

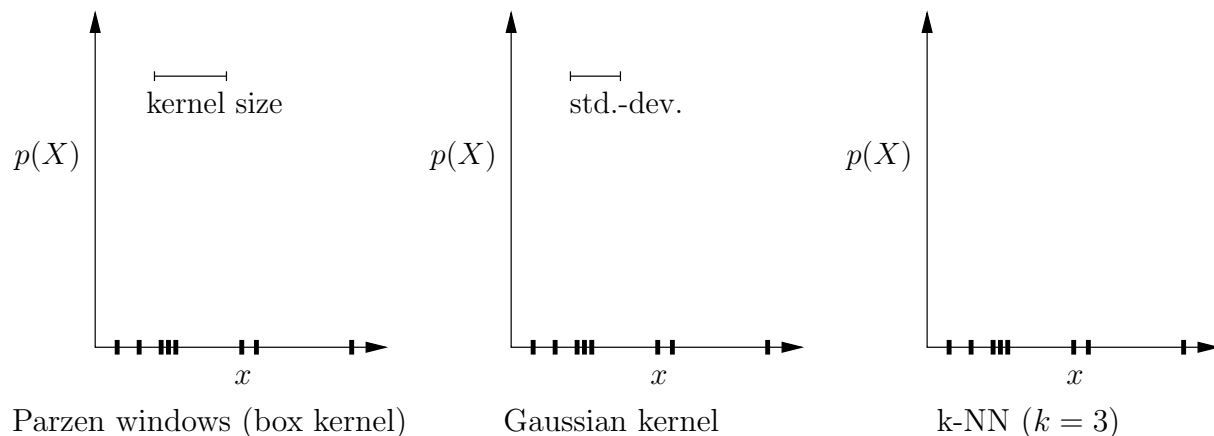


Please watch the video prior to the lecture, and think about the questions below. In the joint meeting, you will have 15 minutes time to discuss the questions with your group. Afterwards, we will jointly discuss your solution proposals.

You can print this sheet and use the space below for your notes.

Task 1: Sketch of Density Estimates

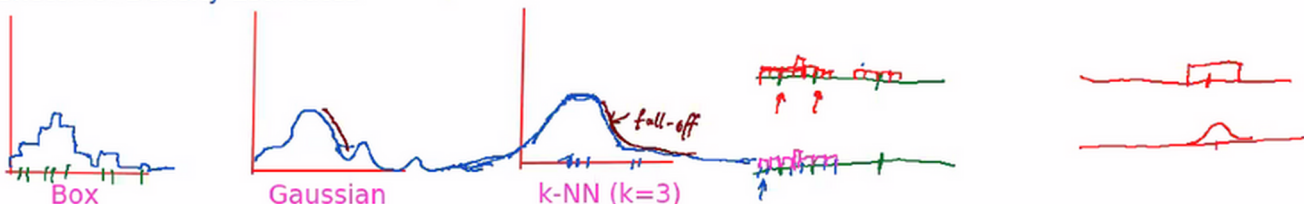
Below are three copies of a 1-D distribution of points along the x -axis. The goal is to perform three density estimations on the paper. Use the y -axis to indicate the probability mass. Of course, approximate ("sketch") solutions are just fine.



Task 2: Density Forensics

Assume that someone provides you with a ready-made density estimate as you produced it in task 1. How can you find out which of the three methods has been used?

Q1: Sketch of Density Estimates



Q2: "Density Forensics"

Proposal: draw samples, recompute density w/ box, gaussian, k-NN, calculate matching error select method with lowest error

Note that if you have an isolated sample, you directly see the box kernel or Gauss kernel

Task 3: Computational Cost and Memory Cost

Given a D -dimensional sample space, each dimension is scaled to a value range between 0 and 1. We have N samples. We would like use a kernel with compact support¹ to estimate the density.

Let us think about the computational cost and memory cost of these variants:

- (a) Discretize the space into a histogram, where each dimension is split into B bins. What is the computational cost of creating such a histogram? What is the memory cost?
- (b) How does the computational and memory cost grow when each dimension is discretized into $2B$ bins?
- (c) Assume that we implement a kernel density estimator with a box kernel of size K^D . If we naively implement this kernel density estimator, what is the computational cost and memory cost?
- (d) Open-ended question: can we do something smarter than the options stated above?
Note 1: This question has *many* possible answers, and is beyond the contents of the lecture. It aims to challenge your algorithmic thinking. Will you pick up the challenge?
Note 2: If a sketch helps to explain your idea, prepare one for the joint meeting.

Q3: Computational Cost & Memory Cost

a) Histogram D dim., B bins per dim:

Memory: $O(B^D)$ (-> query is $O(1)$)

or store only bin-ID per sample: $O(N \cdot D \cdot \log(B))$ (-> query is $O(N \cdot D)$
I would hypothesize that we can get significantly better than that with, e.g. locality sensitive hashing, or with a quad tree or so to expedite nearest neighbor search)

Computational Cost: $O(B^D \cdot K^D)$ if we do this really naively

¹Compact support means that only finitely many entries are non-zero, i.e., we can use the notion of a window size.