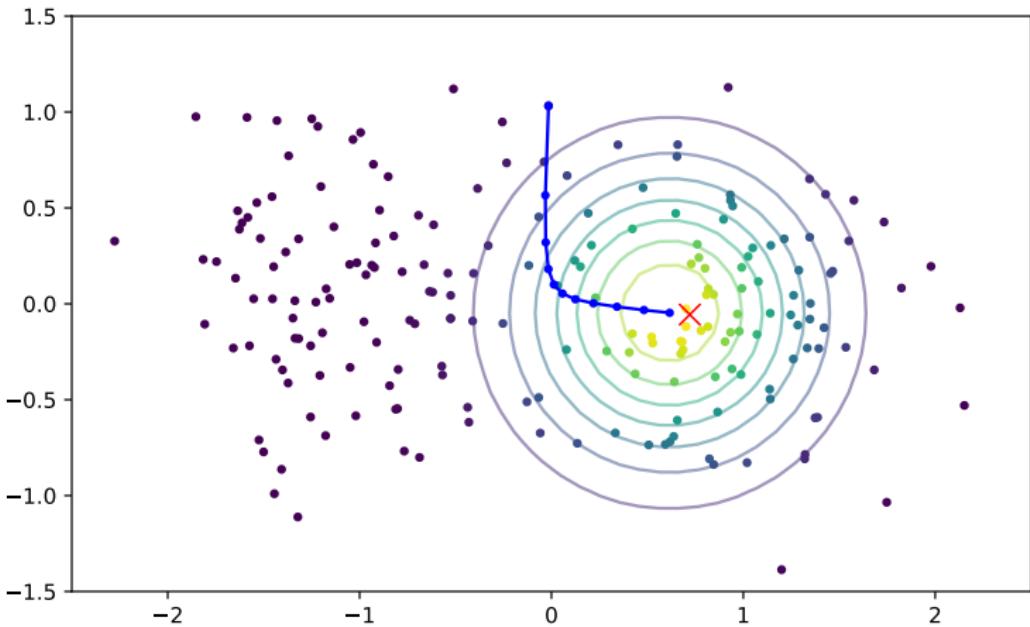
Convergence – Gaussian kernel



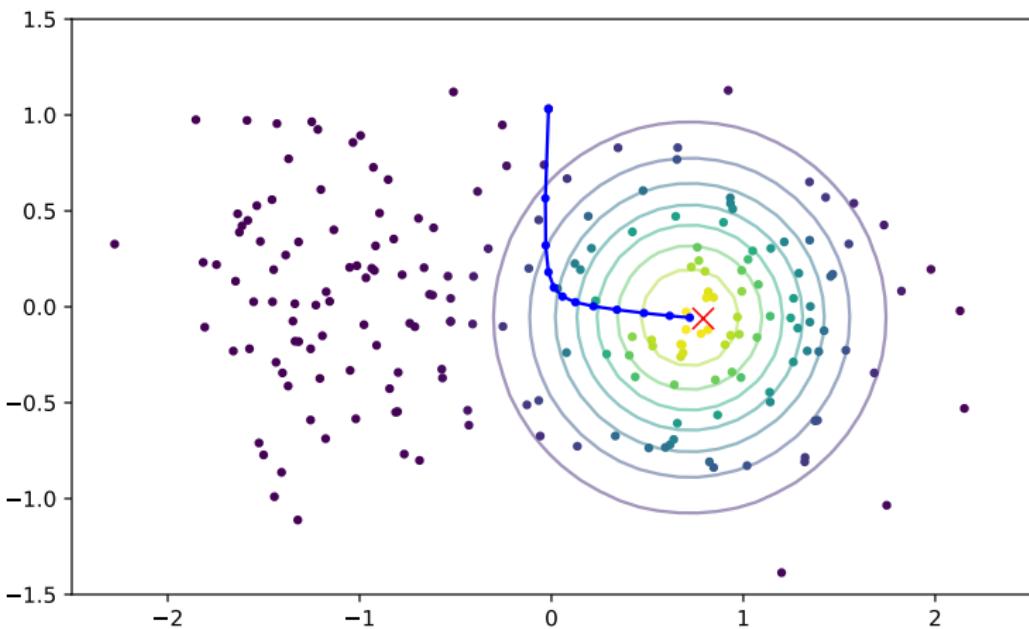
Convergence – Gaussian kernel



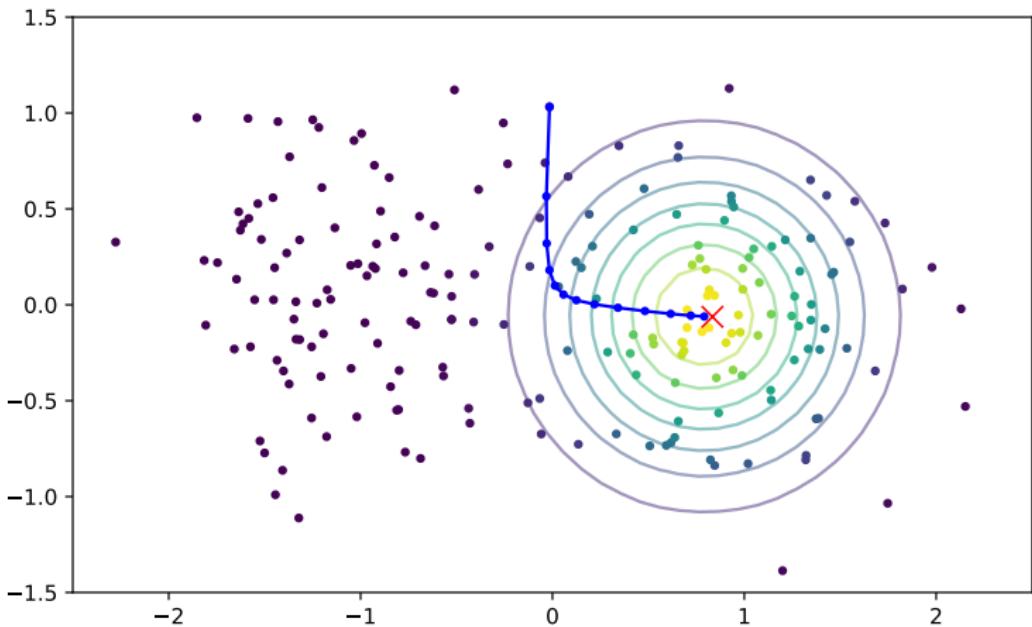
Convergence – Gaussian kernel



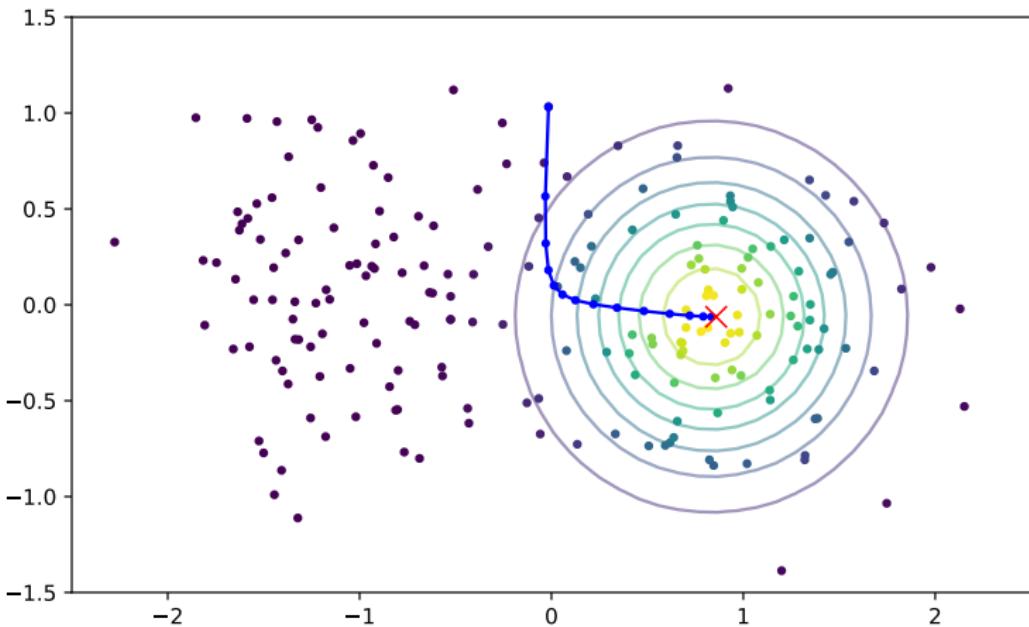
Convergence – Gaussian kernel



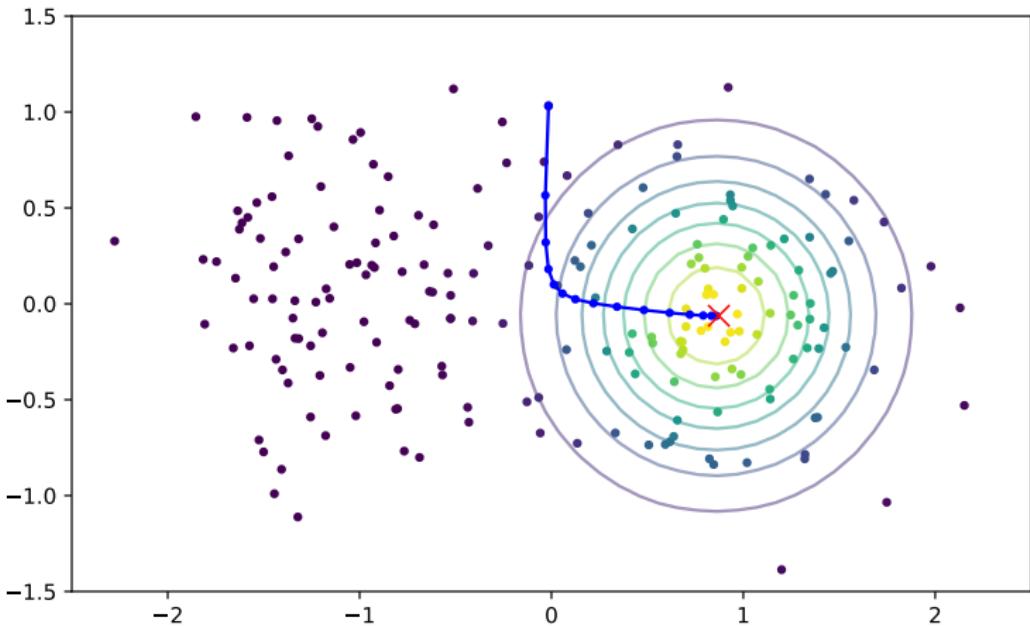
Convergence – Gaussian kernel



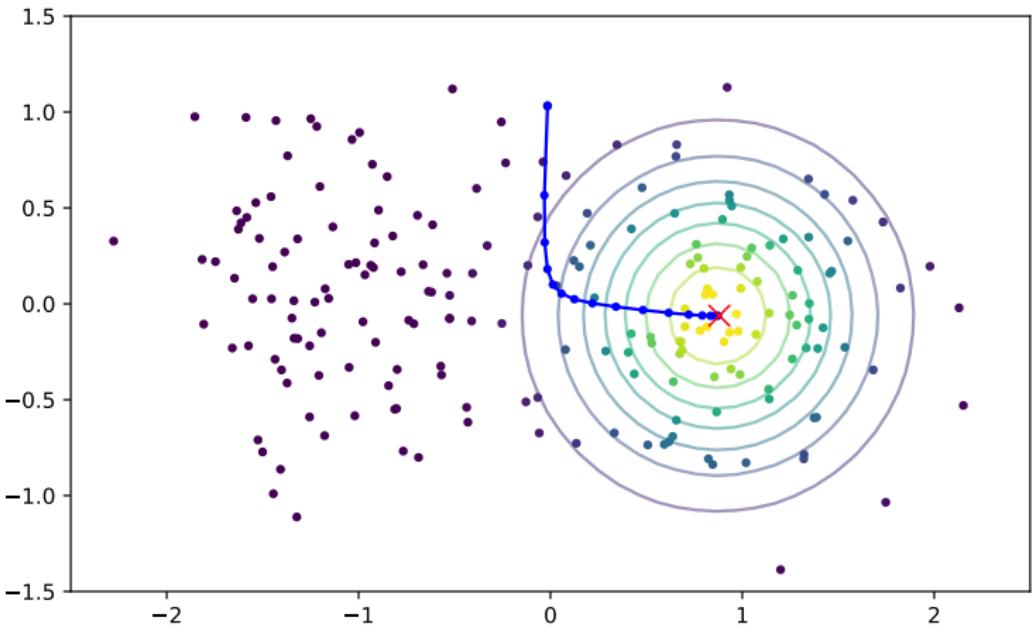
Convergence – Gaussian kernel



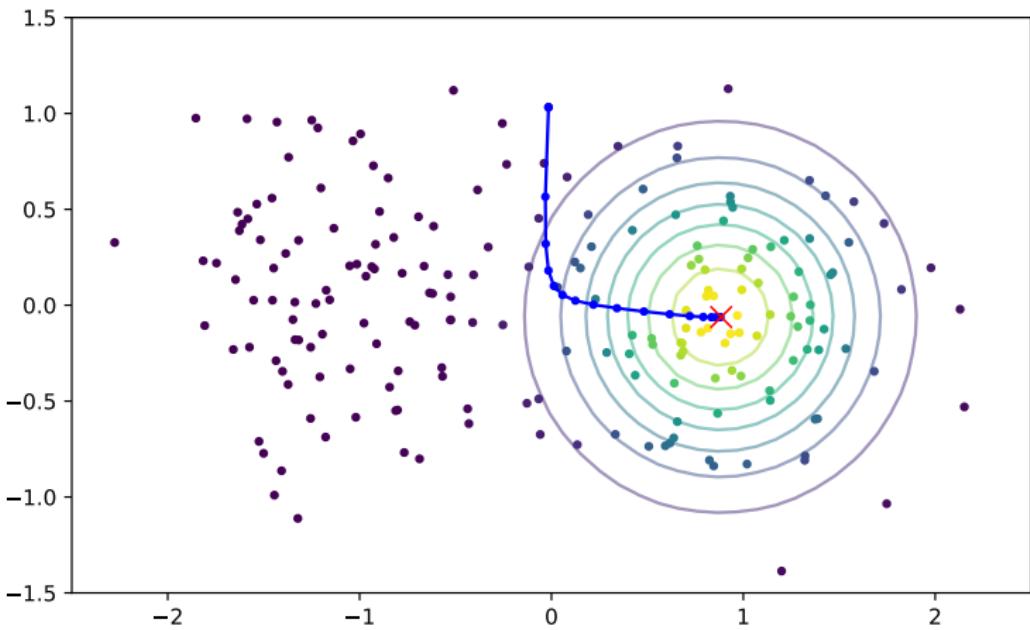
Convergence – Gaussian kernel



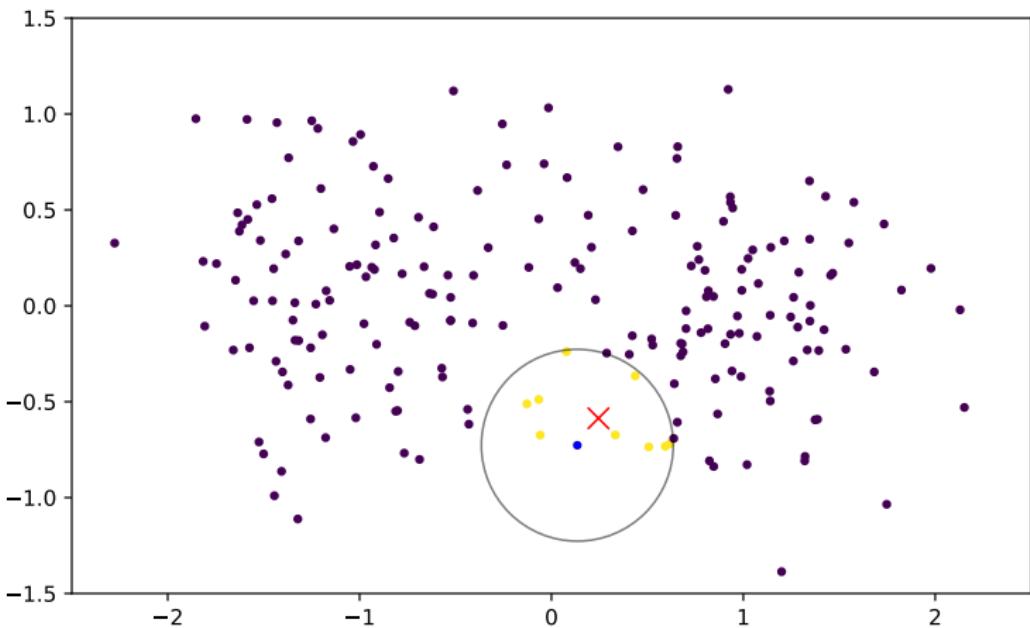
Convergence – Gaussian kernel



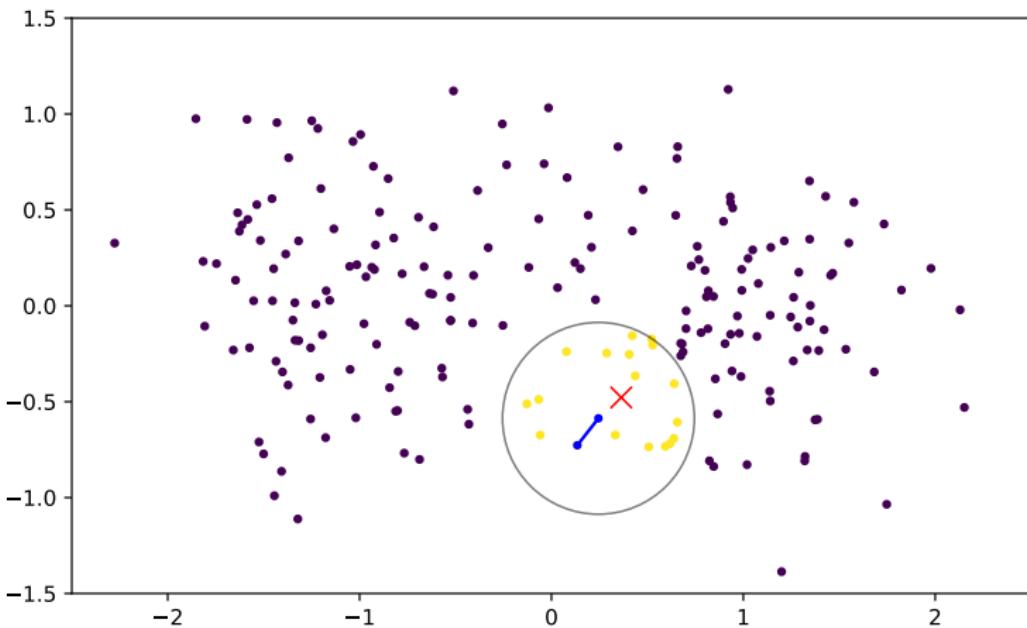
Convergence – Gaussian kernel



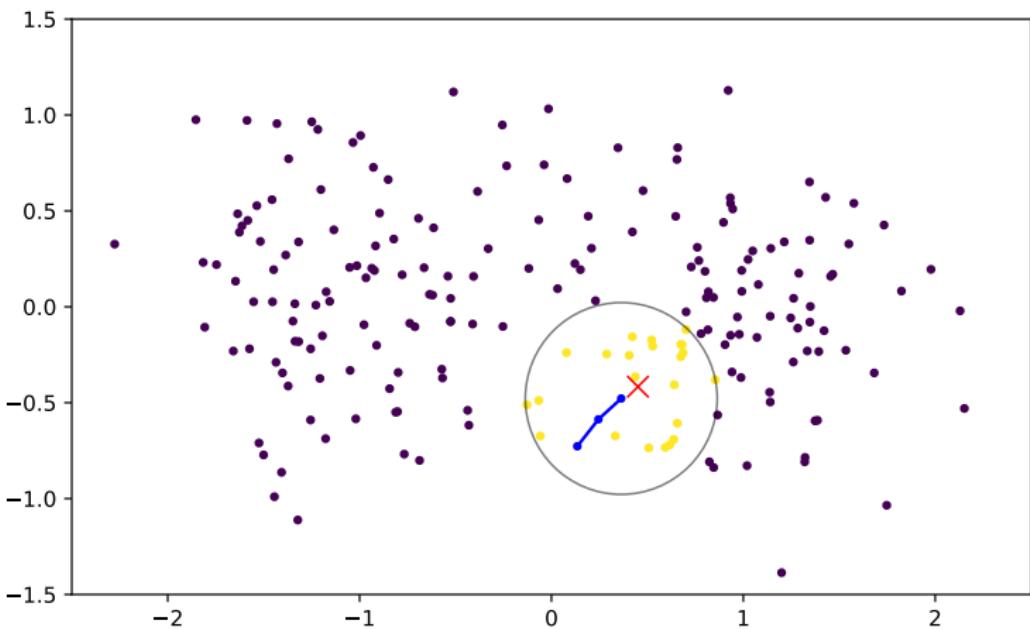
Convergence – Uniform kernel



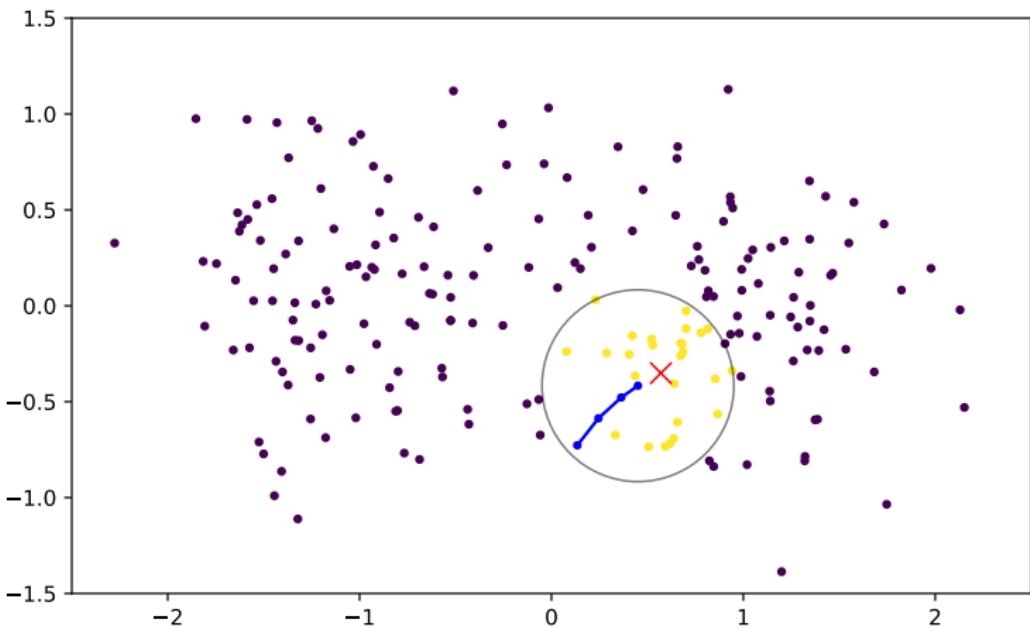
Convergence – Uniform kernel



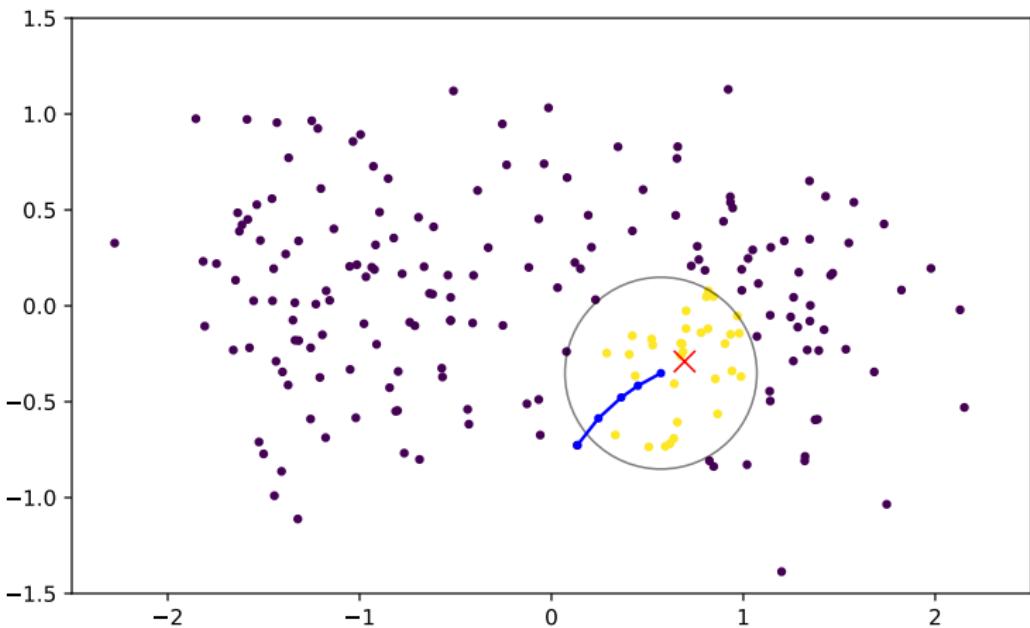
Convergence – Uniform kernel



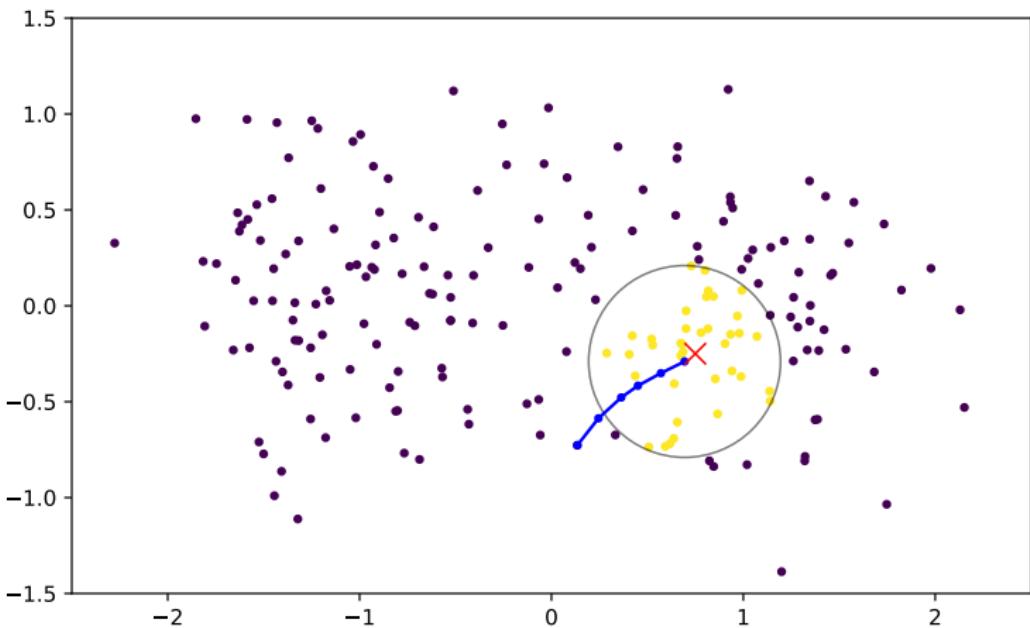
Convergence – Uniform kernel



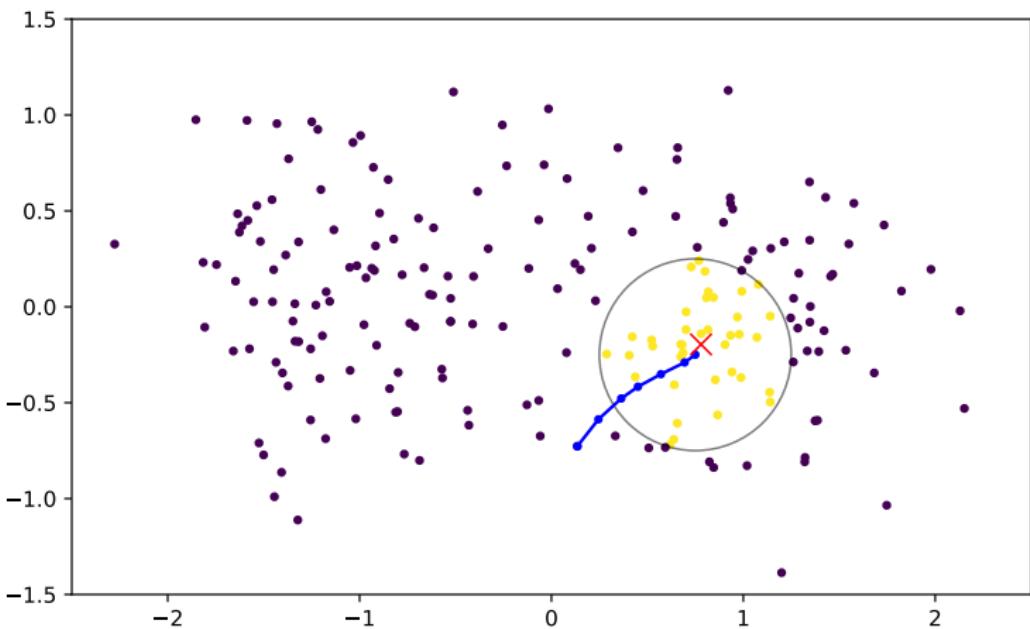
Convergence – Uniform kernel



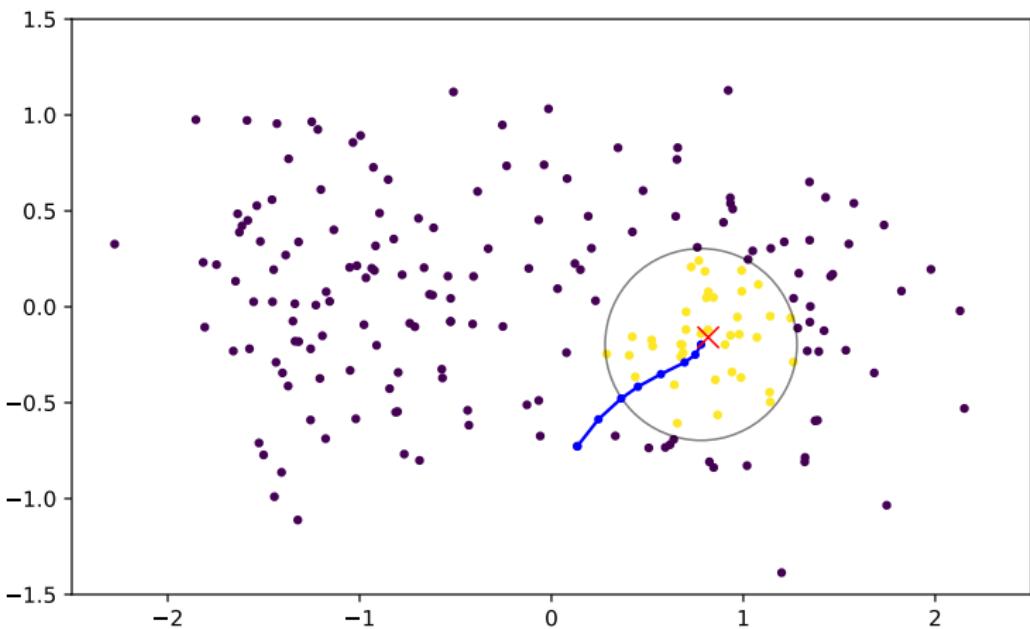
Convergence – Uniform kernel



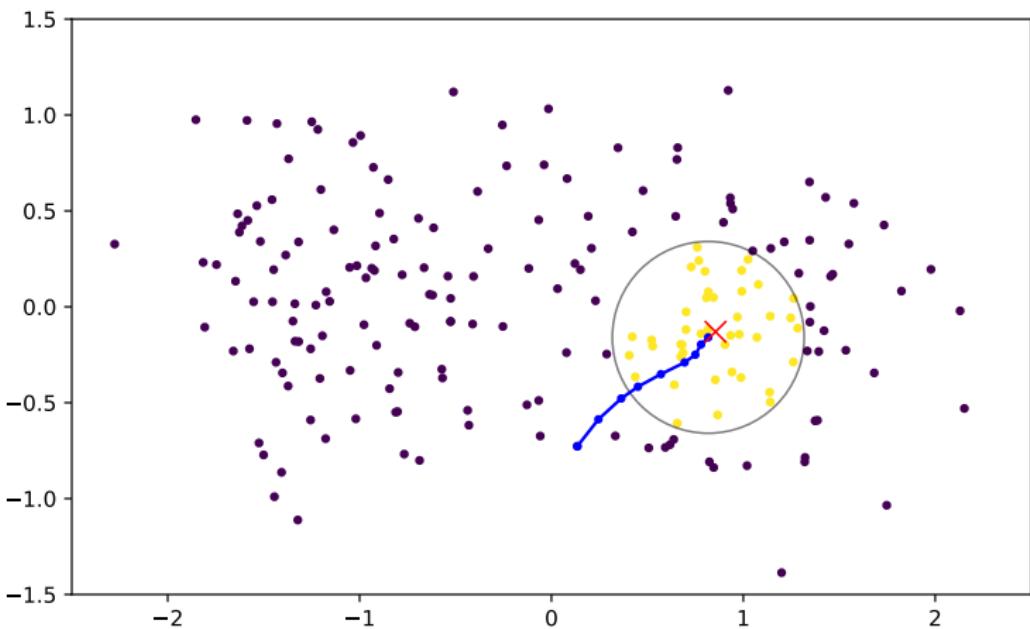
Convergence – Uniform kernel



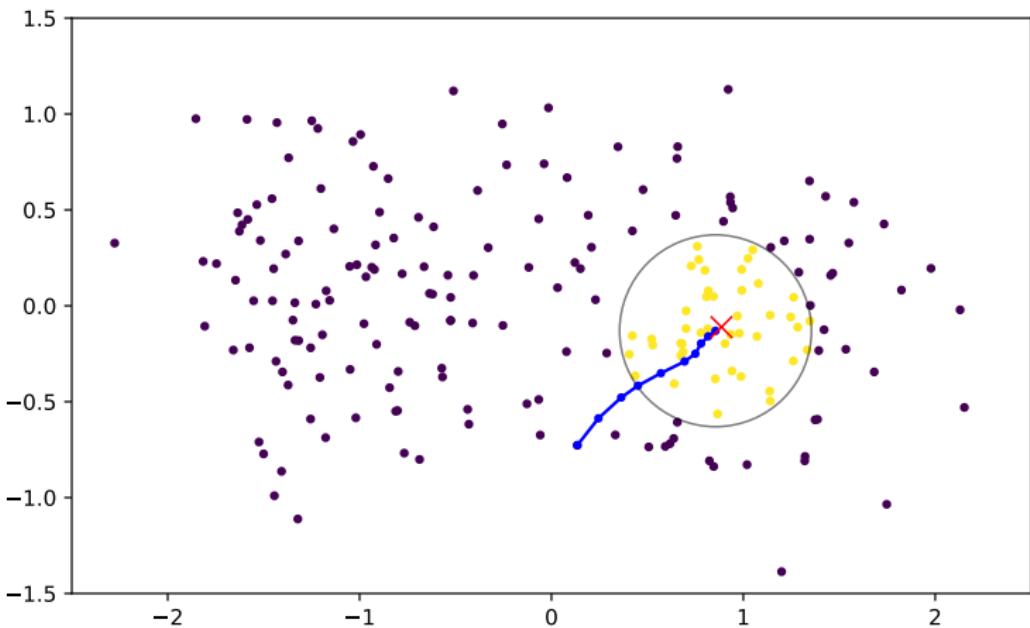
Convergence – Uniform kernel



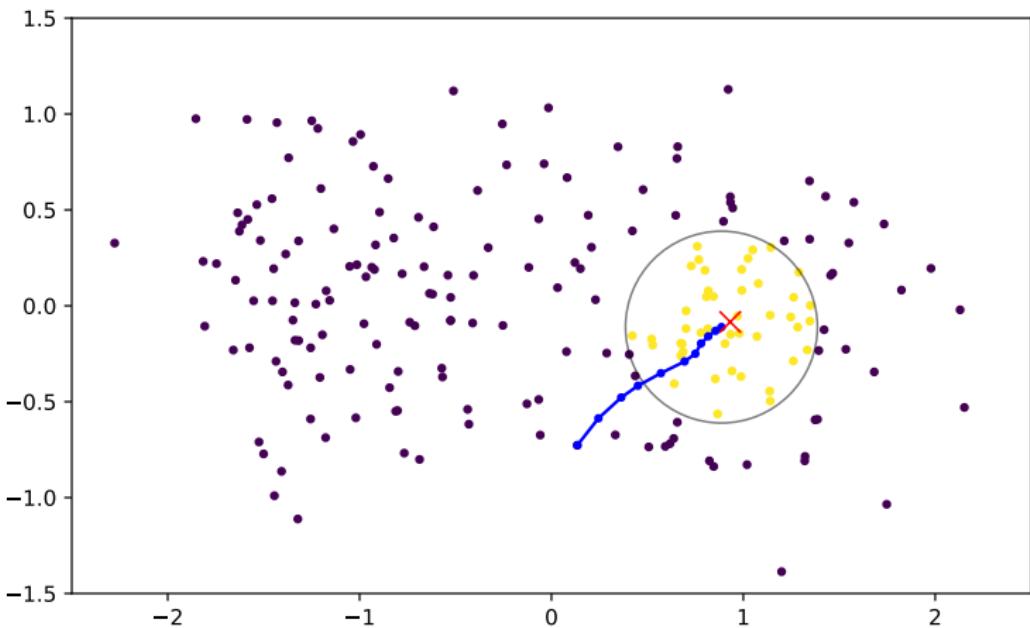
Convergence – Uniform kernel



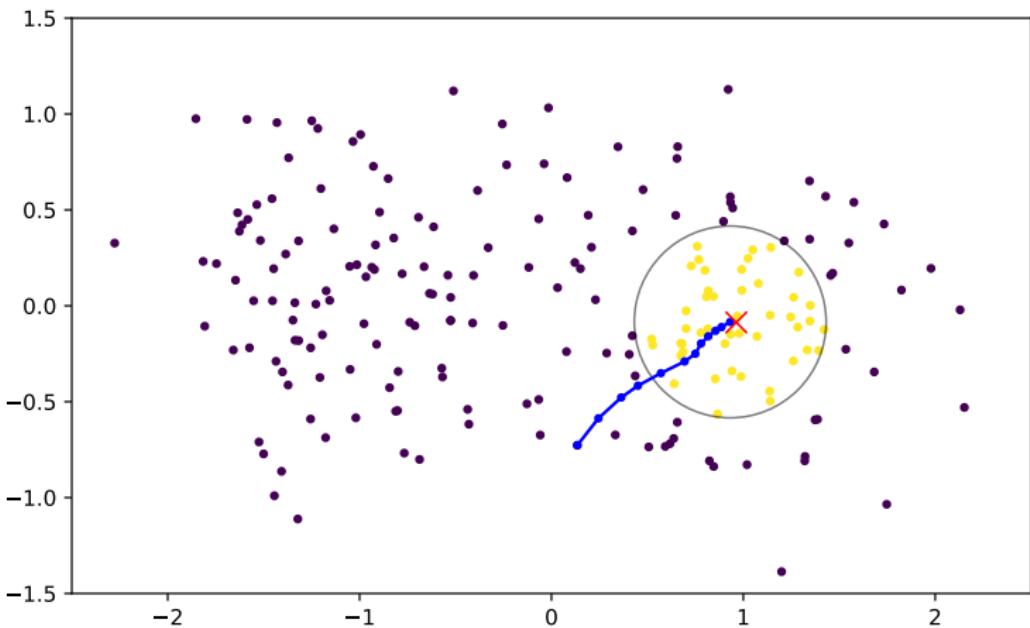
Convergence – Uniform kernel



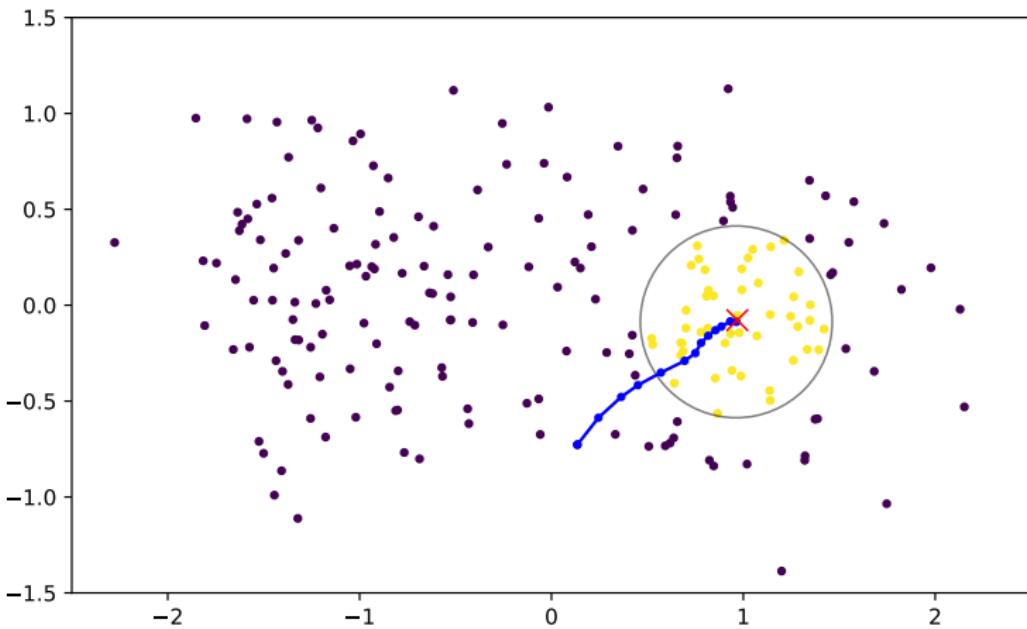
Convergence – Uniform kernel



Convergence – Uniform kernel



Convergence – Uniform kernel



Connecting the mean shift vector and kernel density estimation

- **Result:** the mean shift vector $\mathbf{m}(\mathbf{x})$ points into the gradient direction of a kernel density estimate.

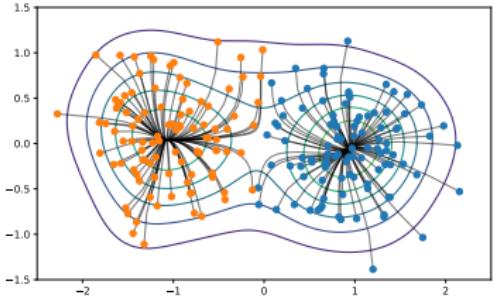
$$\mathbf{m}(\mathbf{x}) = \boldsymbol{\mu}^*(\mathbf{x}) - \mathbf{x} = \frac{h^2 c_{g,d}}{2 c_{k,d}} \frac{\nabla \hat{f}_{h,K}(\mathbf{x})}{\hat{f}_{h,G}(\mathbf{x})}$$

Connecting the mean shift vector and kernel density estimation

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

With $g(u) = -k'(u)$ and $G(\mathbf{x}) = c_{g,d}g(\|\mathbf{x}\|^2)$.

Connecting the mean shift vector and kernel density estimation



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

- For the gaussian kernel $k(u) \propto -k'(u) = g(u)$ and therefore $K = G$.

Connecting the mean shift vector and kernel density estimation

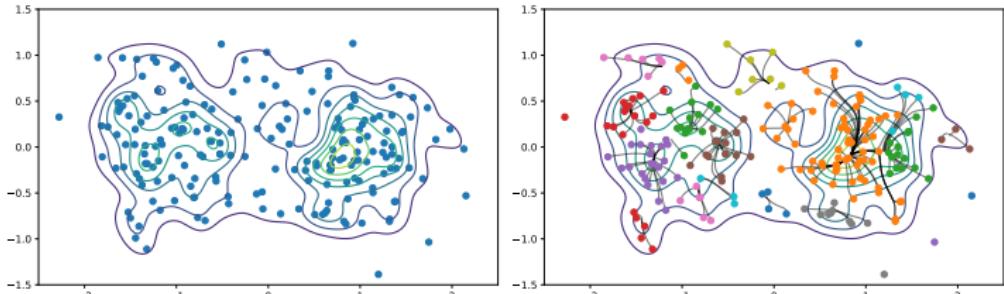
| Name | $k(u)$ | $-k'(u)$ |
|--------------|----------------------------------|---|
| Uniform | 1 | 0 |
| Epanechnikov | $1 - u$ | 1 |
| Biweight | $(1 - u)^2$ | $2(1 - u)$ |
| Triweight | $(1 - u)^3$ | $3(1 - u)^2$ |
| Gaussian | $\exp\left(-\frac{1}{2}u\right)$ | $\frac{1}{2}\exp\left(-\frac{1}{2}u\right)$ |

- A kernel K for which $-k'(u) = g(u)$, is called a **shadow** of G .
- Calculating the mean shift vector $\mathbf{m}(\mathbf{x})$ with G does calculate the gradient direction of the kernel density estimate with K as kernel.

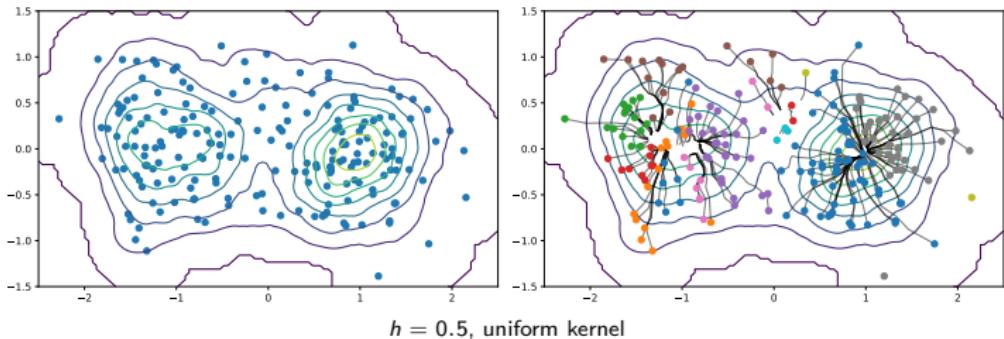
Bandwidth effects

- The choice of bandwidth influences the density estimation and therefore the clustering outcome
- Small bandwidth \Rightarrow many density peaks / clusters
- Large bandwidth \Rightarrow few peaks / clusters

Small bandwidth

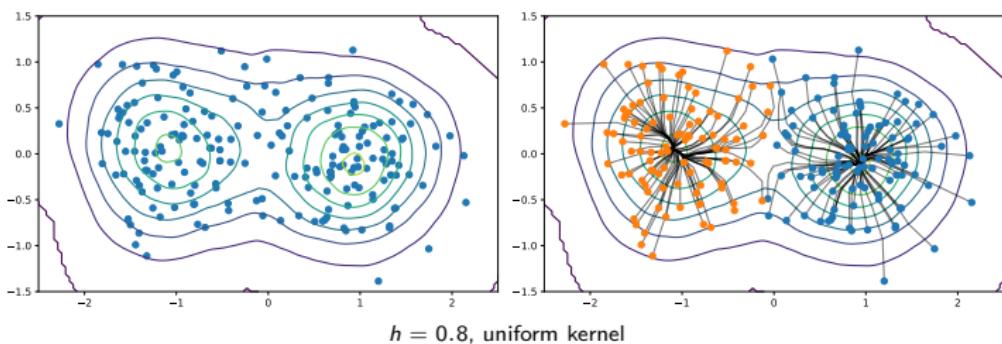
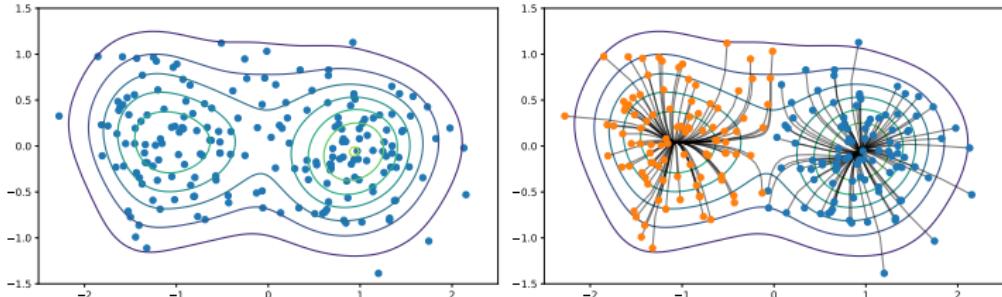


$h = 0.15$, gaussian kernel

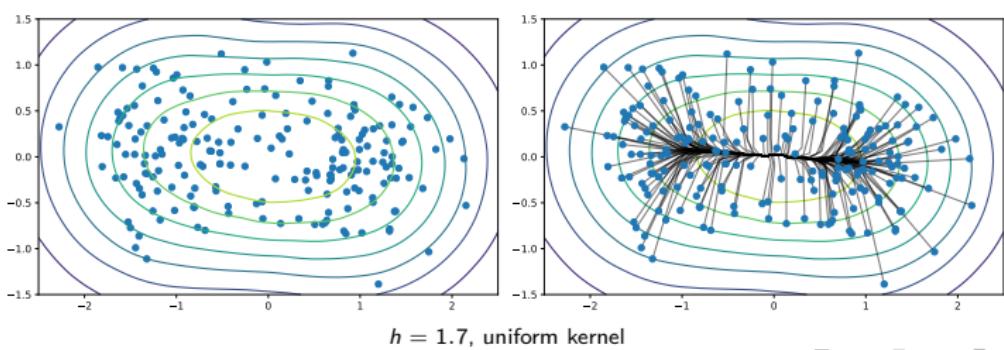
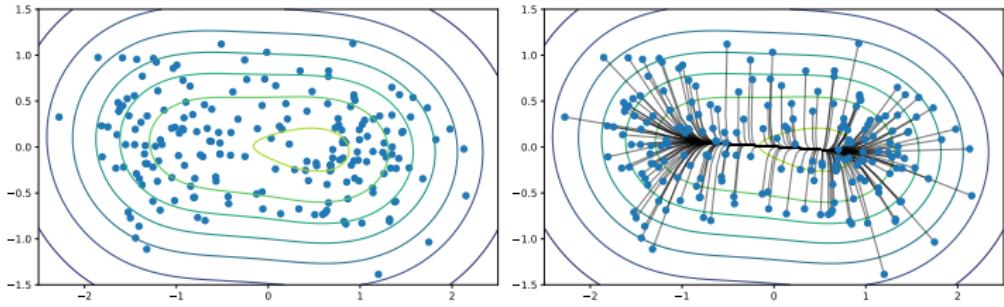


$h = 0.5$, uniform kernel

Medium bandwidth



Large bandwidth



Speedup methods

The mean shift algorithm has computational complexity $\mathcal{O}(Tn^2)$, with T being the number of iterations until convergence.

- **Parallelization** of iteration procedure. The trajectories are computational independent.
- **Sampling** of the dataset. Only perform it on a subset.
- **Bucketing** the dataset. Reduce number of possible positions.
- Use of **spatial data structures**. For kernels with compact support only the data points in the neighborhood are relevant.
- **Merging** of trajectories. When two trajectories get close to each other it is likely that they converge to the same point.

Advantages

- No prior assumption on cluster shapes. Complex and non-convex shapes are possible.
- Only has one tuning parameter, the bandwidth.
- No restriction on number of clusters.
- Outliers do not affect the clustering.

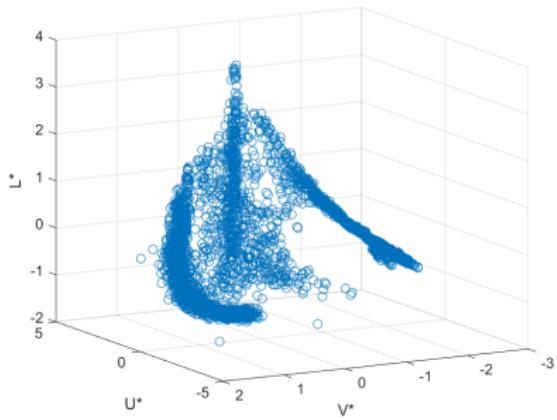
Disadvantages

- Density estimation fails for high dimensions ($\approx d > 5$).
- Bad computational complexity $\mathcal{O}(Tn^2)$.
- Finding a good bandwidth is hard.

Application

- Introduced by Fukunaga & Hostetler (1975), but popularized for computer vision tasks Comaniciu & Meer (2002) and Comaniciu et al. (2003).
- Generic clustering algorithm.
- Image segmentation, image filtering and object tracking.

Application – Image segmentation



- Image data can be represented as data points. The pixels can be clustered and the clustering results in a segmentation of the original image.

Preprocessing

- **Image rescaling** for quicker convergence.
- **Color space transformation**. The CIE 1976 L^*, u^*, v^* color space is built to approximate a perceptually uniform color space.
- **Feature standardization** with mean and standard deviation.
- **Adding spatial features**. Adding (x, y) pixel coordinates to pixel data.

Image segmentation – Gaussian kernel



(a) $h = 0.1$



(b) $h = 0.2$

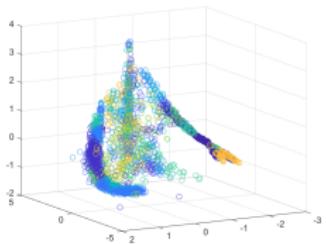


(c) $h = 0.3$

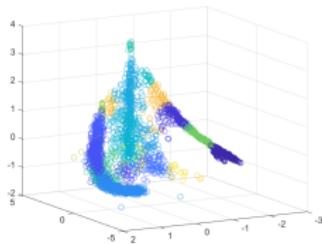


(d) $h = 0.4$

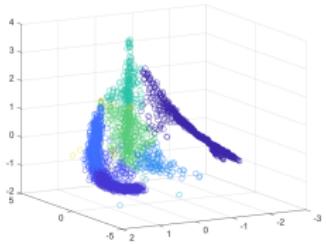
Image segmentation – Gaussian kernel



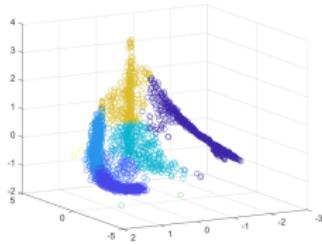
(a) $h = 0.1$



(b) $h = 0.2$



(c) $h = 0.3$



(d) $h = 0.4$

Image segmentation – Uniform kernel



(a) $h = 0.1$



(b) $h = 0.2$

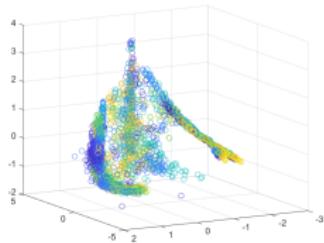


(c) $h = 0.3$

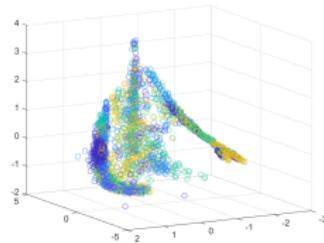


(d) $h = 0.4$

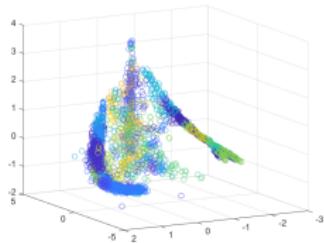
Image segmentation – Uniform kernel



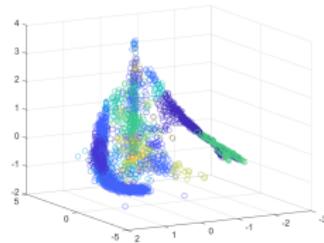
(a) $h = 0.1$



(b) $h = 0.2$

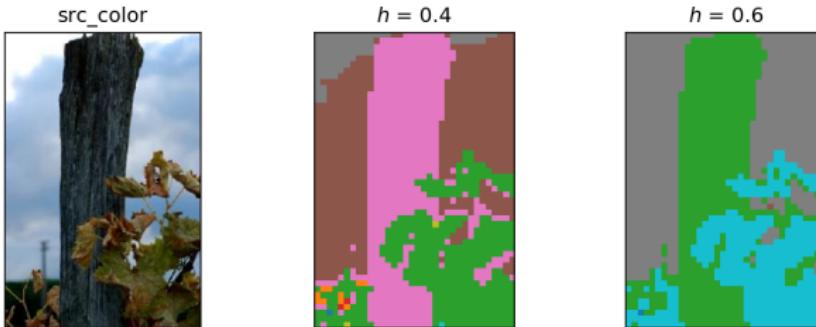
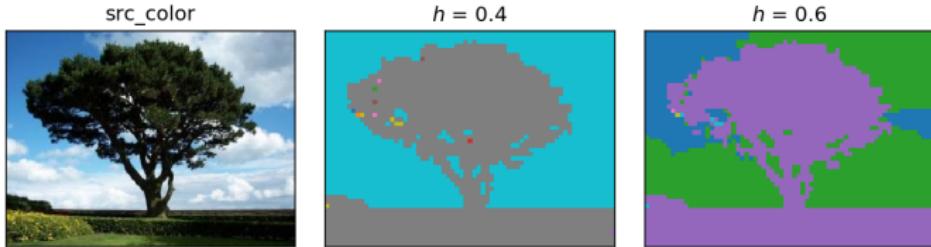


(c) $h = 0.3$

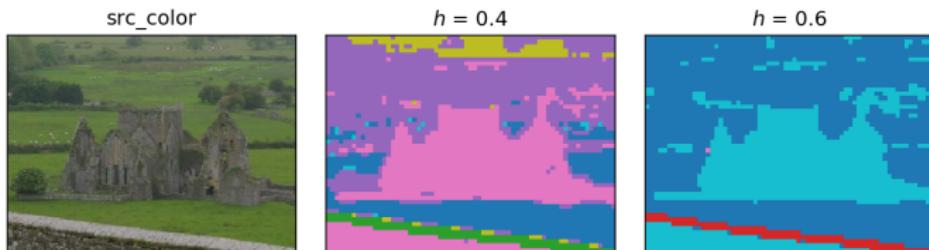
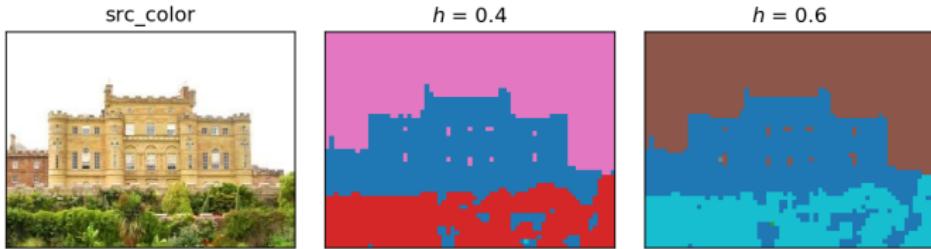


(d) $h = 0.4$

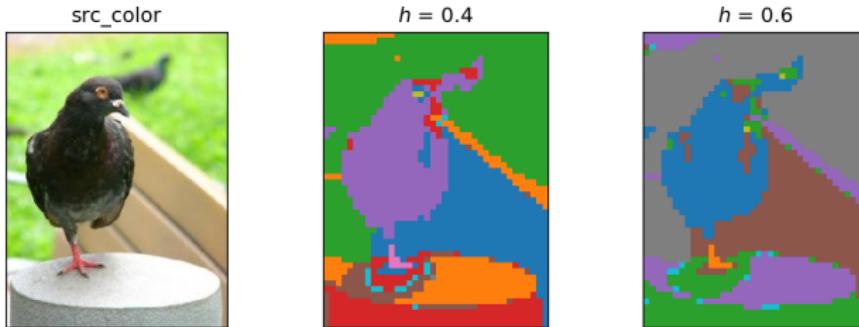
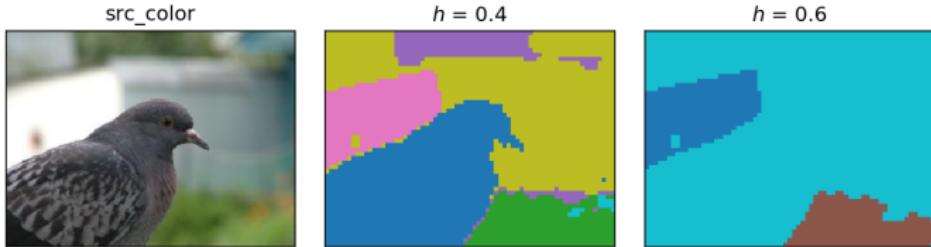
Results



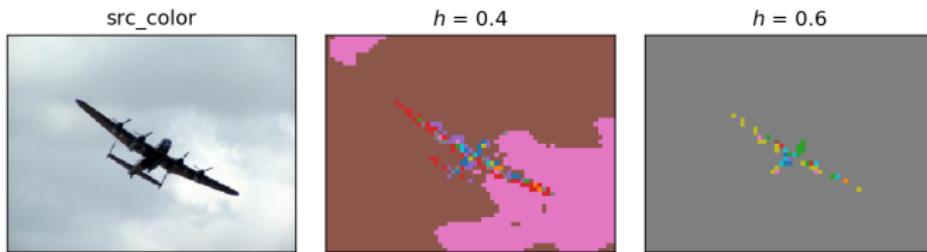
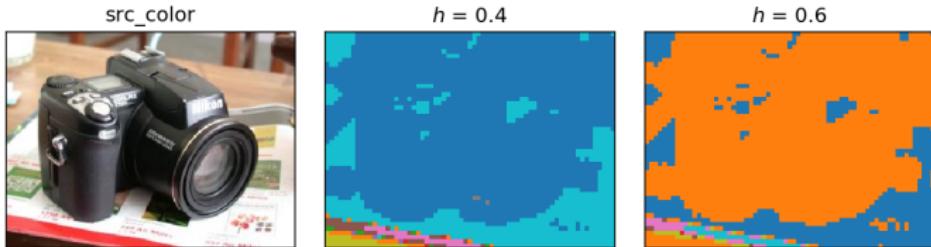
Results



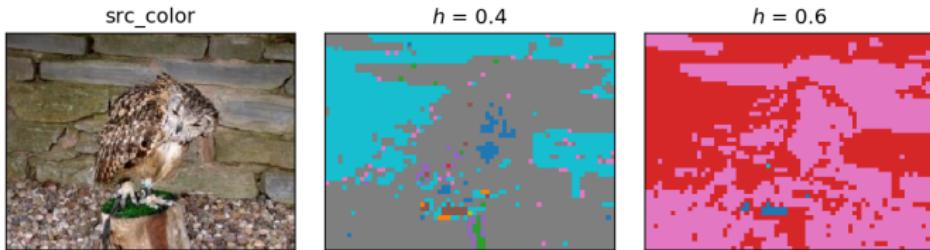
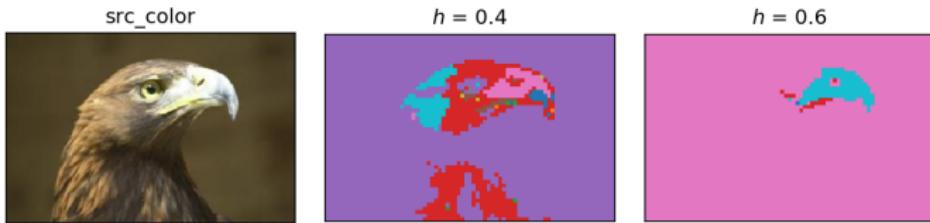
Results



Limitations



Limitations



References I

- 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.

Figures

- “KIT chemical faculty building” picture. Copyright by KIT.
- Segmentation pictures. “segmentation evaluation database” from Alpert et al. (2012).
- All other illustrations were done by the author.