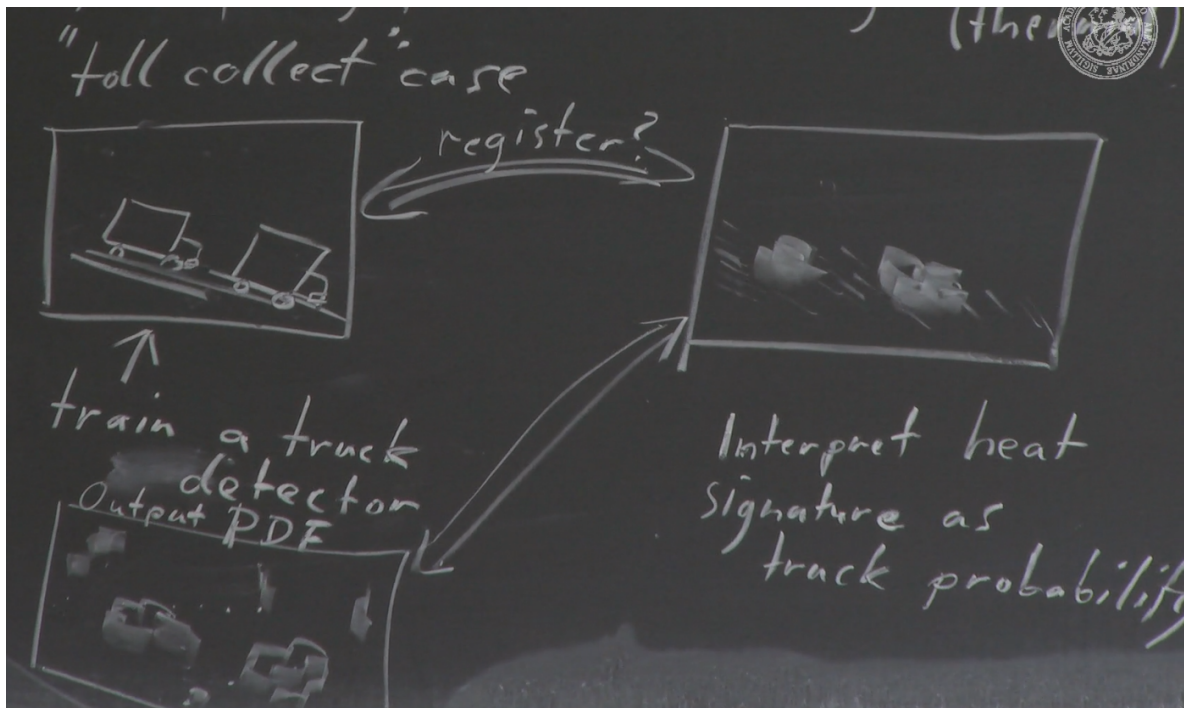


PA-2018 - 03: Non-parasitic density function

why might i want to estimate a probability density function (PDF) from some discrete observation?

- compute statistical measures, for example in Mutual information
 - image from modality 1
 - e.g. a photograph
 - image from modality 2
 - e.g. infrared image



- sample(draw) new observations with the same distribution as the actual(measurement)

Parzen Window Estimator

Idea: Given a set of discrete observations "smear" them out to obtain a PDF

Let $S = \{x_1, x_2, x_3, \dots, x_N\}$ denote the set of observations.

Let PR denote the probability that x is falling into region R .

$$PR = \int_R p(x) \partial x \quad (1)$$

if we assume $p(x)$ is approximately constant in R ,

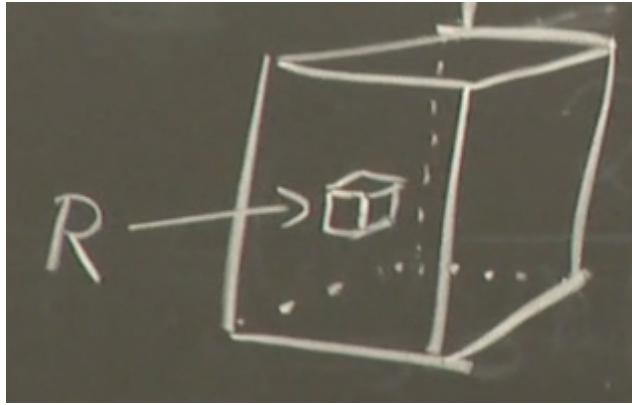
then $PR \approx p(x) * \int_R \partial x$

Here $\int_R \partial x = V_R$ is the "Volume" of R

$$\Rightarrow PR = p(x) * V_R = \frac{k_R}{N} * V_R$$

N is the number of features that fall in R over # of all features.

That called relative frequency feature in R .



For example, let R be a d -dimensional hypercube and let h denote the side-length of the hypercube $V_R = h^d$

The kernel window function is

$$K(x_c, x) = \begin{cases} 1 & \text{if } \frac{x_{i,k} - x_k}{h} \leq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

$$\Rightarrow \text{Rewrite } P(x) = \frac{1}{N} * \sum_{i=1}^N K(x_i, x) \quad (2)$$

Alternatively, $K(x_i, x)$ can be any other kernel for example Gaussian.

$$K_{\Sigma}(x_i, x) = \frac{1}{\sqrt{\det(\Sigma) * 2\pi}} * e^{-\frac{1}{2}(x_i - x)^T \Sigma^{-1} (x_i - x)}$$

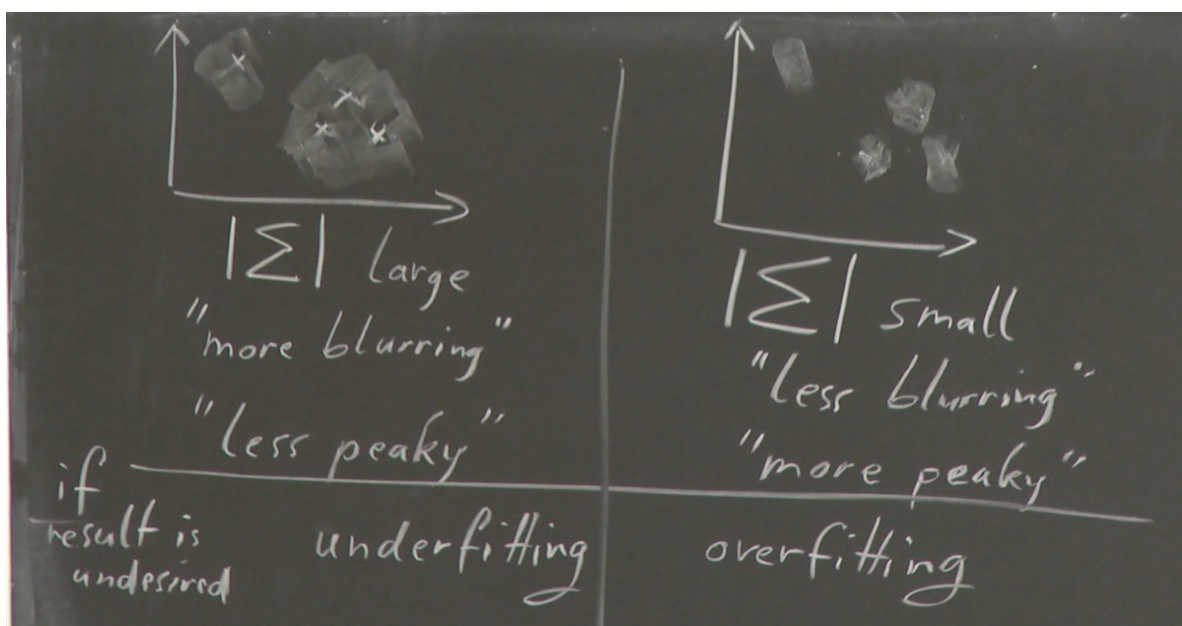
where Σ is covariance

Question tree

Question tree: how do we obtain Σ ?

equivalent: how do we obtain h for the hypercube?

Qualitatively, how does the result change if $|\Sigma|$ become larger or smaller?



Estimation of the covariance Σ or the window width h can be done via ML estimation in the case of limited training data additionally with cross-validation

Cross-validation

Let $P_{\lambda, N-1}^i$ be the PDF defined by $S = \frac{S}{x_i}$
 $\lambda = \sum, h$

Then we consider the objective function

$$\hat{\lambda} = \underset{\lambda}{\operatorname{argmax}} \prod_{i=1}^N P_{\lambda, N-1}^i(x_i) \quad (3)$$

$P_{\lambda, N-1}^i$ is "Trained model for all samples except of x_i "
 (x_i) is Test sample

$$= \underset{\lambda}{\operatorname{argmax}} \sum_{i=1}^N \log P_{\lambda, N-1}^i(x_i) \text{ log-likelihood}$$

For a differentiable kernel, we can now compute the gradient and look for an optimum, if the kernel is non-differentiable, we have to "brute force" the solution (Note: Gaussian=differentiable, Hypercube=non-diff.)