

Nonparametric Models (KNN and KDE)

COMP3314 — Week 7

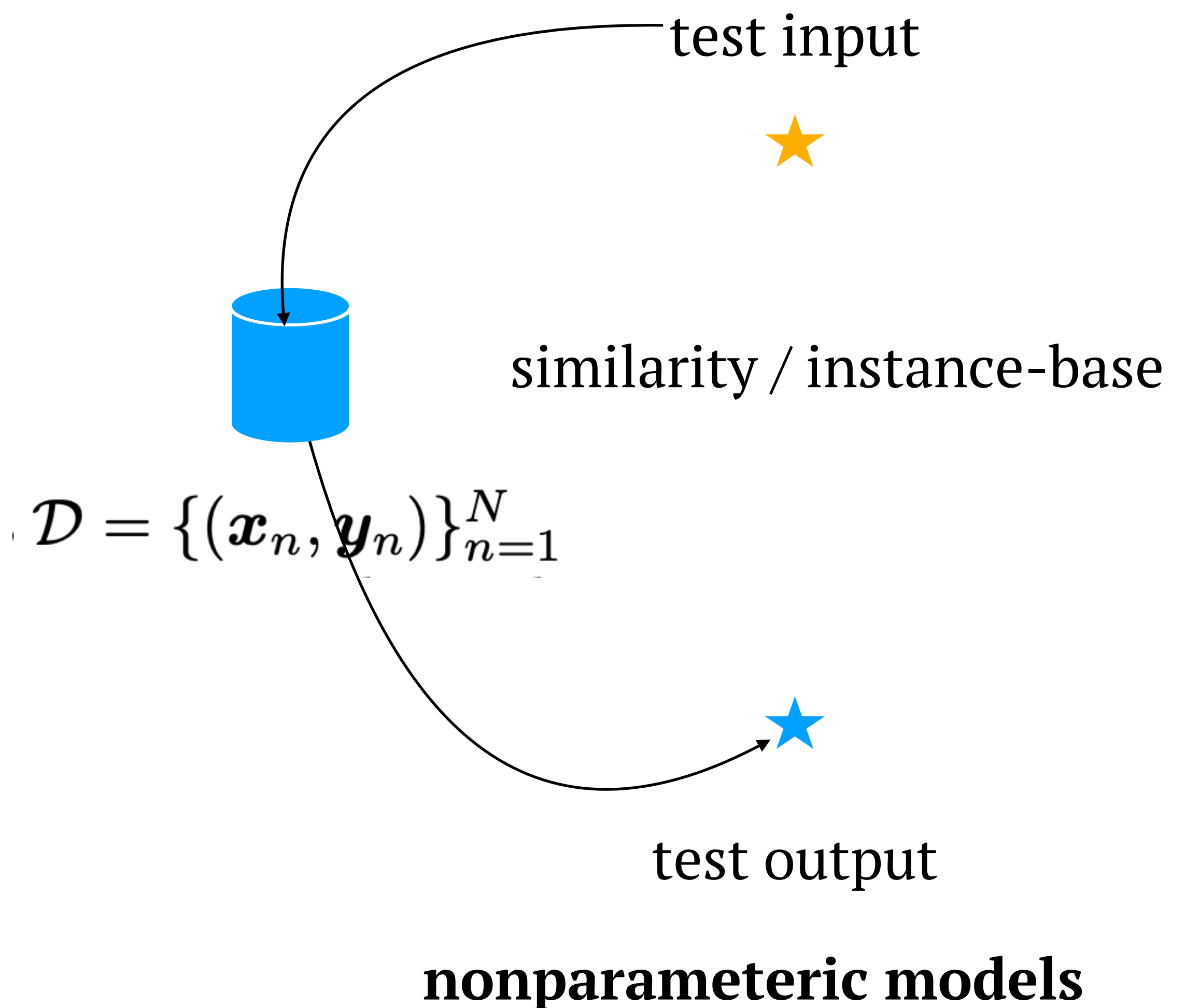
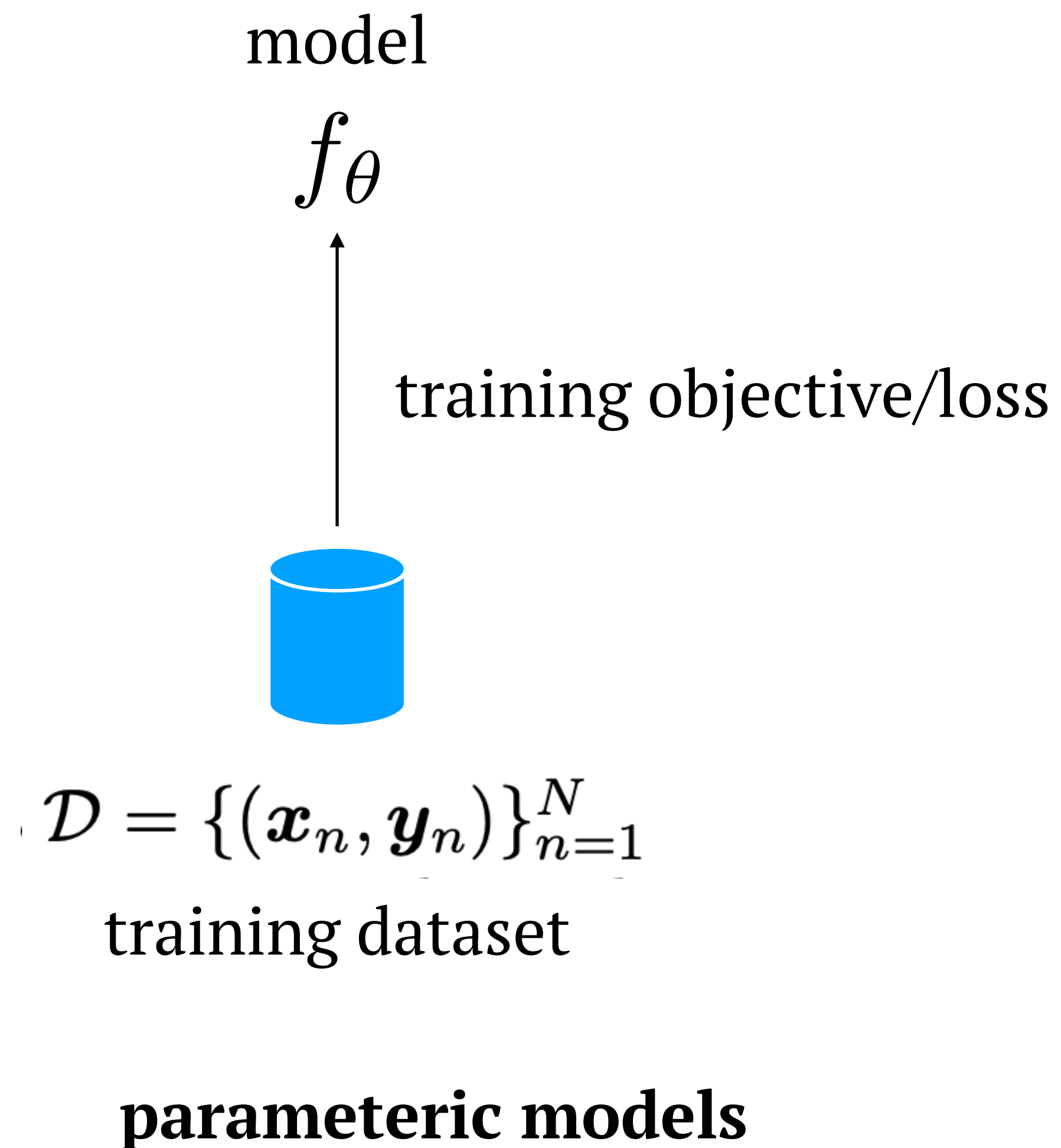
Lingpeng Kong

Department of Computer Science, The University of Hong Kong

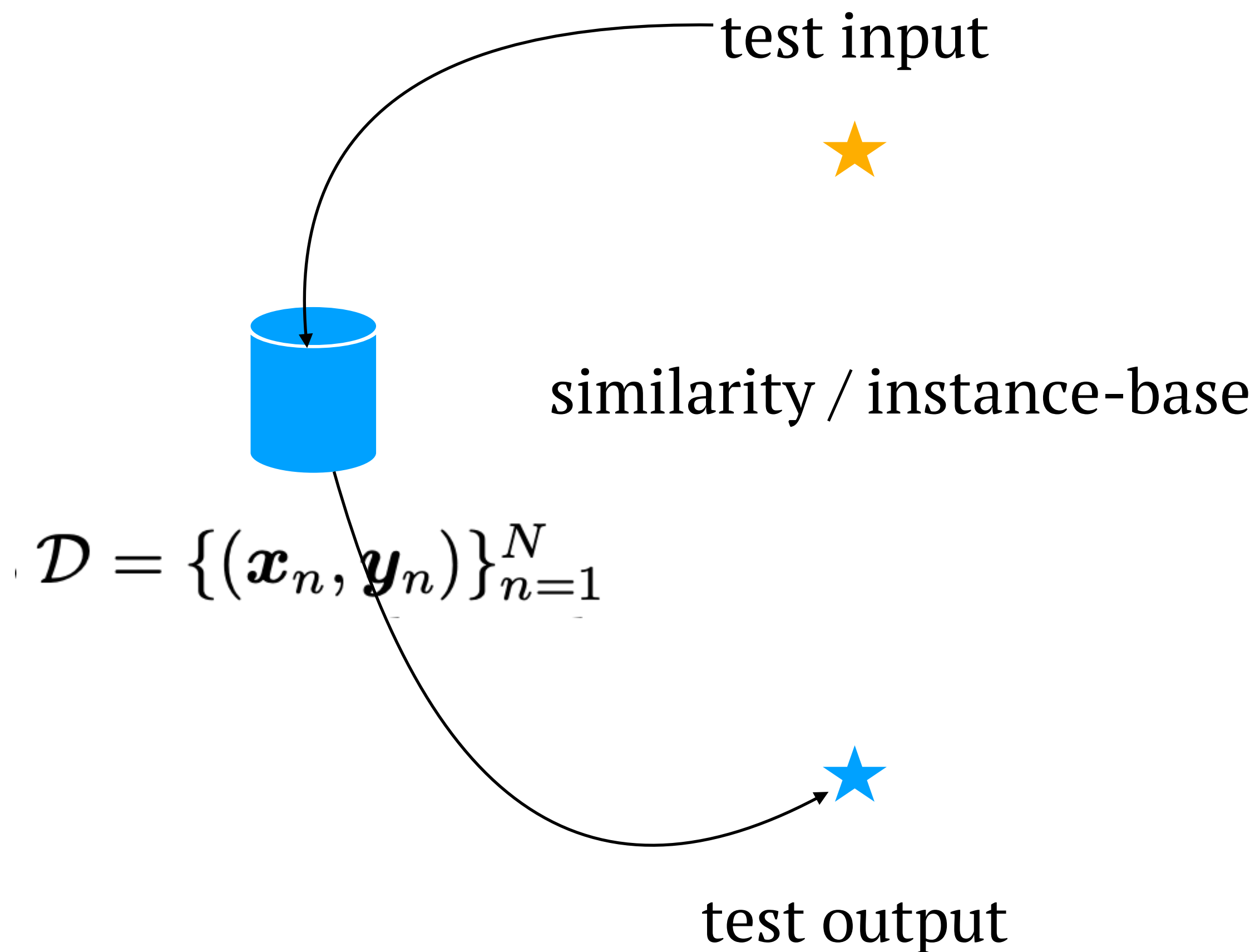
Based on: Probabilistic Machine Learning by Kevin Murphy

Slides from: Saw Shier Nee with special thanks!

Parametric Models v.s. Nonparametric models

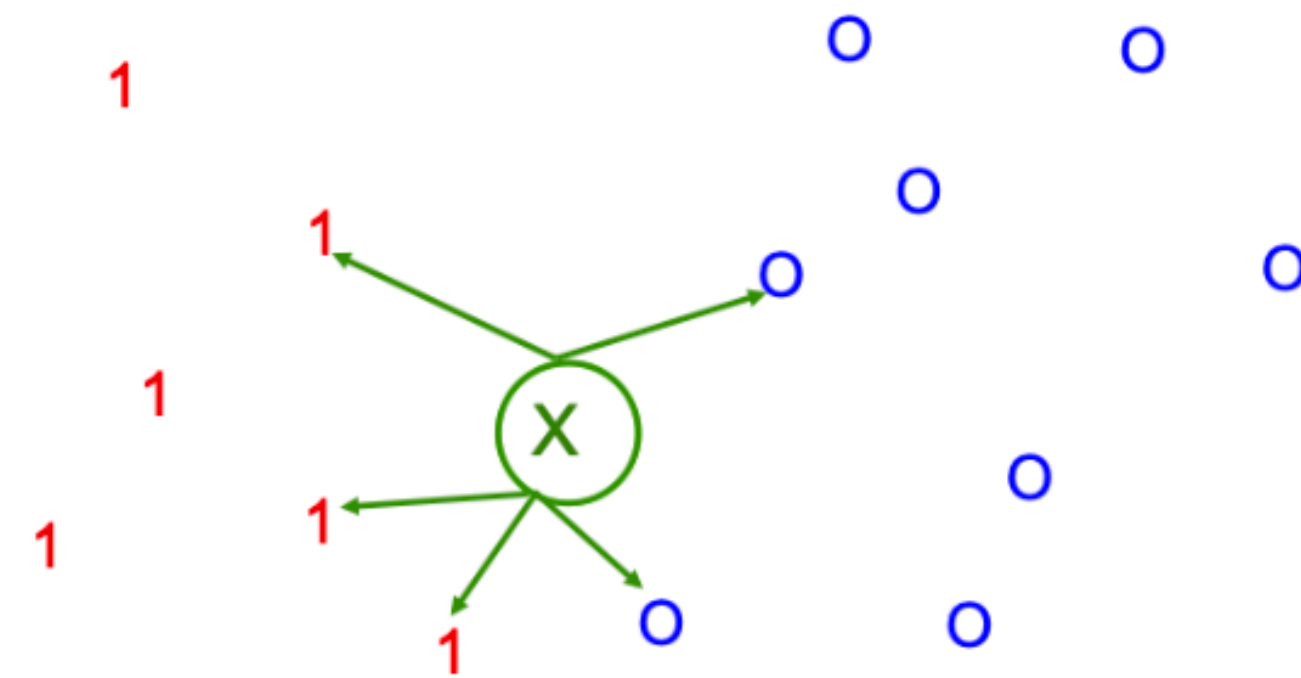


K Nearest Neighbor (KNN) Classification



$$p(y = c | \mathbf{x}, \mathcal{D}) = \frac{1}{K} \sum_{n \in N_K(\mathbf{x}, \mathcal{D})} \mathbb{I}(y_n = c)$$

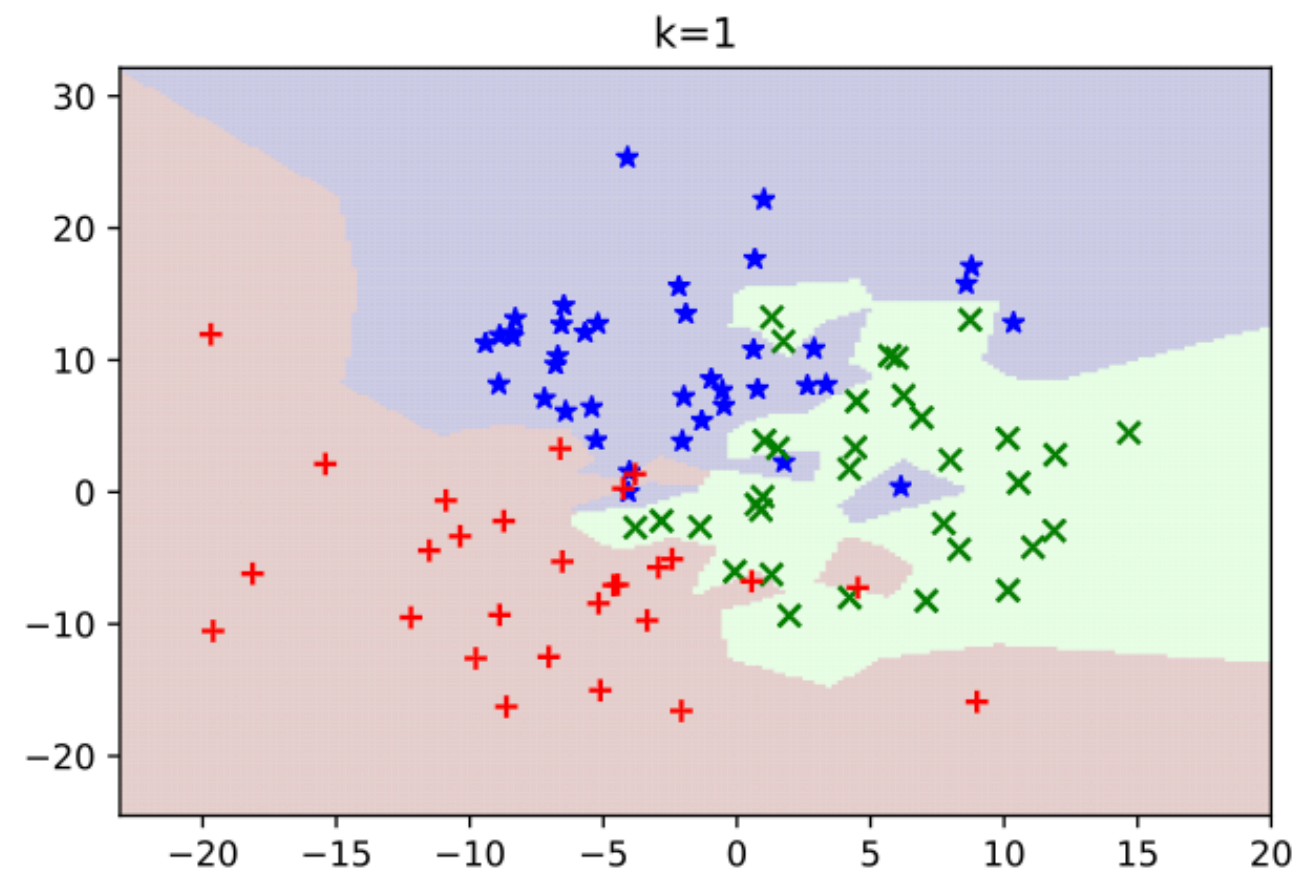
K $d(\mathbf{x}, \mathbf{x}')$



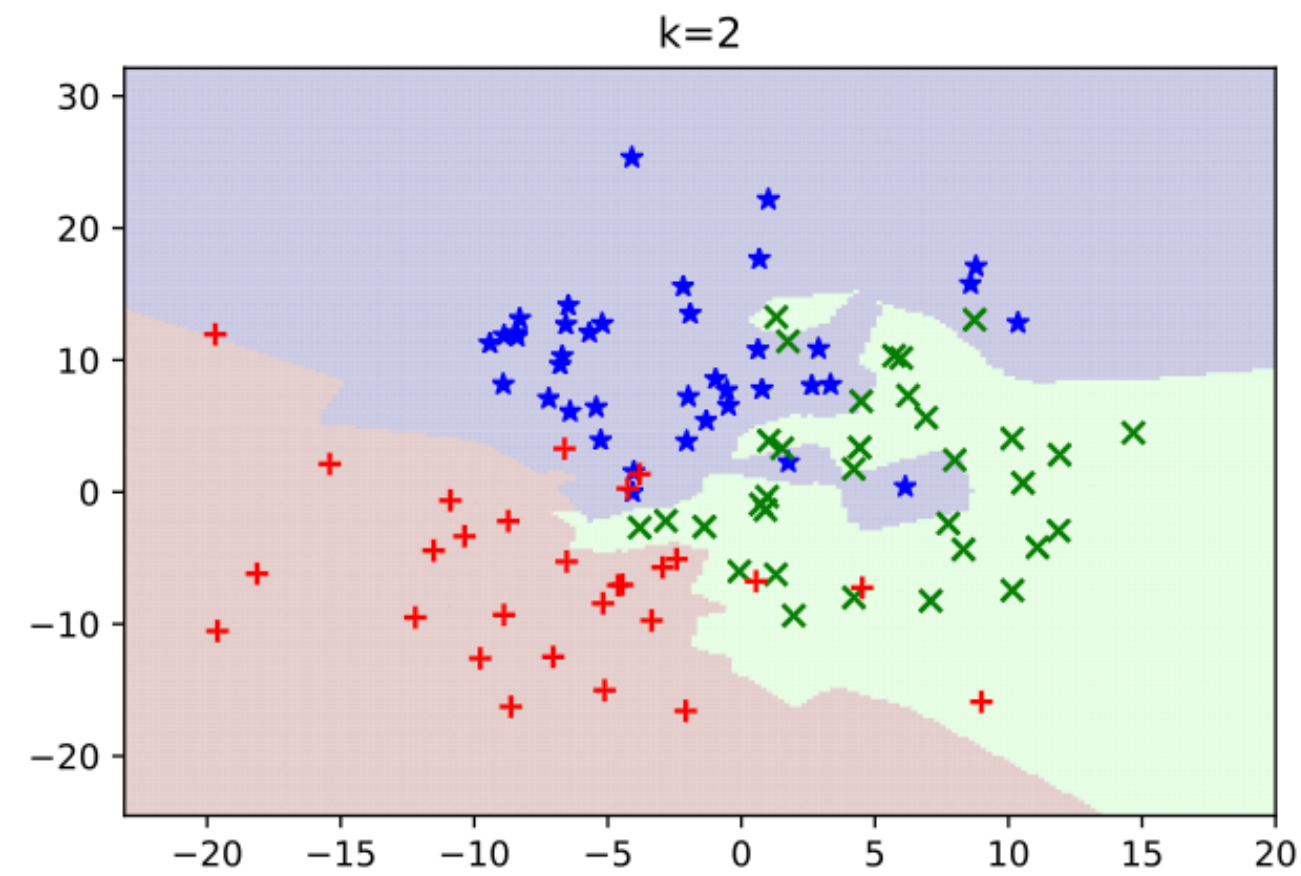
$$d_M(\mathbf{x}, \boldsymbol{\mu}) = \sqrt{(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{M} (\mathbf{x} - \boldsymbol{\mu})}$$

Mahalanobis distance

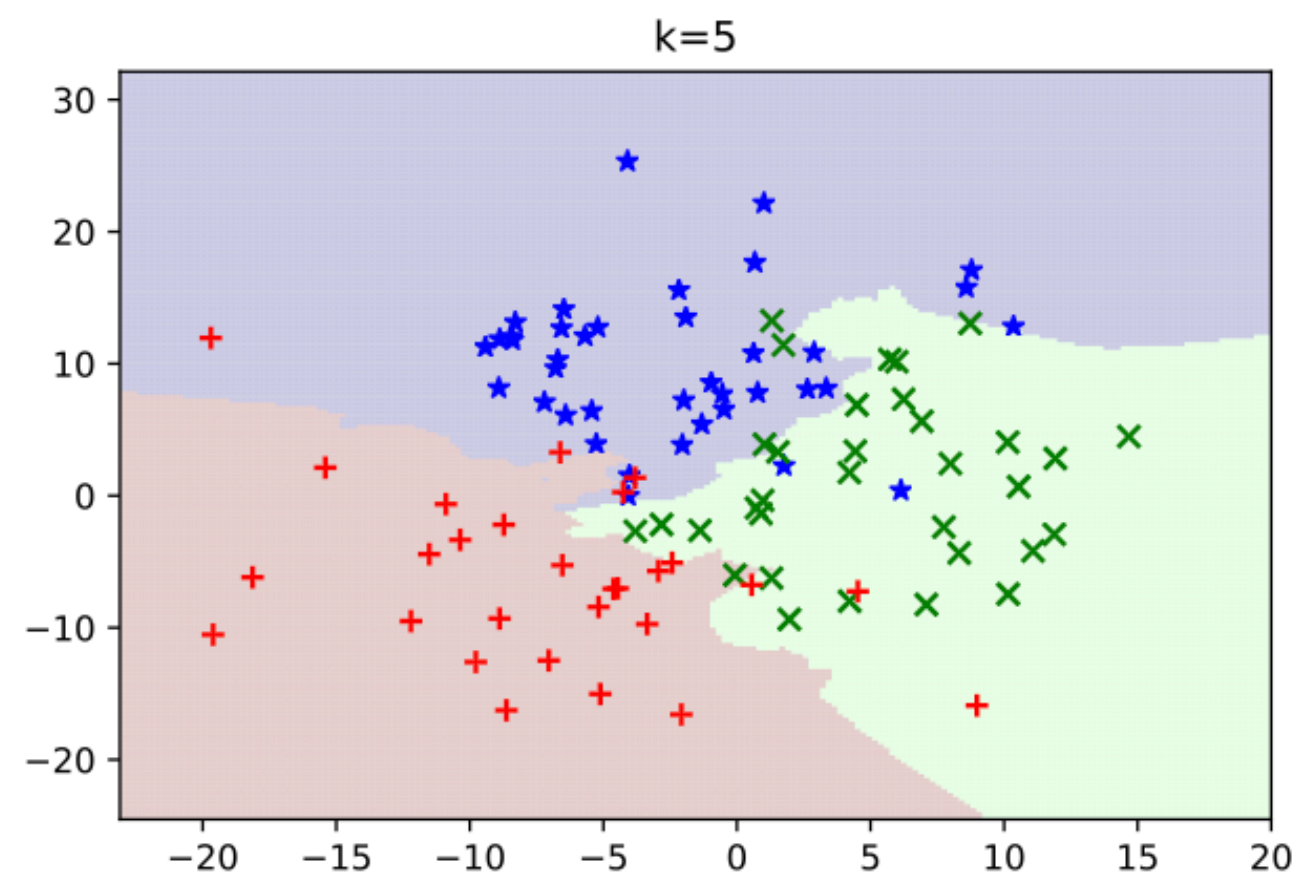
Overfitting in KNN



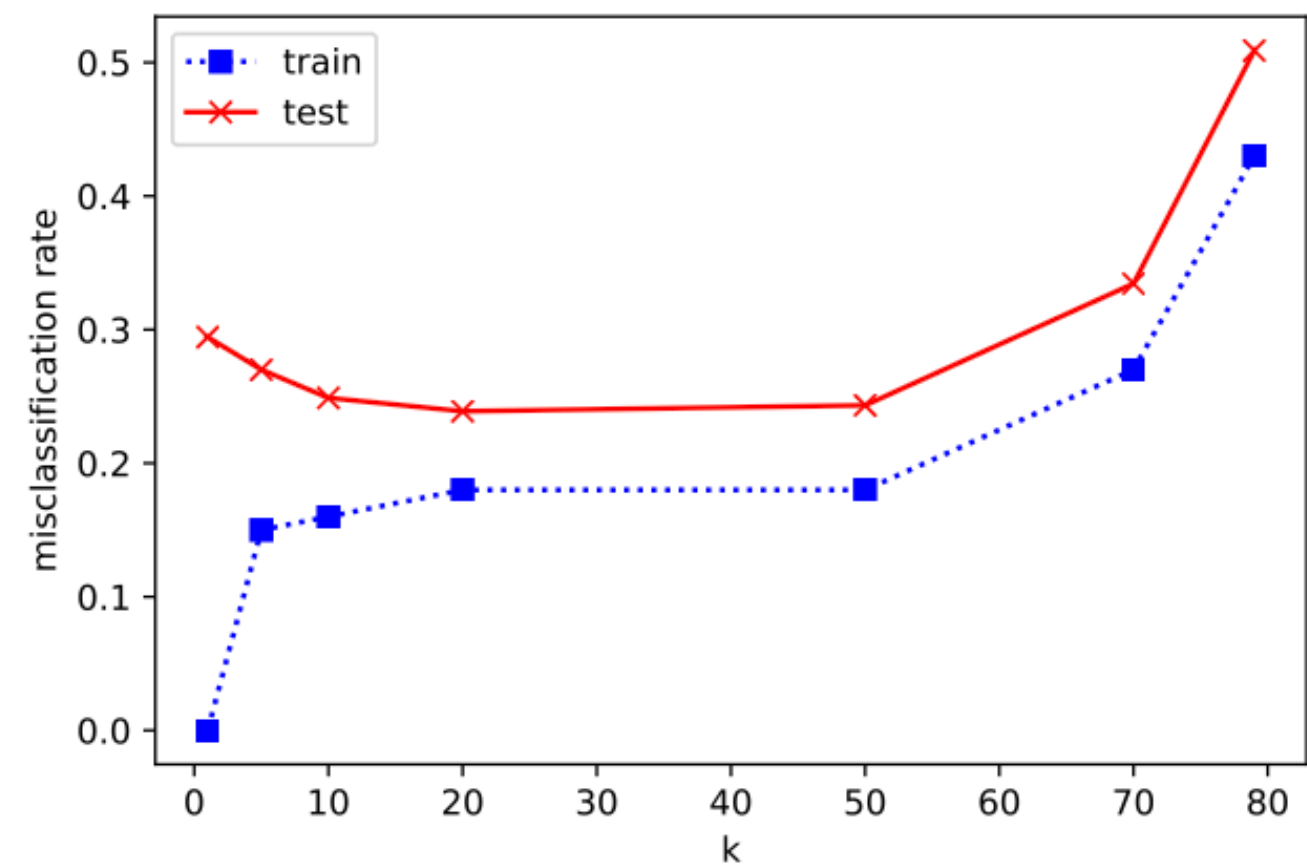
(a)



(b)



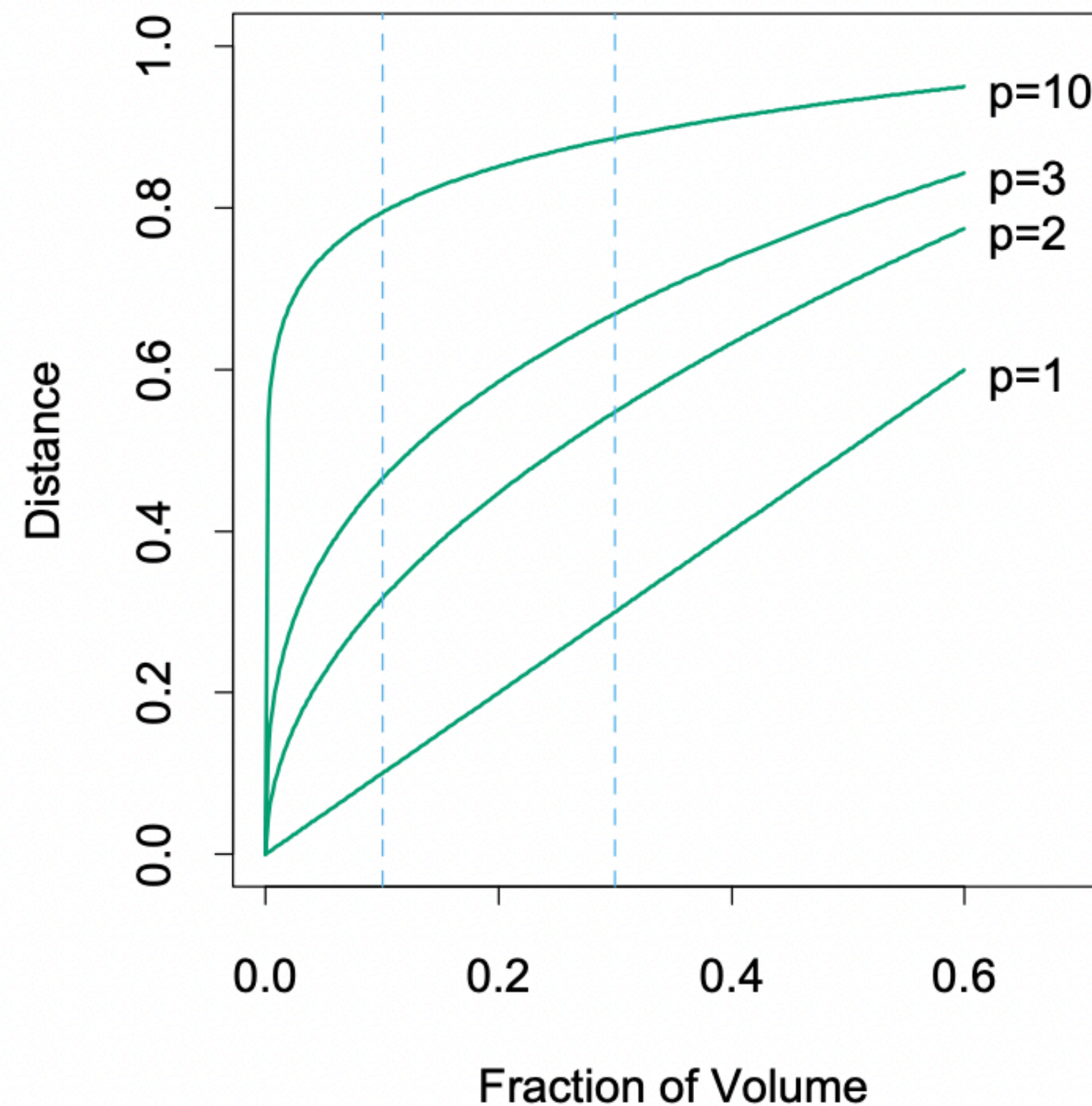
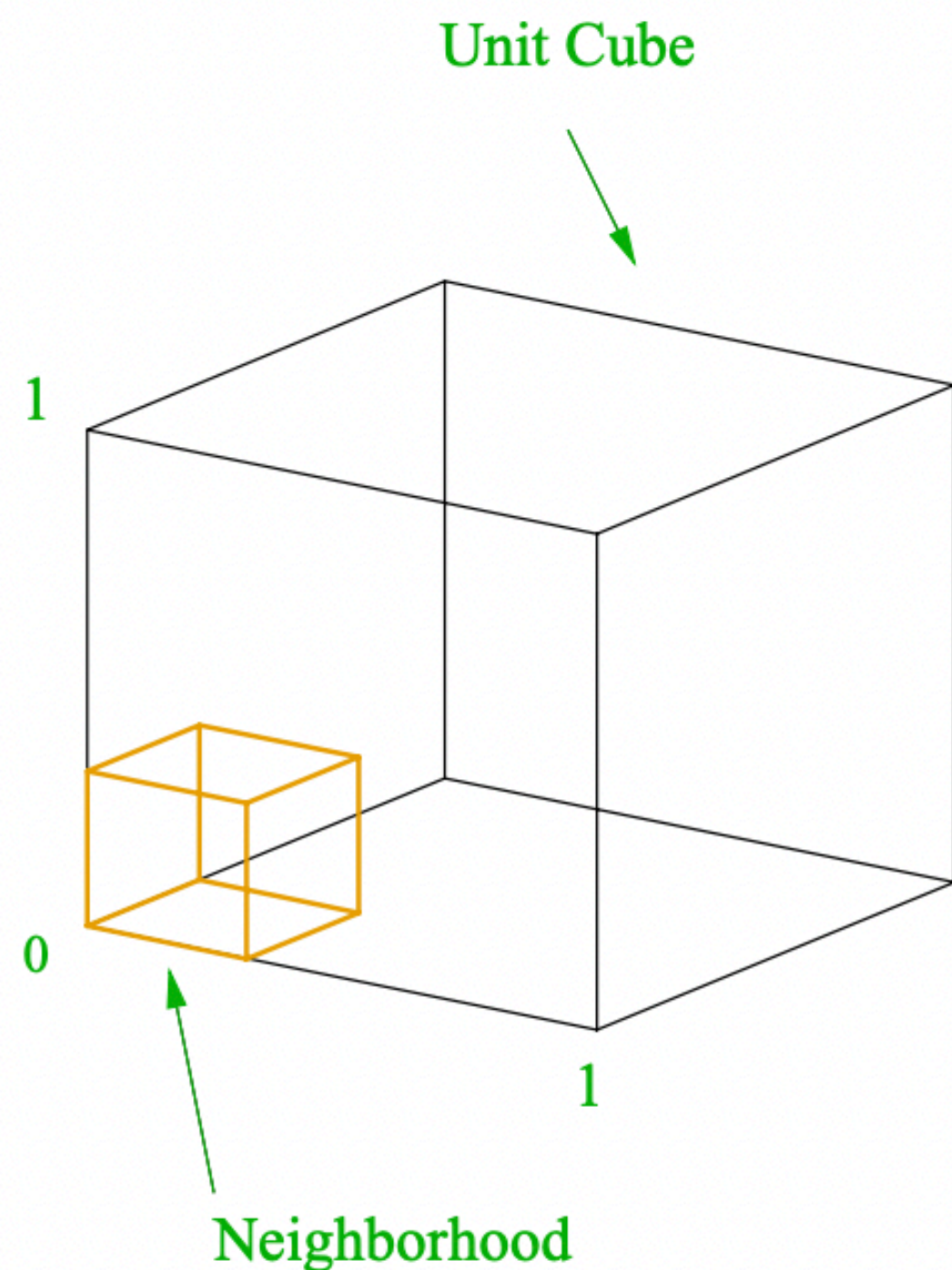
(c)



(d)

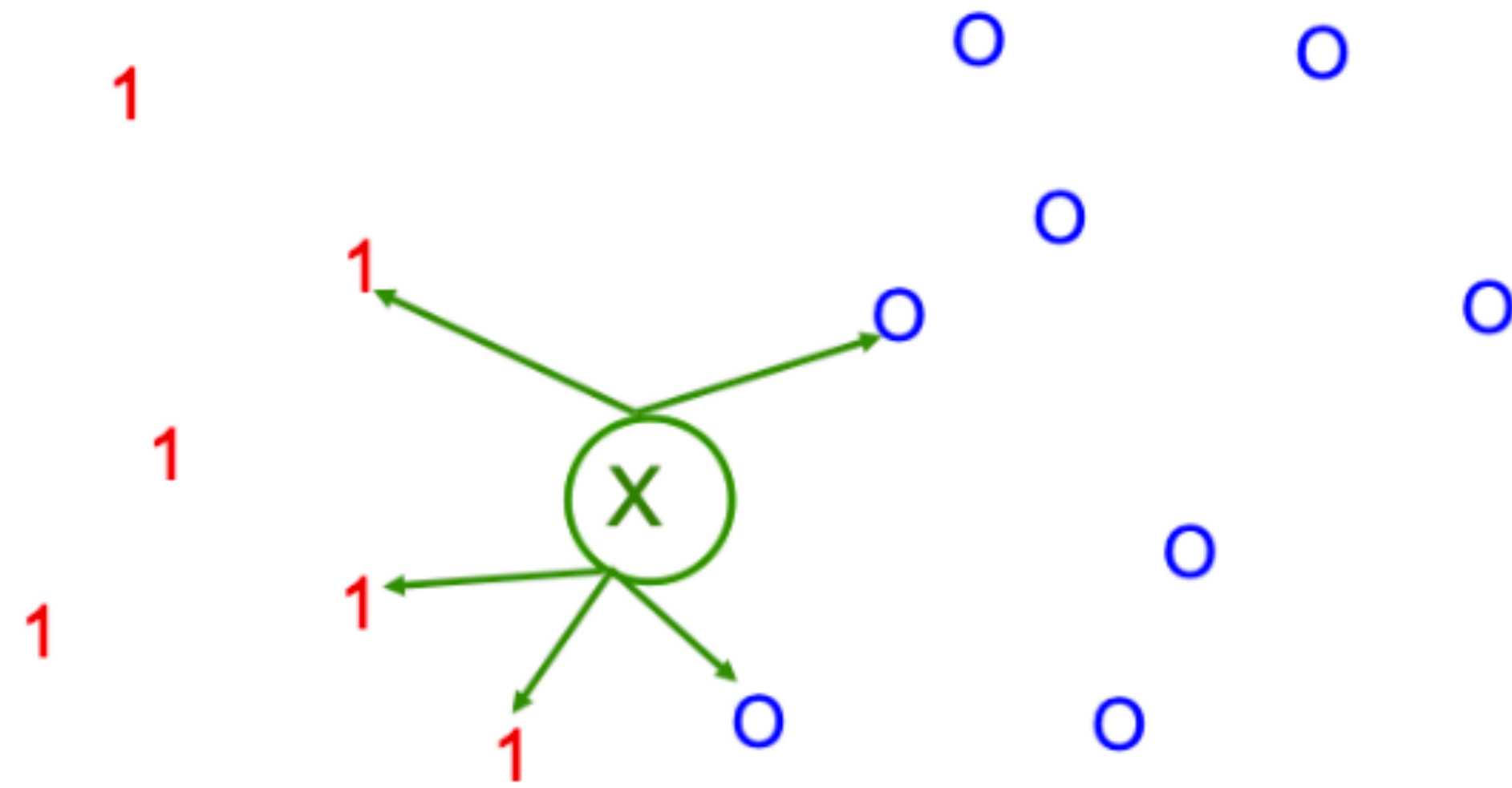
Which one is more likely to overfit, smaller K or larger K? Why?

The Curse of Dimensionality



The curse of dimensionality is well illustrated by a subcubical neighborhood for uniform data in a unit cube. The figure on the right shows the side-length of the subcube needed to capture a fraction r of the volume of the data, for different dimensions p . In ten dimensions we need to cover 80% of the range of each coordinate to capture 10% of the data.

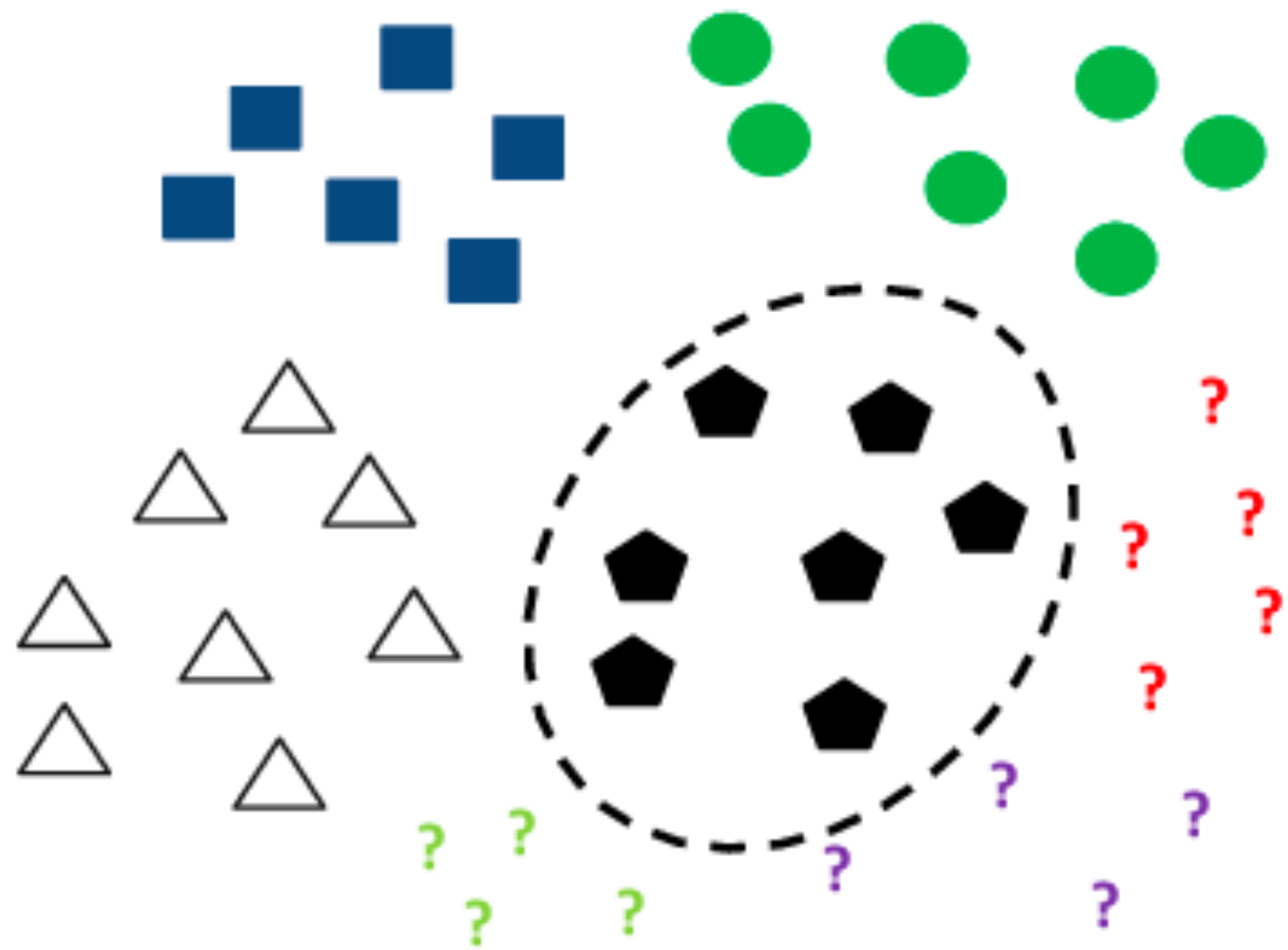
Efficient Speed and Memory Methods



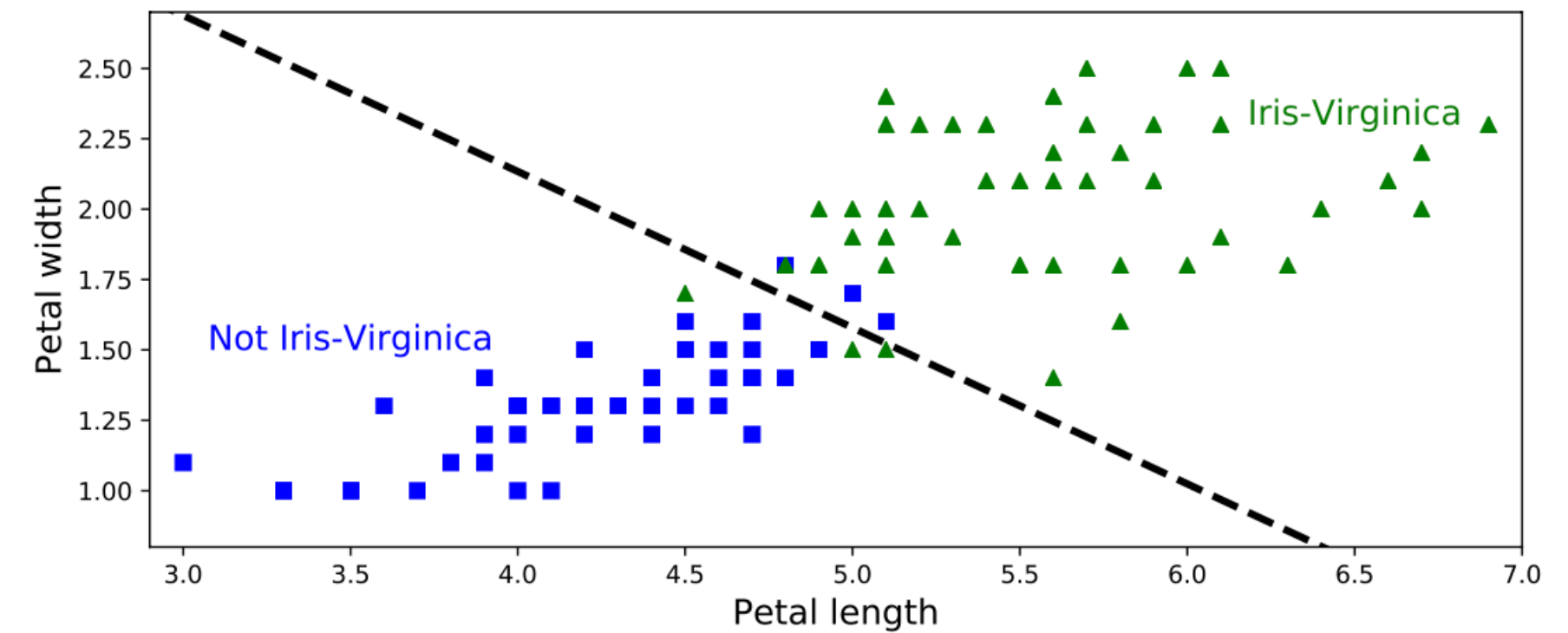
(1) Store only points that affect the decision boundaries.

(2) Hashing method, e.g., Locality-sensitive hashing.

Open Set Recognition



Open world / novelty detection / lifelong learning



$$f(\mathbf{x}; \boldsymbol{\theta}) = b + \mathbf{w}^T \mathbf{x} = b + \sum_{d=1}^D w_d x_d$$



$$\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$$

Closed world assumption

Learning Distance Metrics

$$d_{\mathbf{M}}(\mathbf{x}, \boldsymbol{\mu}) = \sqrt{(\mathbf{x} - \boldsymbol{\mu})^{\top} \mathbf{M} (\mathbf{x} - \boldsymbol{\mu})}$$

$\mathbf{M} = \mathbf{I}$ Euclidean Distance

↑
metrics learning

$e = f(\mathbf{x})$ When f is a DNN — deep metric learning

Large margin nearest neighbors

$$\mathcal{L}_{\text{pull}}(\mathbf{M}) = \sum_{i=1}^N \sum_{j \in N_i} d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_j)^2$$

$$\mathcal{L}_{\text{push}}(\mathbf{M}) = \sum_{i=1}^N \sum_{j \in N_i} \sum_{l=1}^N \mathbb{I}(y_i \neq y_l) [m + d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_j)^2 - d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_l)^2]_+$$

semidefinite programming

Deep Metric Learning

$$\mathbf{e} = f(\mathbf{x}; \boldsymbol{\theta}) \in \mathbb{R}^L$$

learn an embedding to a lower dimensional “semantic” space, suffer less from the “curse of dimensionality”

$$d(\mathbf{x}_i, \mathbf{x}_j; \boldsymbol{\theta}) = \|\hat{\mathbf{e}}_i - \hat{\mathbf{e}}_j\|_2^2$$

normalized Euclidean distance

$$\hat{\mathbf{e}} = \mathbf{e} / \|\mathbf{e}\|_2$$

$$d(\mathbf{x}_i, \mathbf{x}_j; \boldsymbol{\theta}) = \hat{\mathbf{e}}_i^\top \hat{\mathbf{e}}_j$$

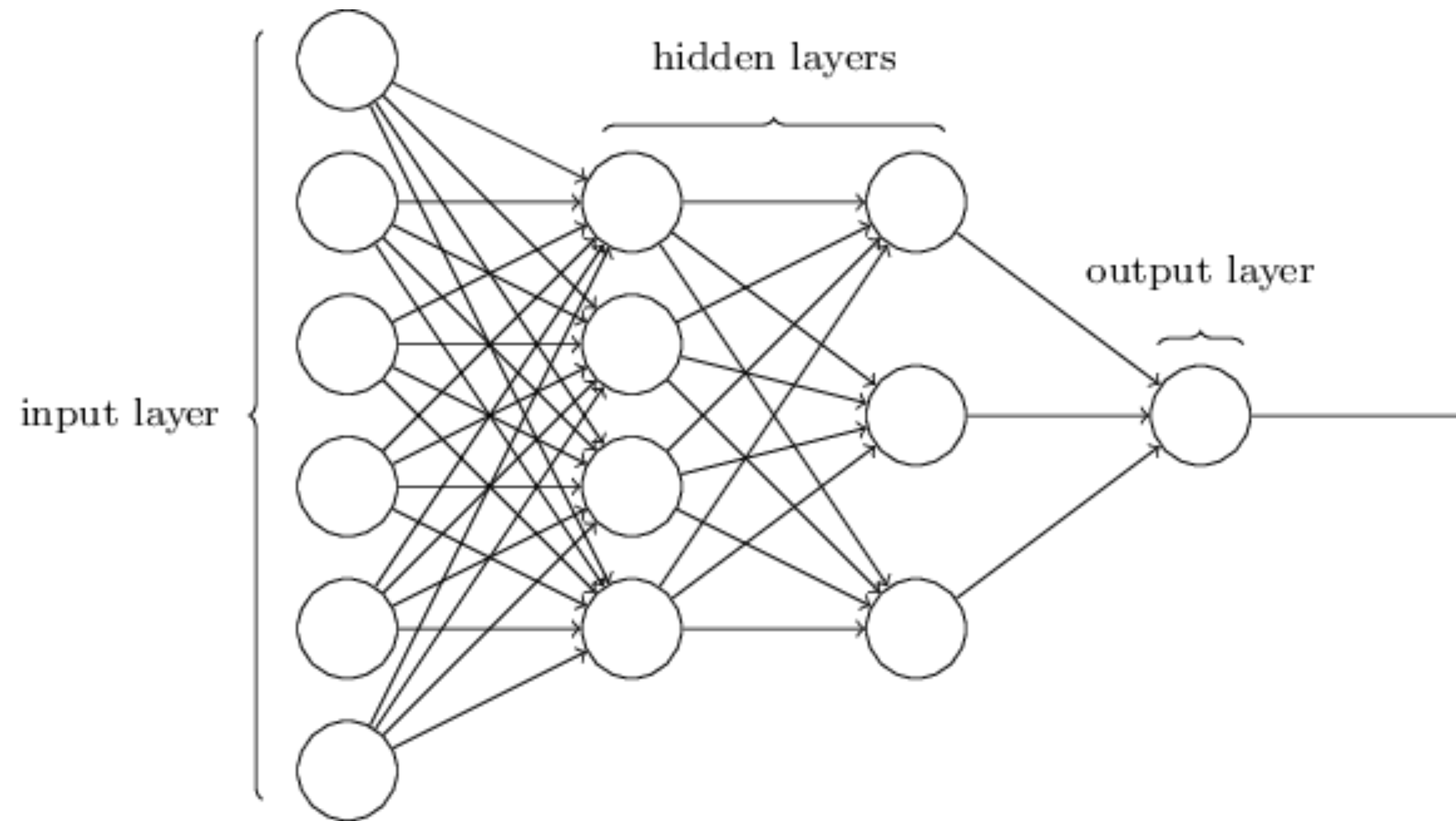
confine similarity

$$\mathcal{S} = \{(i, j) : y_i = y_j\}$$

i, j closer, i, k far away.

Classification Losses

$$e = f(\mathbf{x}; \boldsymbol{\theta}) \in \mathbb{R}^L$$



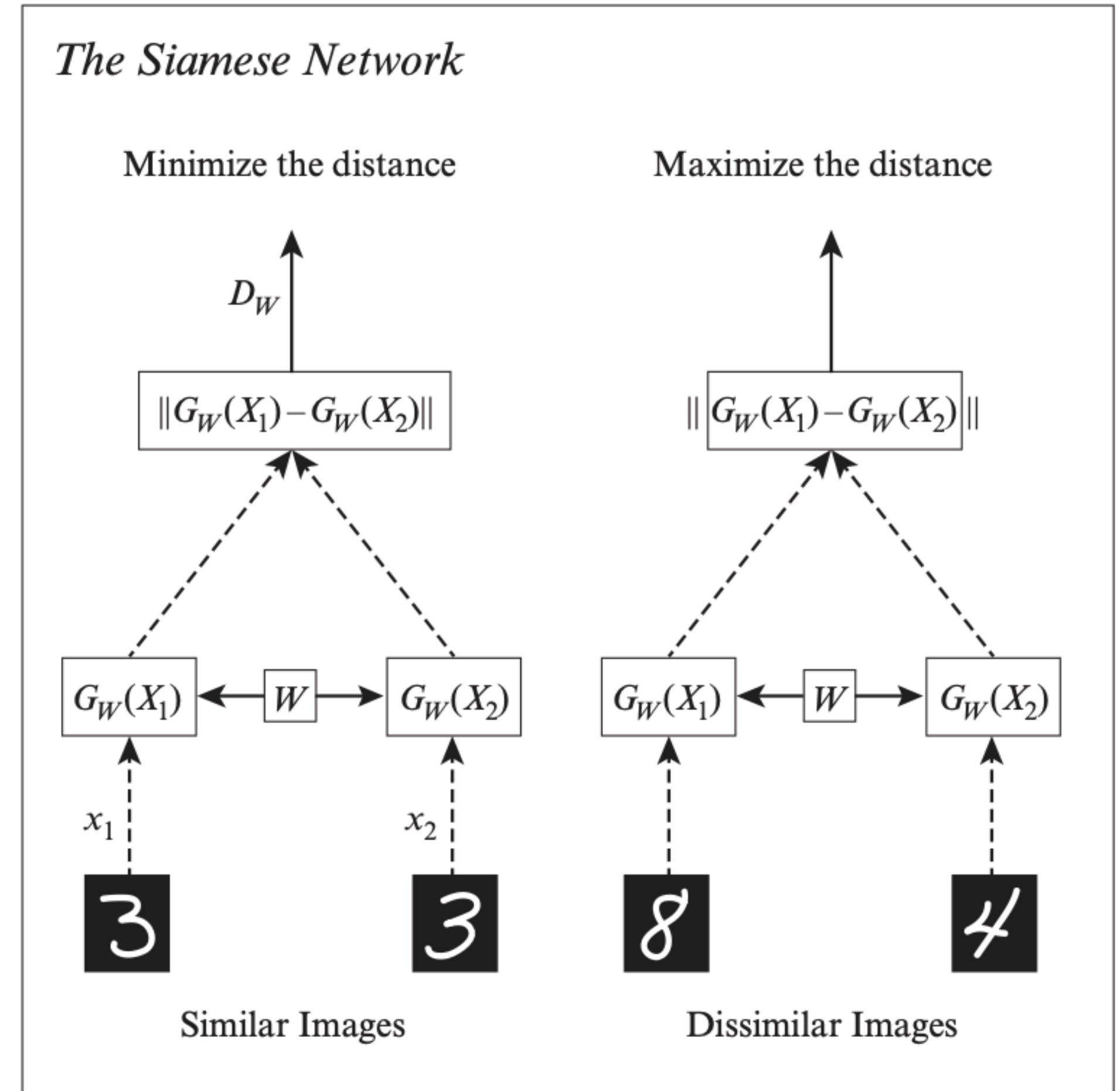
Learn a classifier and use the (second-to-last) layer as an embedding function.

Ranking Losses

$$\mathbf{e} = f(\mathbf{x}; \boldsymbol{\theta}) \in \mathbb{R}^L$$

Pairwise (contrastive) loss

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{x}_i, \mathbf{x}_j) = \mathbb{I}(y_i = y_j) d(\mathbf{x}_i, \mathbf{x}_j)^2 + \mathbb{I}(y_i \neq y_j) [m - d(\mathbf{x}_i, \mathbf{x}_j)]_+^2$$



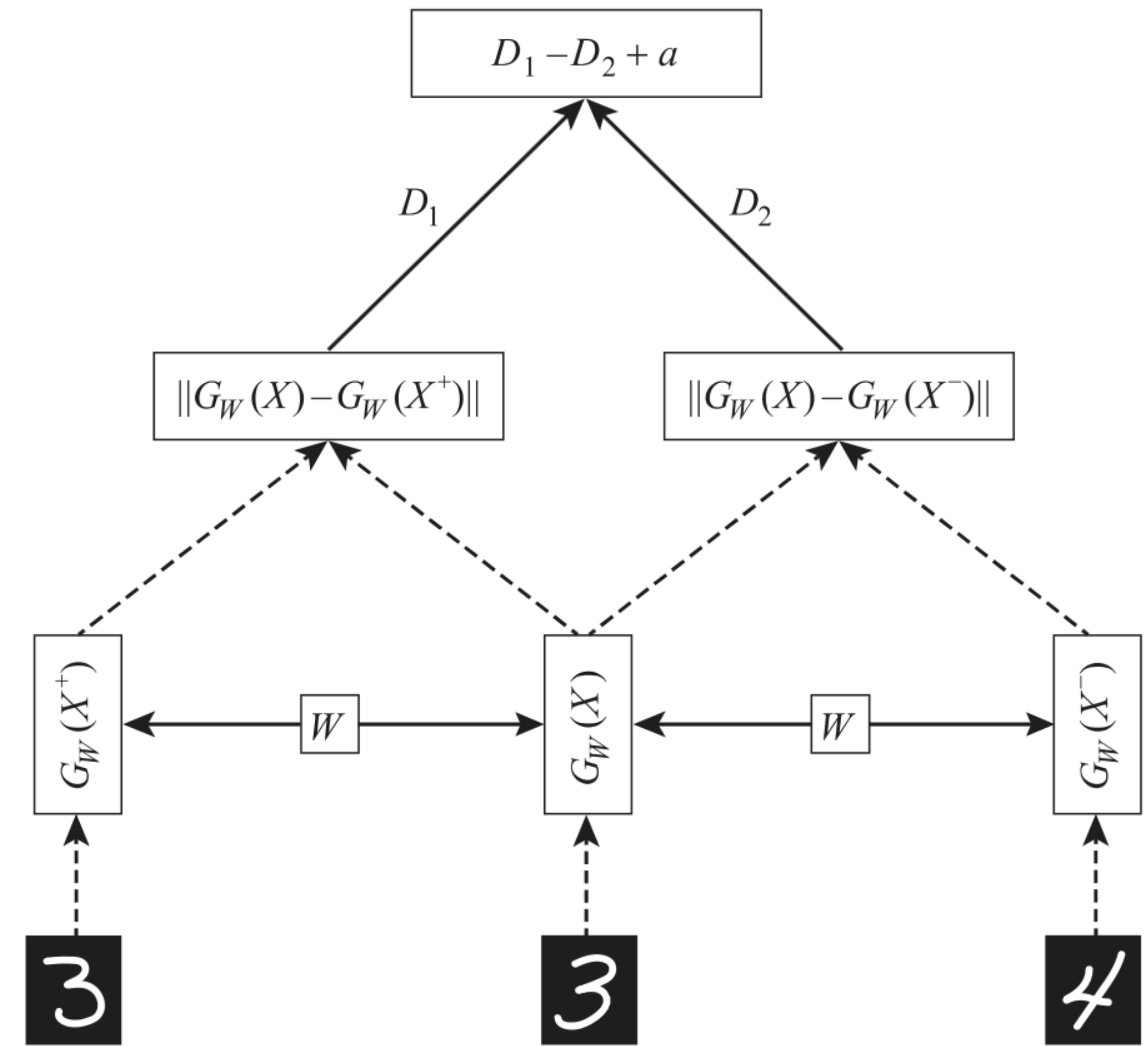
Triple Losses

$$e = f(\mathbf{x}; \boldsymbol{\theta}) \in \mathbb{R}^L$$

optimization of the positive pairs is independent of the negative pairs, which can make their magnitudes incomparable

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{x}_i, \mathbf{x}_i^+, \mathbf{x}_i^-) = [d_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{x}_i^+) - d_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{x}_i^-) + m]_+$$

The Triplet Network



N-pairs Loss (InfoNCE)

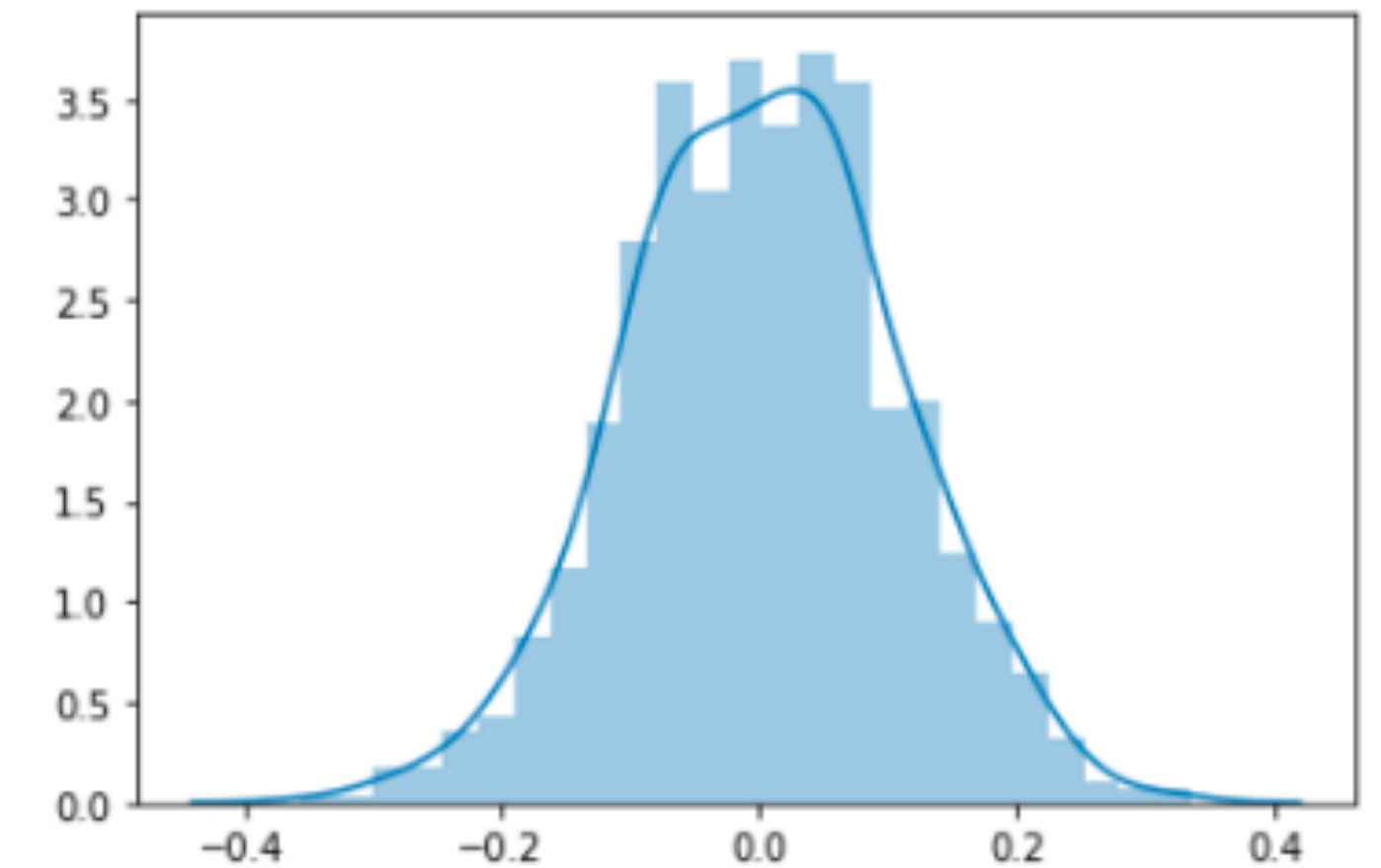
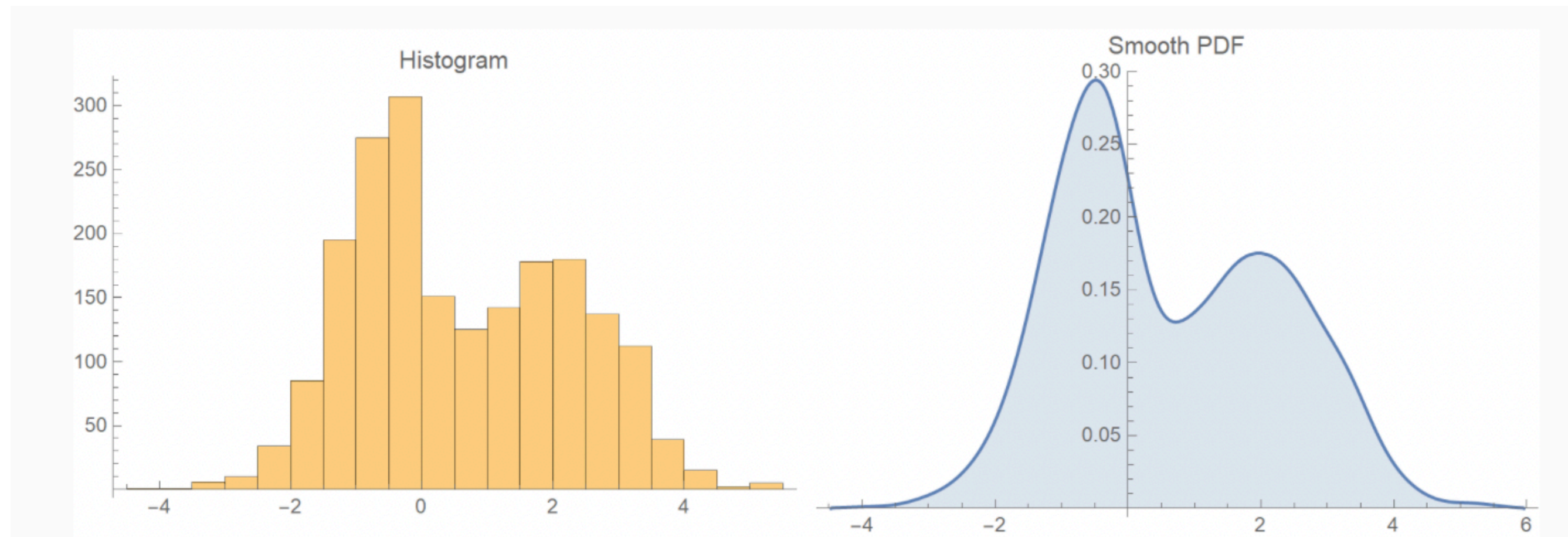
$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}; \mathbf{x}_i, \mathbf{x}_j^+, \{\mathbf{x}_k^-\}_{k=1}^{N-1}) &= \log \left(1 + \sum_{k=1}^{N-1} \exp(\hat{\mathbf{e}}_{\boldsymbol{\theta}}(\mathbf{x}_i)^{\top} \hat{\mathbf{e}}_{\boldsymbol{\theta}}(\mathbf{x}_k^-) - \hat{\mathbf{e}}_{\boldsymbol{\theta}}(\mathbf{x}_i)^{\top} \hat{\mathbf{e}}_{\boldsymbol{\theta}}(\mathbf{x}_j^+)) \right) \\ &= \log \frac{\exp(\hat{\mathbf{e}}_{\boldsymbol{\theta}}(\mathbf{x}_i)^{\top} \hat{\mathbf{e}}_{\boldsymbol{\theta}}(\mathbf{x}_j^+))}{\exp(\hat{\mathbf{e}}_{\boldsymbol{\theta}}(\mathbf{x}_i)^{\top} \hat{\mathbf{e}}_{\boldsymbol{\theta}}(\mathbf{x}_j^+)) + \sum_{k=1}^{N-1} \exp(\hat{\mathbf{e}}_{\boldsymbol{\theta}}(\mathbf{x}_i)^{\top} \hat{\mathbf{e}}_{\boldsymbol{\theta}}(\mathbf{x}_k^-))}\end{aligned}$$

Softmax function:

$$\text{softmax}(z_j) = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \text{ for } j = 1, \dots, K$$

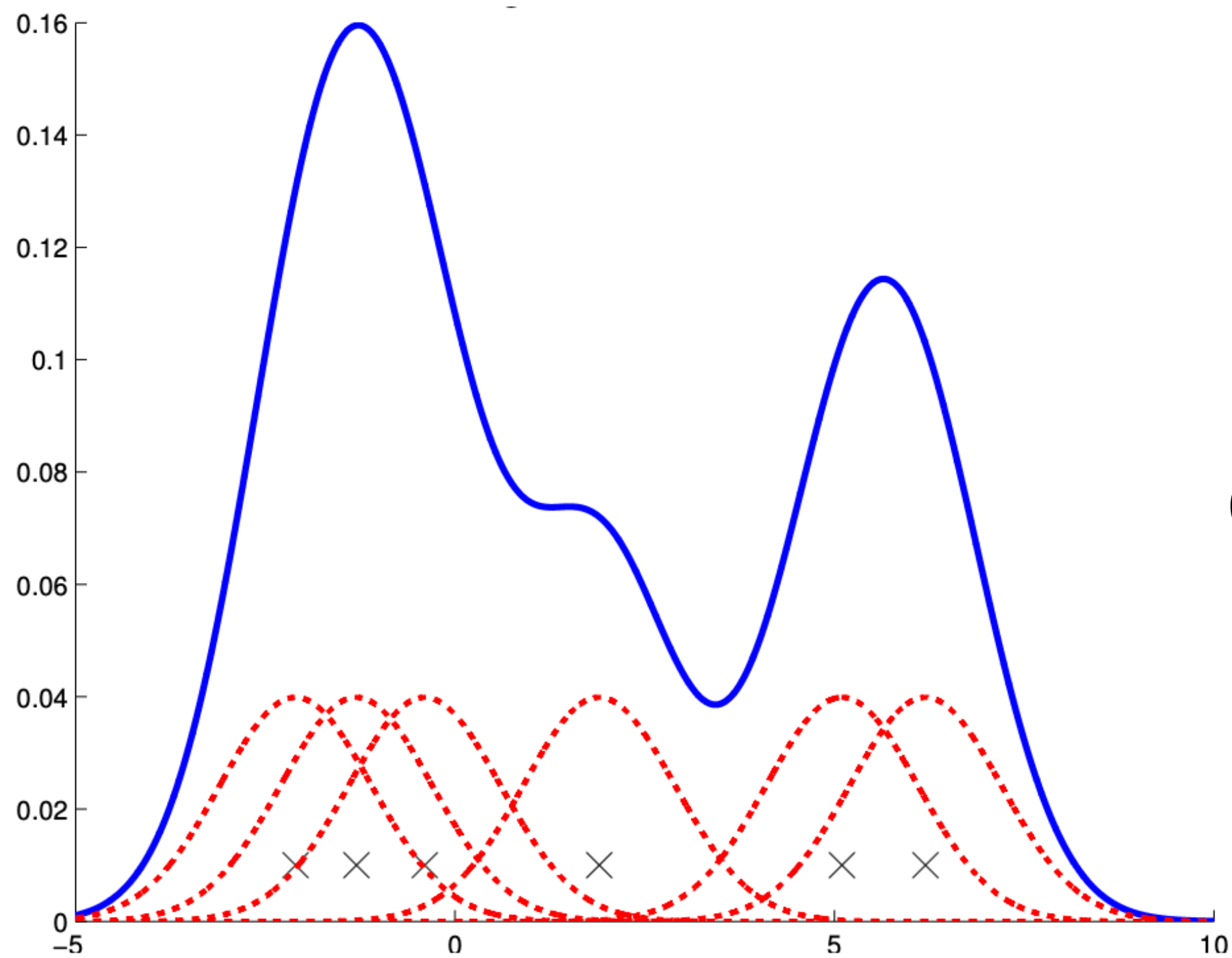
Kernel Density Estimation (KDE)

Density Estimation



Estimate $p(\boldsymbol{x})$ from seen data points

Kernel Density Estimation (KDE)



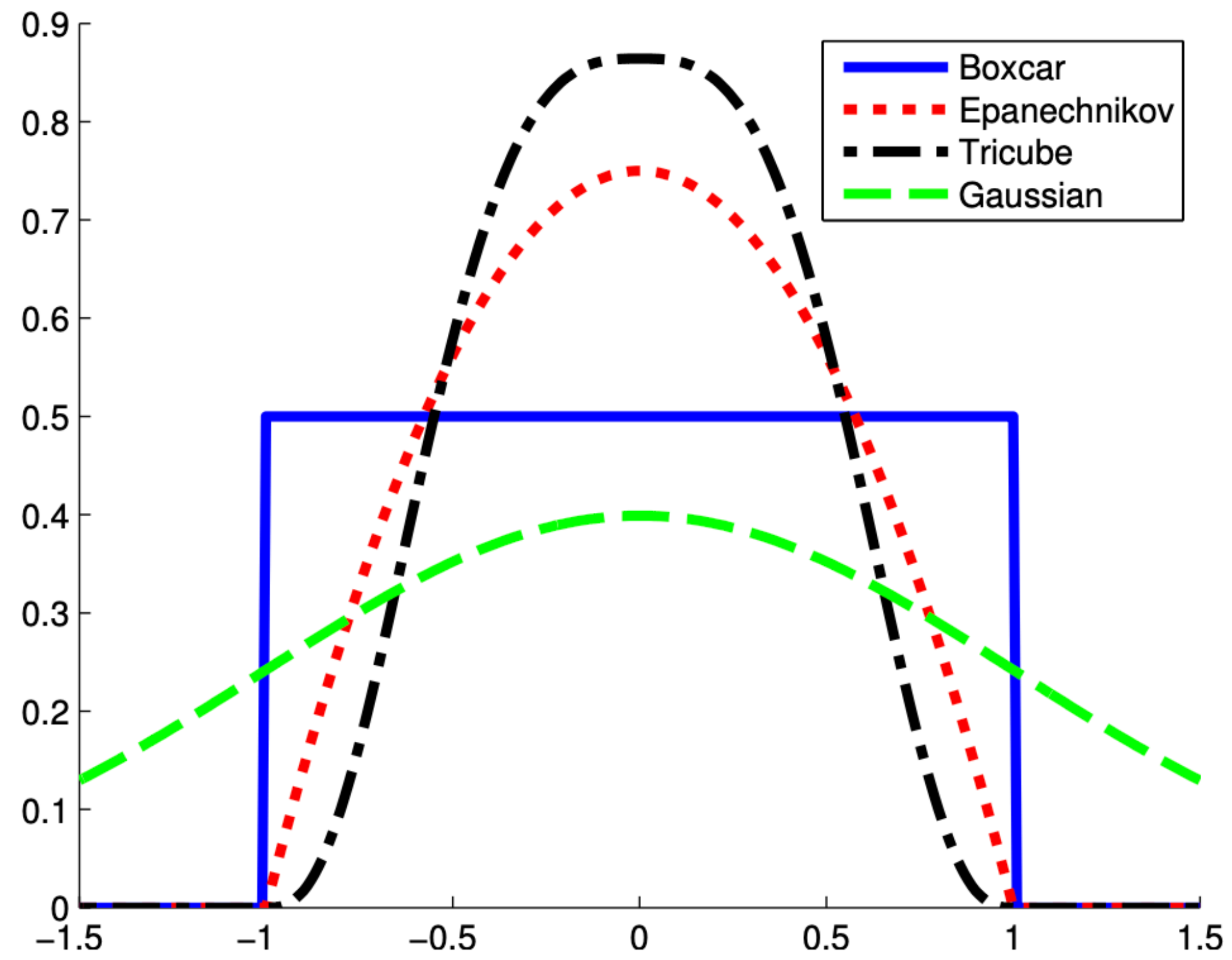
Parzen window density estimator:

$$p(\mathbf{x}|\mathcal{D}) = \frac{1}{N} \sum_{n=1}^N \mathcal{K}_h(\mathbf{x} - \mathbf{x}_n)$$

(In comparison) Gaussian mixture model:

$$p(\mathbf{x}|\boldsymbol{\theta}) = \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \sigma^2 \mathbf{I})$$

Density Kernels

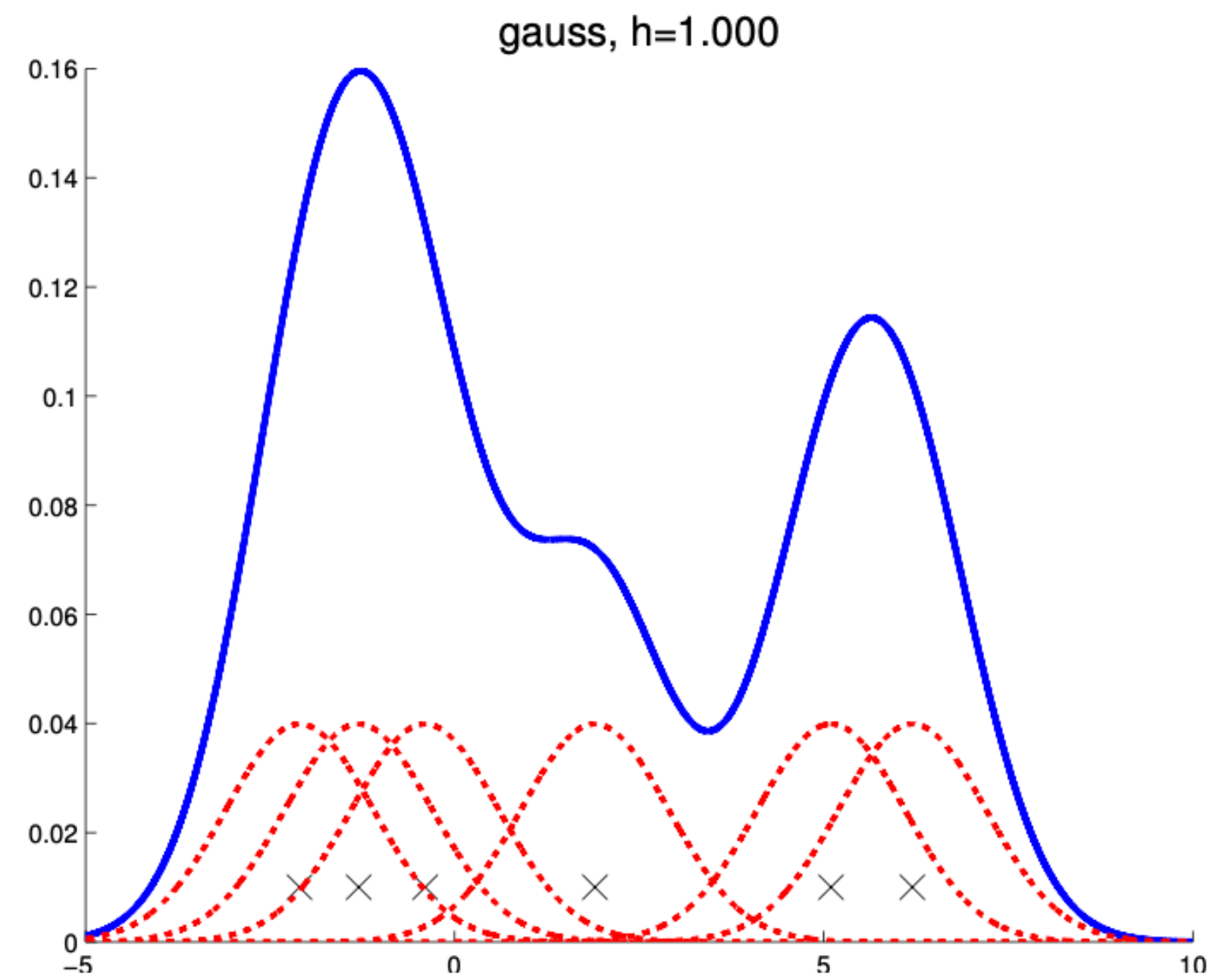
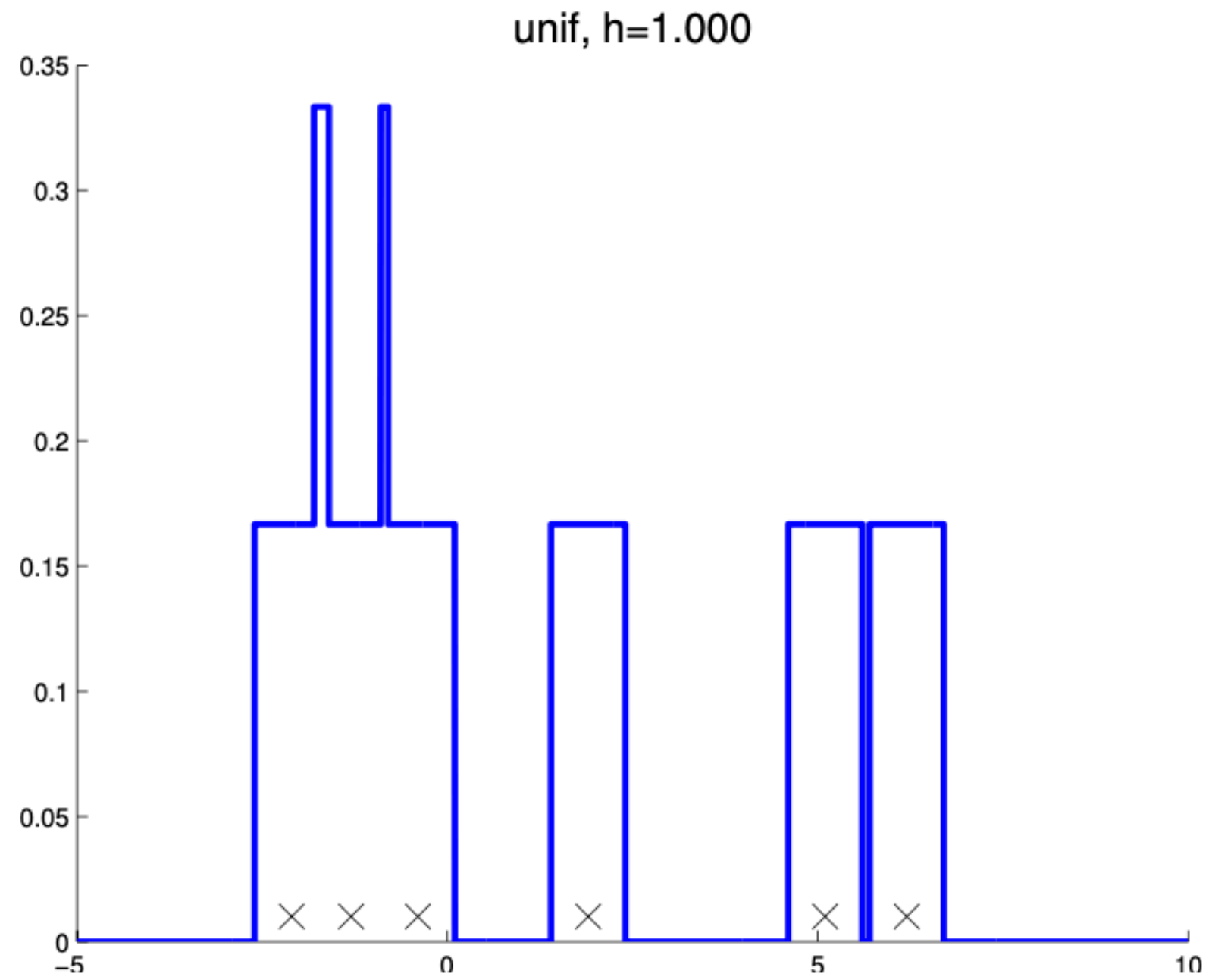


$$\mathcal{K} : \mathbb{R} \rightarrow \mathbb{R}_+$$

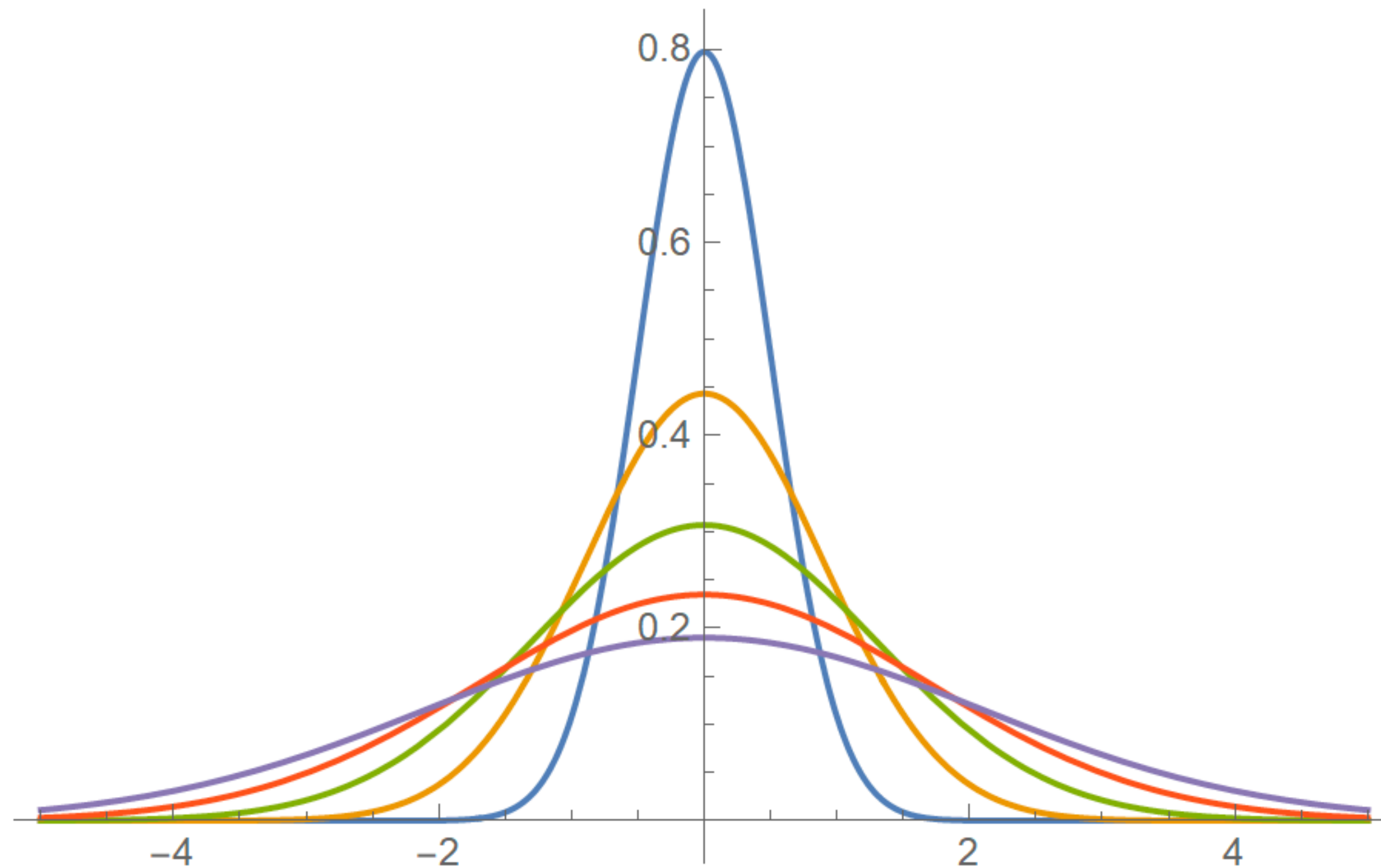
$$\int \mathcal{K}(x) dx = 1$$

$$\mathcal{K}(-x) = \mathcal{K}(x)$$

Nonparametric (Parzen) Density Estimator



Bandwidth

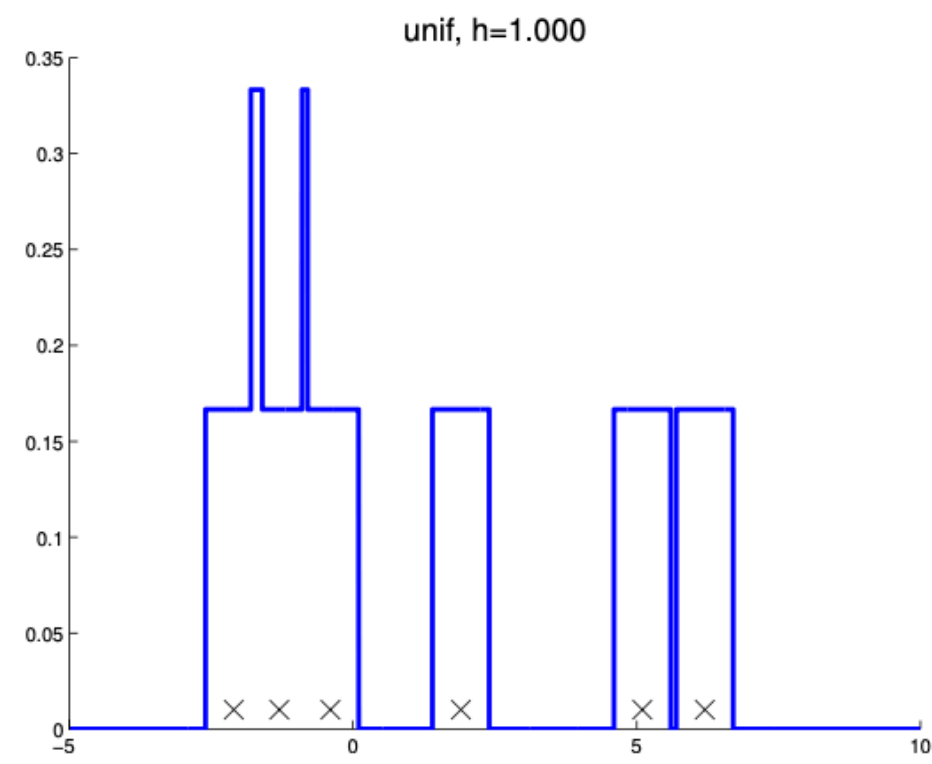


- h=0.5
- h=0.9
- h=1.3
- h=1.7
- h=2.1

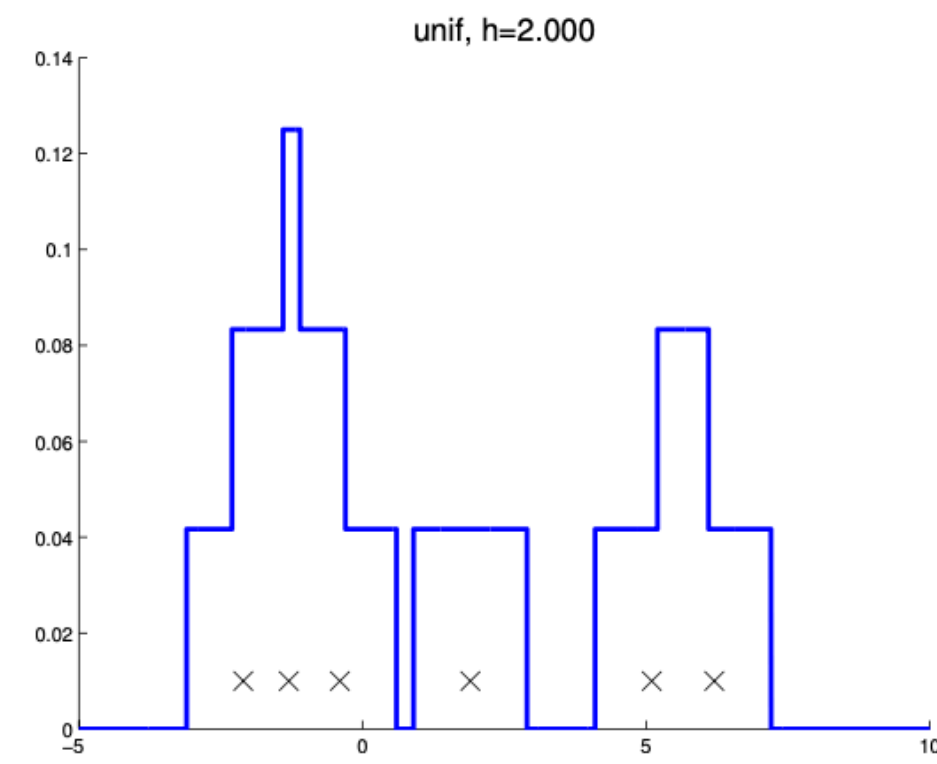
$$\mathcal{K}(x) = \frac{1}{(2\pi)^{\frac{1}{2}}} e^{-x^2/2}$$

$$\mathcal{K}_h(x) \triangleq \frac{1}{h} \mathcal{K}\left(\frac{x}{h}\right)$$

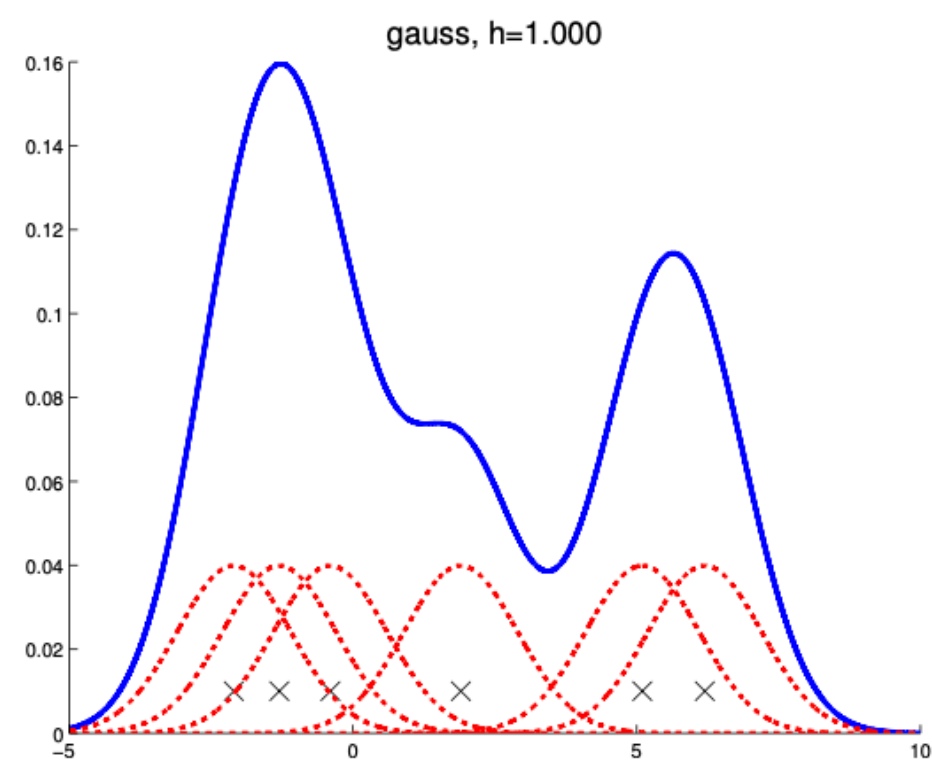
How to Choose Bandwidth



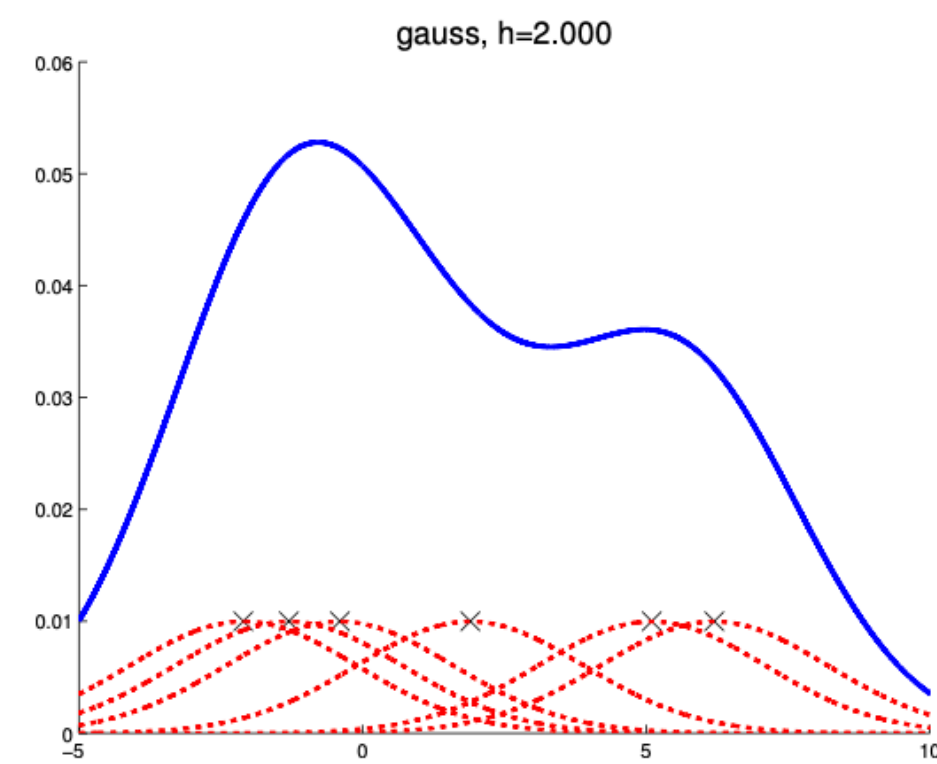
(a)



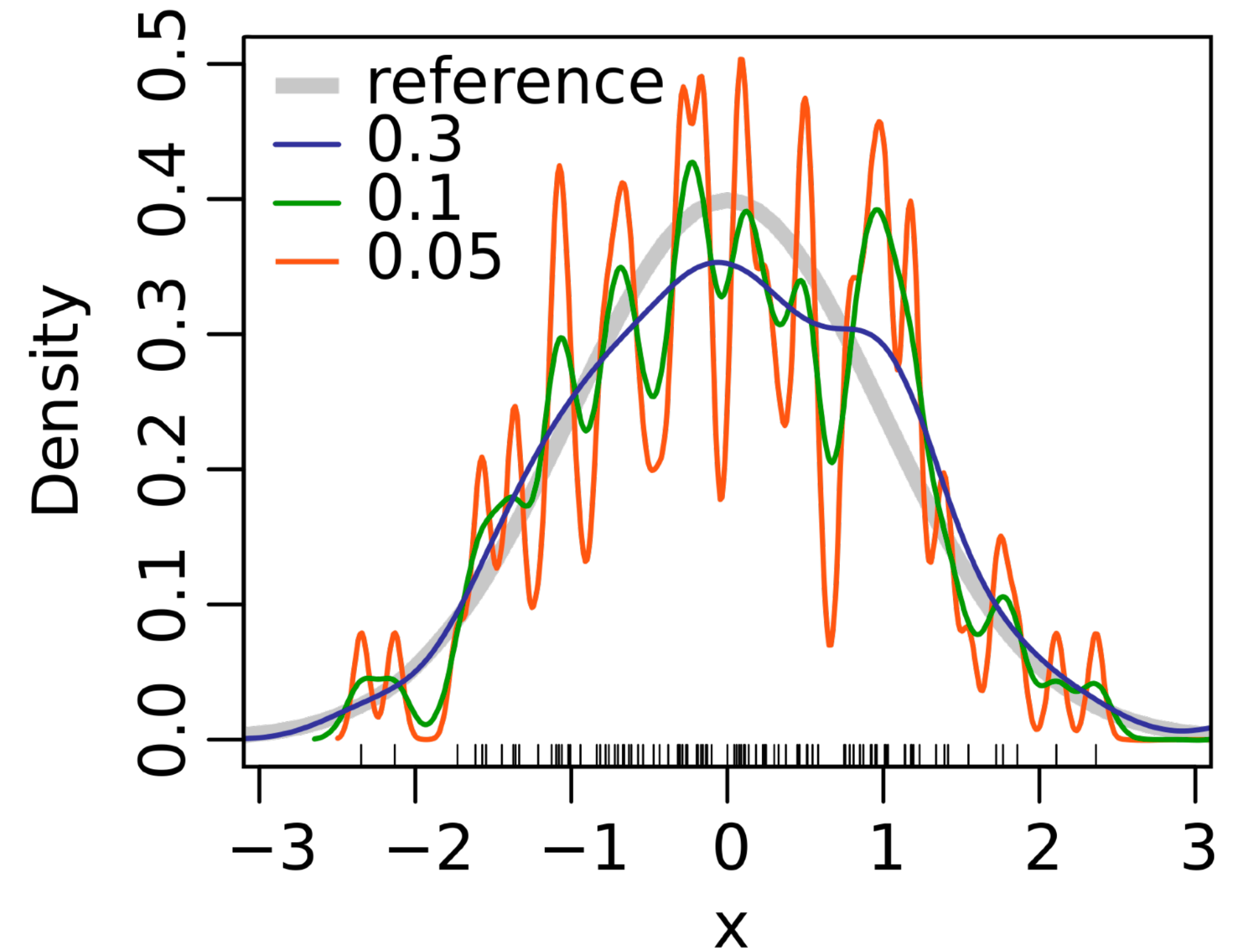
(b)



(c)



(d)



How to Choose Bandwidth

