



university of  
 groningen

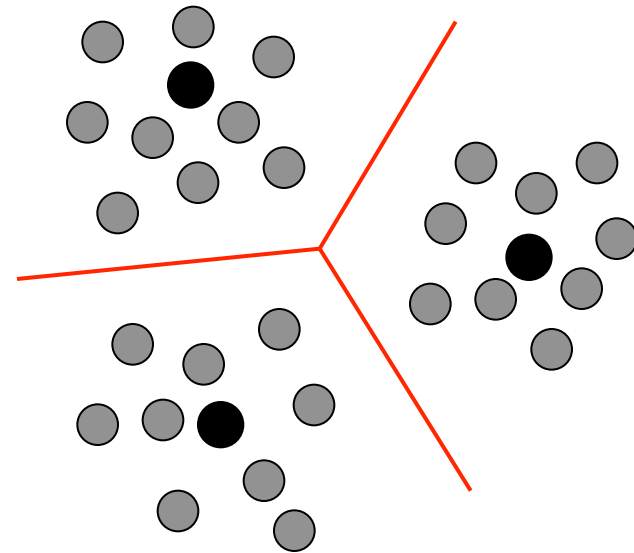
# Unsupervised Learning: Vector Quantization Competitive Learning and K-means algorithm

Introduction to Intelligent Systems  
Michael Biehl

aim:  
representation of large amounts  
of data by (few) **prototype vectors**

often used for:  
identification and grouping  
in *clusters* of similar data

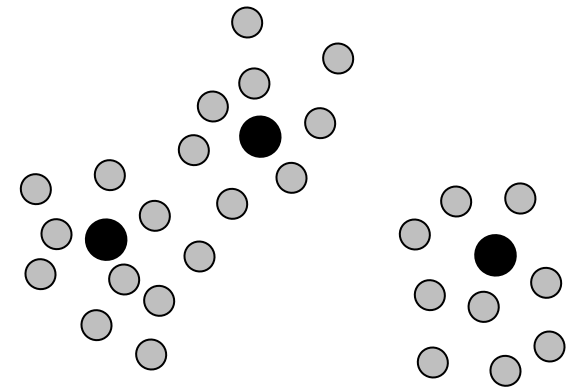
assignment of feature vector  $\xi$   
to the **closest prototype  $w$**   
(similarity or distance measure,  
e.g. Euclidean distance )





- initialize  $K$  prototype vectors

- present a single example
- identify the closest prototype, i.e the so-called *winner*
- move the winner even closer towards the example



intuitively clear, plausible procedure

- places prototypes in areas with high density of data
- identifies typical representatives of similar data
- similar to the K-means algorithm (see lectures by Nicolai Petkov)

VQ system: set of prototypes

$$w^1, w^2, \dots, w^K$$

$$w^k \in \mathbb{R}^N$$

data: set of feature vectors

$$\xi^1, \xi^2, \dots, \xi^P$$

$$\xi^i \in \mathbb{R}^N$$

assignment to prototypes:

based on dis-similarity/distance measure

$$d[w, \xi] \geq 0$$

given feature vector  $\xi$ , determine the *winner*

$$w^{i^*} = \operatorname{argmin}_j \{ d[w^j, \xi] \}$$

→ assign  $\xi$  to prototype  $i^*$

just one popular example: (squared) Euclidean distance

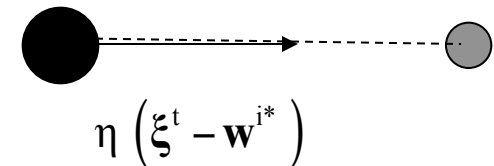
$$d[w, \xi] = \sum_{j=1}^N (w_j - \xi_j)^2$$

initially: randomized  $\mathbf{w}^k$ , e.g. equal to randomly selected data points

sequential presentation of example data  $\{\xi^t\} \quad t = 1, 2, \dots, P, 1, 2, \dots$

... *the winner takes it all*:  $\mathbf{w}^{i*} = \underset{j}{\operatorname{argmin}} \left\{ d[\mathbf{w}^j, \xi^t] \right\}$

$$\mathbf{w}^{i*} \rightarrow \mathbf{w}^{i*} + \eta \left( \xi^t - \mathbf{w}^{i*} \right)$$



$\eta$  ( $< 1$ ): learning rate, step size of update

competitive VQ (and K-means) aim at optimizing a cost function:

$$H_{VQ} = \sum_{j=1}^K \sum_{\mu=1}^P \underbrace{(\xi^\mu - w_j)^2}_{d_j^\mu} \underbrace{\prod_{k \neq j}^K \Theta(d_k^\mu - d_j^\mu)}_{w_j \text{ is the winner !}}$$

prototypes      data

here:  
Euclidean distance

- assign each data to one prototype
- measure the corresponding (squared) distance
- **quantization error** (sum over all data points)  
quantifies the quality of the representation

$$\Theta(x) = \begin{cases} 1 & \text{for } x \geq 0 \\ 0 & \text{else} \end{cases}$$

defines a (one) criterion to evaluate / compare the quality of different prototype configurations (important: only for fixed K, see below!)

# Lloyd's algorithm (*K-means*)

(0) initialization:

place vectors  $w^k$ , e.g. equal to randomly selected points

(1) assignment of data to centers:

assign every data point to nearest center  
(e.g. according to Euclidean distance)

(2) re-compute centers:

compute the means of the K clusters

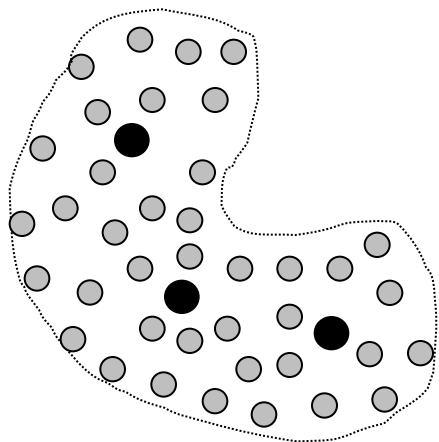
comparison:

K-means: updates all prototypes, considers all data at a time  
( batch- or offline-optimization )

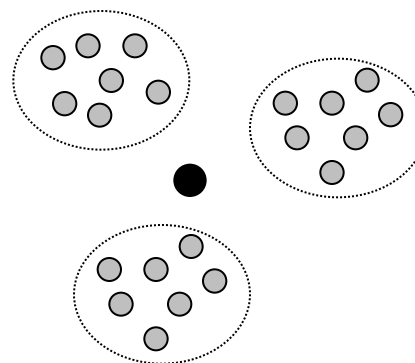
VQ-alg. : updates only the winner, random sequential presentation of  
single examples (stochastic or on-line optimization)



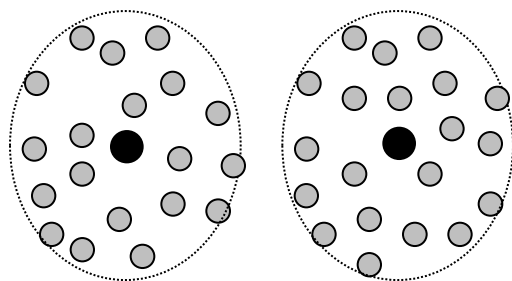
Remark 1: VQ / K-means may be related to clustering  
but it is not quite the same



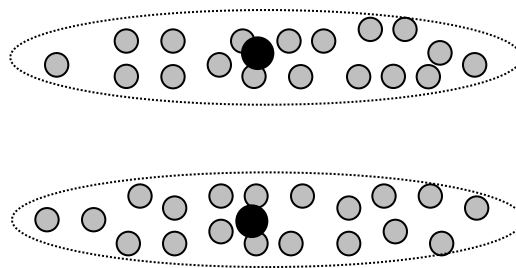
one “cluster” represented  
by several prototypes



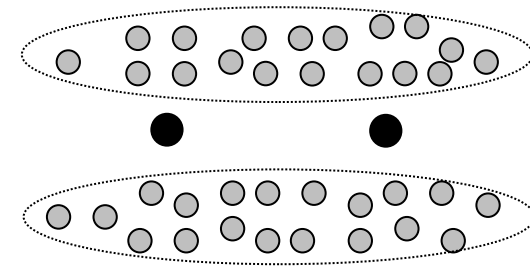
several “clusters”  
represented  
by one prototype



intuitive clustering  
= low  $H_{VQ}$



intuitive clustering

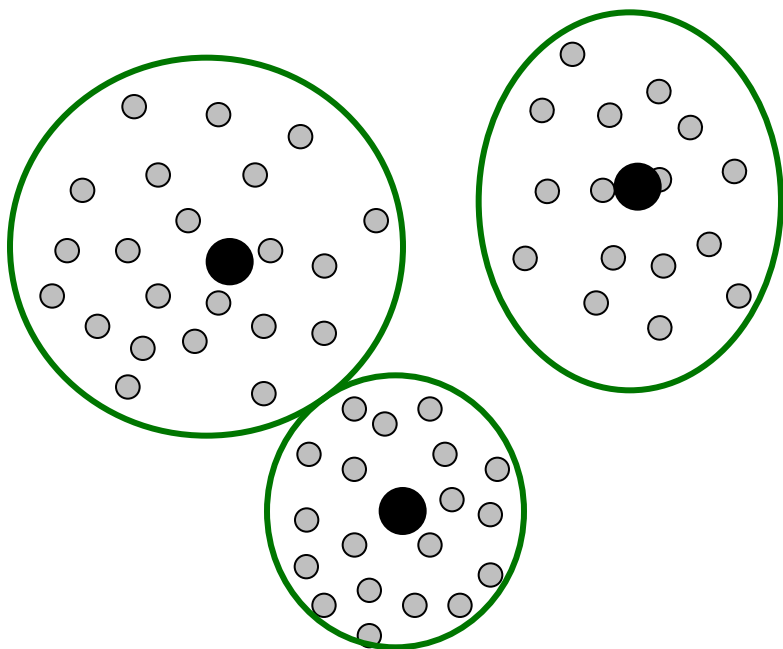


lower  $H_{VQ}$  !

VQ/K-means result depends on

- the **number of prototypes** (predefined!)
- shape and relative position of clusters
- the **distance measure** / metric used
- **representation of the data**, e.g. scaling, transformations, normalizations etc. !!!

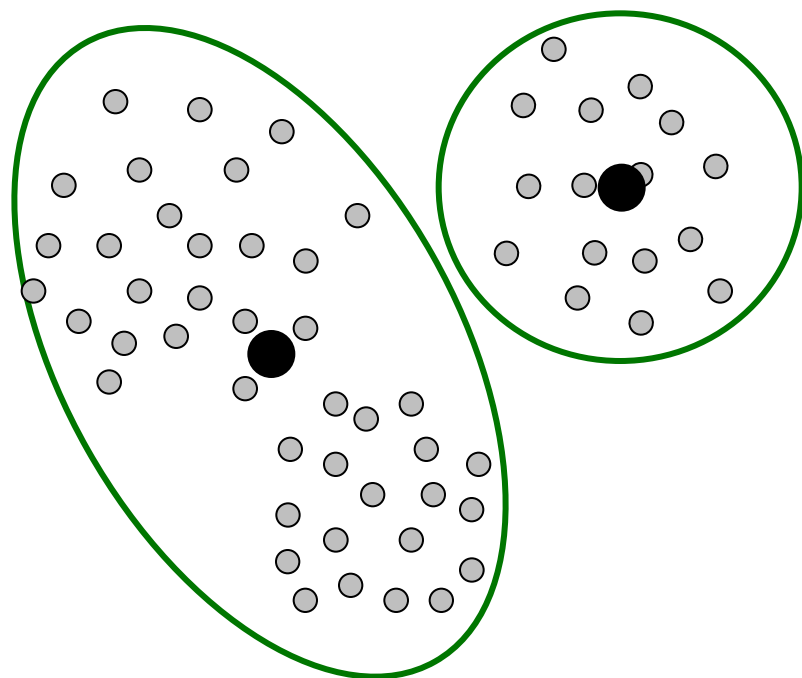
Remark 2: clustering is an ill-defined problem!



“obviously three clusters”

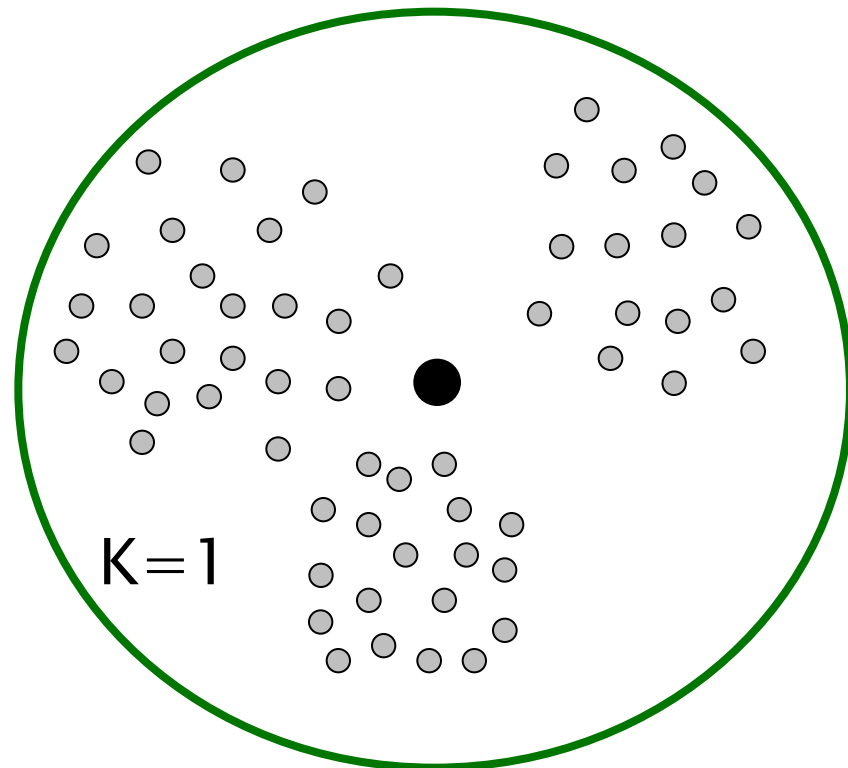
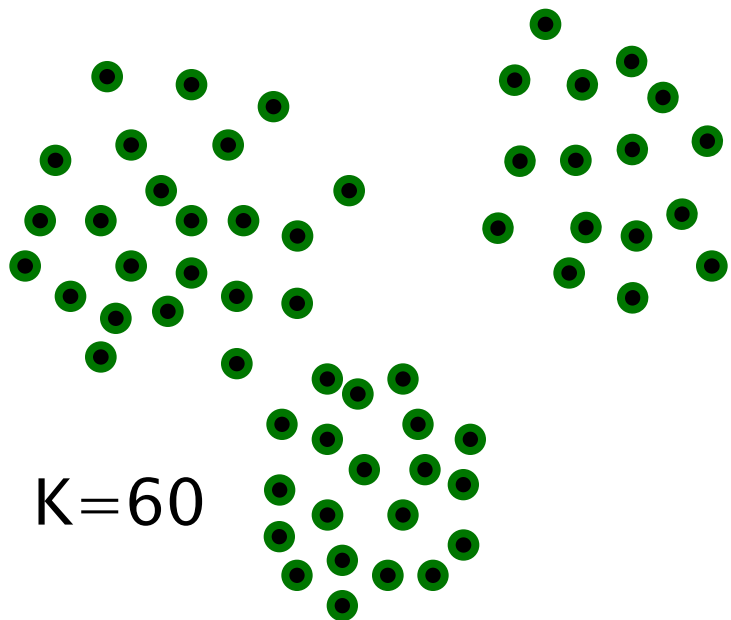
our criterion: lower  $H_{VQ}$

→ “better clustering” ???



“well, maybe only two?”

higher  $H_{VQ}$



$$H_{VQ} = 0$$

→ “the best clustering” ???

the simplest “clustering” ...

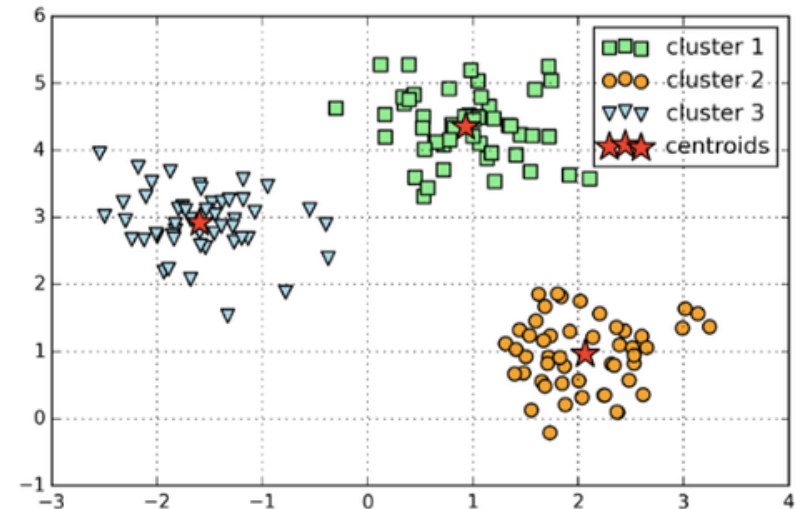
$H_{VQ}$  (and similar criteria) allow only to compare VQ with the same  $K$  !

more general: heuristic compromise between “error” and “simplicity”

popular heuristics: **elbow method**

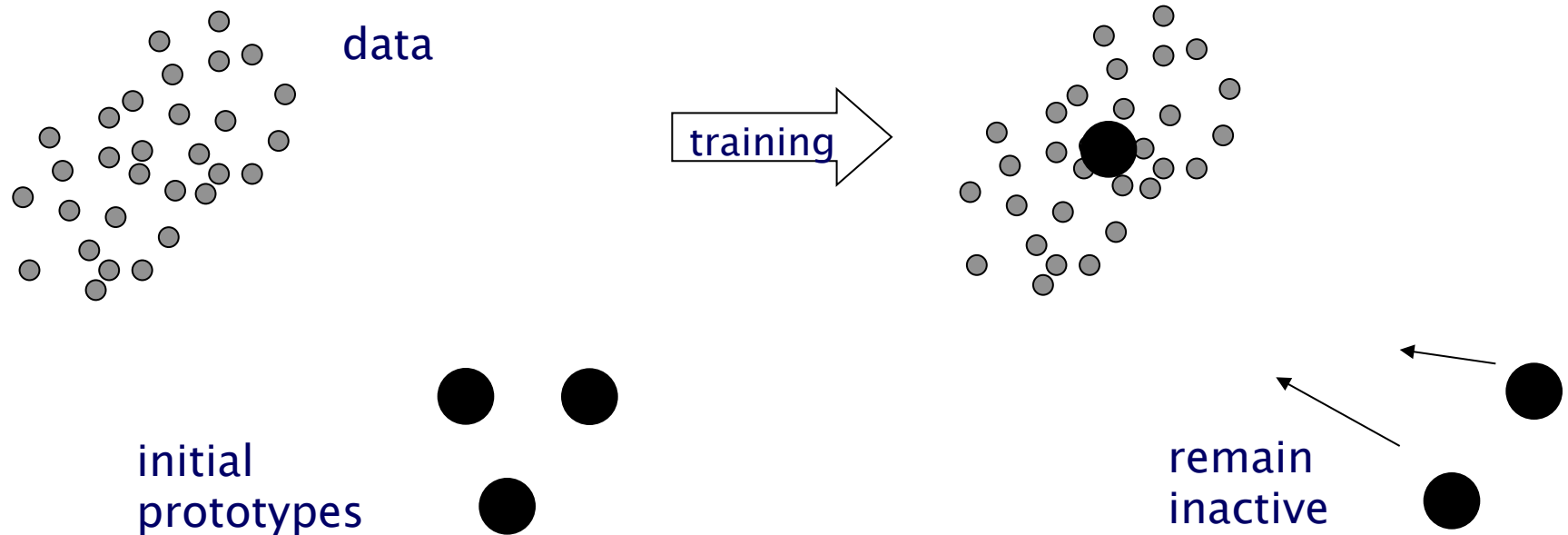
- run VQ for different K to convergence
- determine quantization error as a function of K
- identify “**elbow**” – characteristic values of K

very clear example (in practice usually less pronounced)



<http://sebastianraschka.com>  
Machine Learning Blog

one problem of *winner takes all* prescriptions



solution: rank-based updates (winner, second, third, ...)

more general problem (also K-means!)

- **local minima** of the quantization error
- initialization-dependent outcome of training