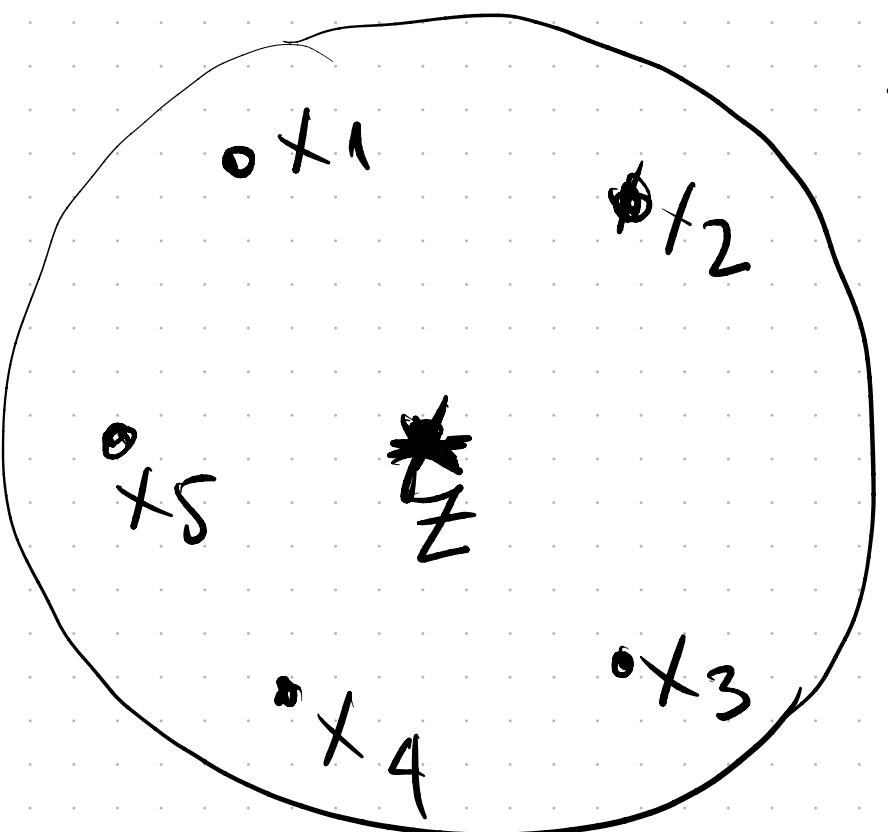


Lecture 7/18

- **KNN**
 - optional: tSNE
- Kernel PCA
- HW5 demo

K Nearest Neighbor: For each point $z \in \text{Test}$
define neighbourhood $N_z = \{\text{set of closest points from training to } z\}$



- predict label / class / target by estimating
from N_z

classification: predict class y_z that is
most present in N_z $y_z = \underline{\text{mode}}\{N_z\}$
frequent

quantity/regression: predict $\begin{array}{l} \text{AVG}\{N_z\} \\ \text{(label)} \end{array}$

ex $y = \text{house price} \Rightarrow \text{pred } y_z = \begin{array}{l} \text{AVG } \{y_i\} \\ x_i \in N_z \end{array}$

Not really training, just pred $f_Z = \text{estimate from closest neighbors to } Z$

• Need dist/similarity (kernel)

across datapoints

$$k_{ij} = \text{sim}(x_i, x_j)$$

$$k_{iz} = \text{sim}(x_i, z)$$

appropriate for data
and for task

ex $x_i = \text{patients} \Rightarrow k_{ij} = \text{similarity of patients w.r.t. diabetes}$

$\underline{x_i = \text{images}} \Rightarrow k_{ij} = \text{similarity of image prod}$
 $y = \text{prod. price}$ w.r.t. price (willingness to pay)

$\underline{x_i = \text{email}} \Rightarrow k_{ij} = \text{similarity of emails w.r.t. spam}$

$\underline{x_i = \text{movies}} \Rightarrow k_{ij} = \text{similarity of movies w.r.t. user satisfaction}$

$$y = \text{rating}$$

KNN 3 variants

- $K = \text{fixed}$; ex $\Rightarrow K=5$. N_Z always has 5 training points.

rank all training points
by k_{zi} similarity.
- select top $K \rightarrow N_Z$

$$N_Z = \{ \text{closest 5 to } z \}$$

disadvantage: some z don't have 5 close neighbors (low density)

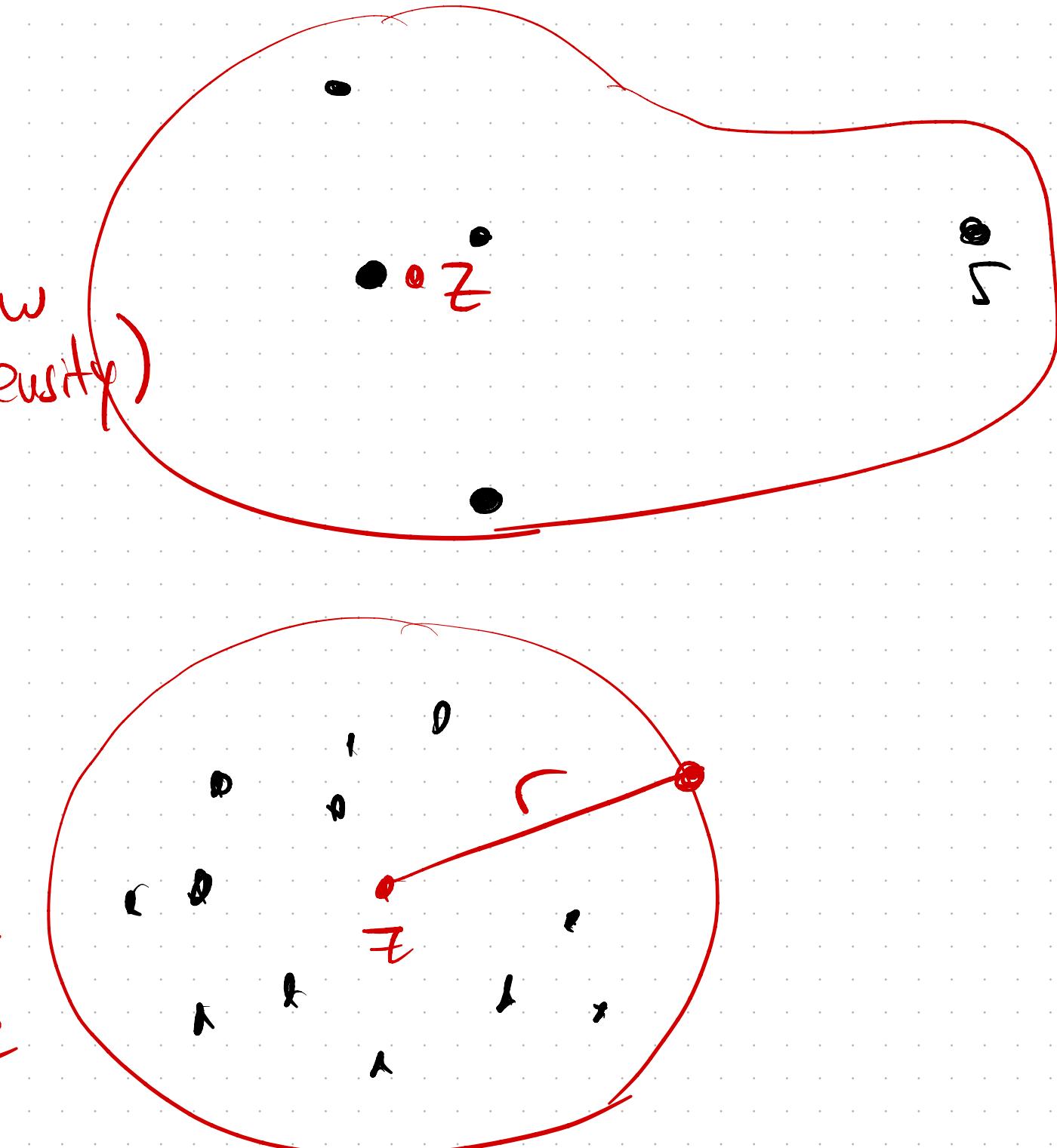
range

- range = fixed & variable

$$N_Z = \{ x_i \mid \begin{array}{l} k_{iz} \leq r \\ \text{training} \\ \text{dist}(i, z) \leq r \end{array} \}$$

disadvantage: some z few neighbors in range

some z : many neighbors in range



• Use all training points weighted by similarity $k_{iZ} = \text{sim}(x_i, z)$

regression predict (z) = $\frac{\sum_{i=1}^N k_{iZ} \cdot y_i}{\sum_{i=1}^N k_{iZ}} \rightarrow \text{labels / targets}$

classification

score (z, y_e) = $\sum_{i=1}^N k_{iZ} \cdot [1_{[y_i = y_e]}]$ filter

class label \uparrow

N_Z = all training set, but weighted.

$k_{iZ} = \text{high} \Rightarrow x_i \approx z \Rightarrow y_i$ counts a lot

$k_{iZ} = \text{low} \Rightarrow x_i$ not similar $\Rightarrow y_i$ doesn't count much

Kernel PCA

rewrite PCA primal \Rightarrow dual form.

watty

$$\text{PCA word} = X \cdot \begin{bmatrix} e_1 & e_2 & \dots & e_p \end{bmatrix}$$

eigen vectors (Σ_{covar})
sorted by desc eigen val.

?

$$= \underbrace{X X^T}_{K} \cdot \underbrace{B}$$

dual variables.

choose
↳ can use any valid kernel K (not just linear X^T)

$$- \|x_i - x_j\|^2 / 2\sigma^2$$

$$\text{ex } K = \text{gaussian} = e$$

compute B for that kernel.

$$B = \text{Eigenvect}(K)$$

compute

Intuition we want eigen vectors
 (primal var) = $X^T \boxed{B}$ dual.

detail (skip) : show that $e = X^T B$ possible for every e

calculate/find B ? e eigen vector of Σ $\Rightarrow \Sigma e = \lambda e$
 CORAD eigen val

$$\Sigma e = \lambda e$$

$$\left(\frac{1}{N} X^T X \right) X^T B e = \lambda X^T B e$$

$$\Sigma = \frac{1}{N} X^T X$$

$$e = X^T B$$

CORAD estimation

want

$$X^T X X^T B = N \lambda X^T B$$

$$X^T X X^T B = N \lambda \boxed{X^T B}$$

left

$$K \cdot K^a B = N \lambda K \cdot B$$

$$N \lambda \text{ scalar}$$

$B = \text{eigenvector}(K)$
 $N \lambda = \text{corresp. eigenval}$

t-SNE representation of data $D\text{-dim} \Rightarrow G\text{-dim } (G \ll D)$

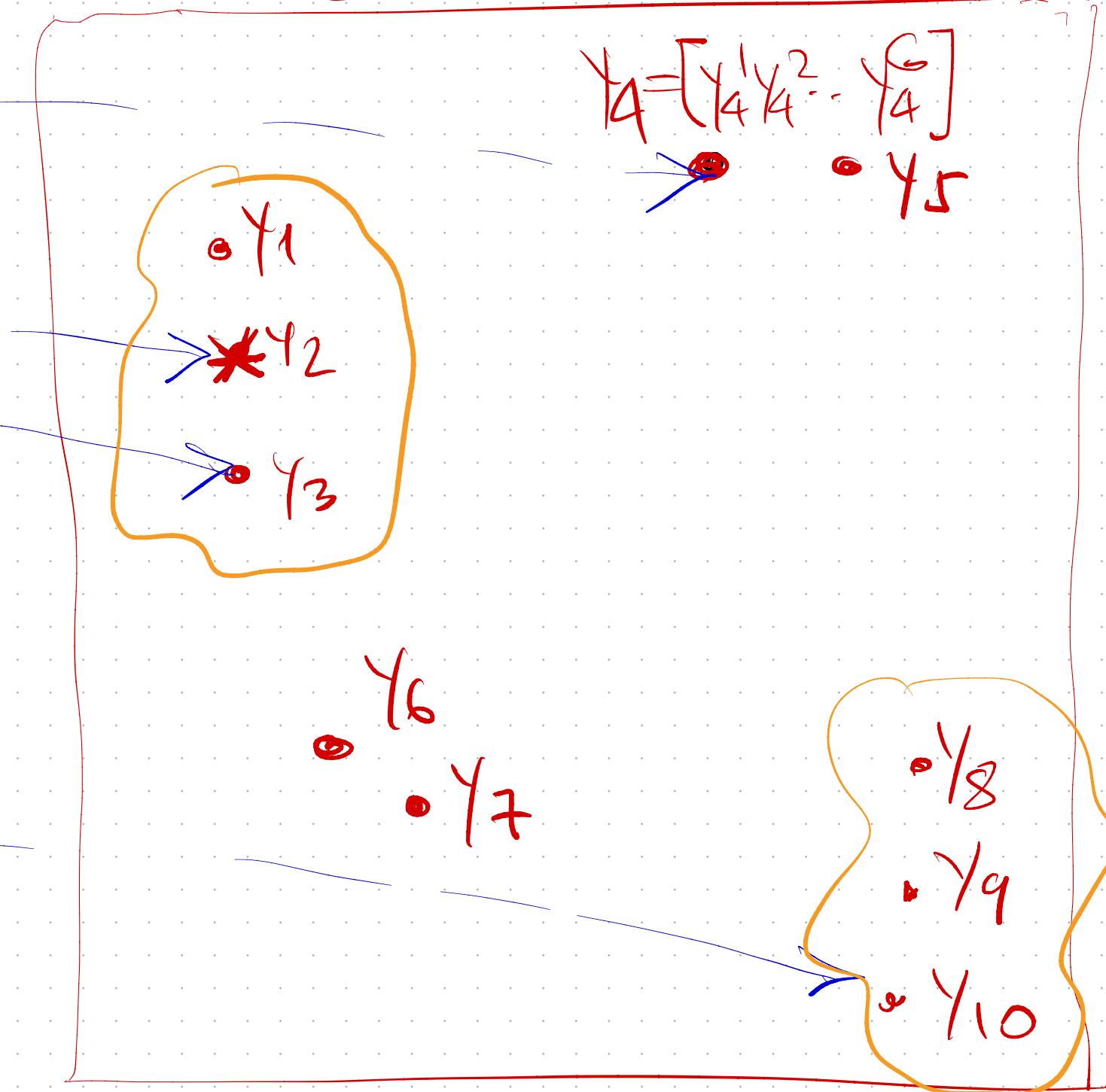
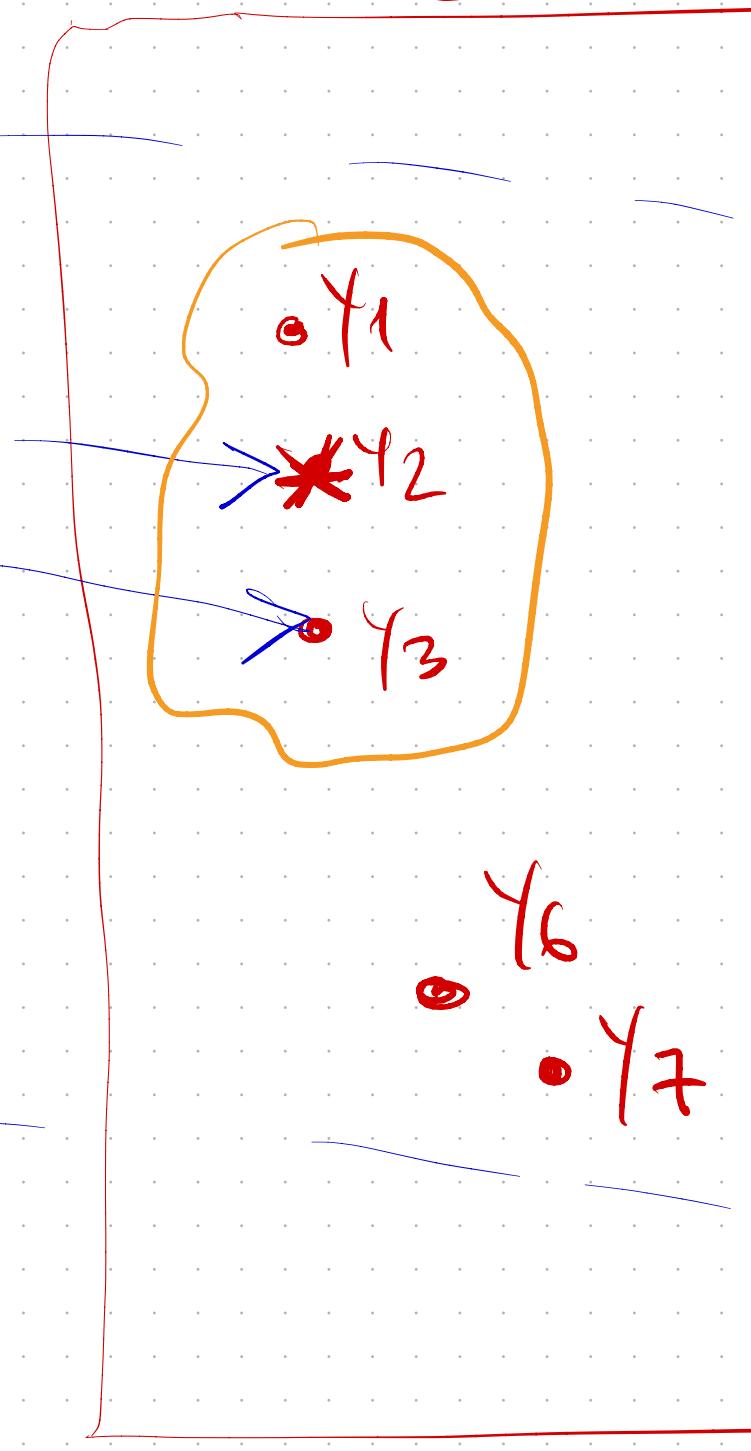
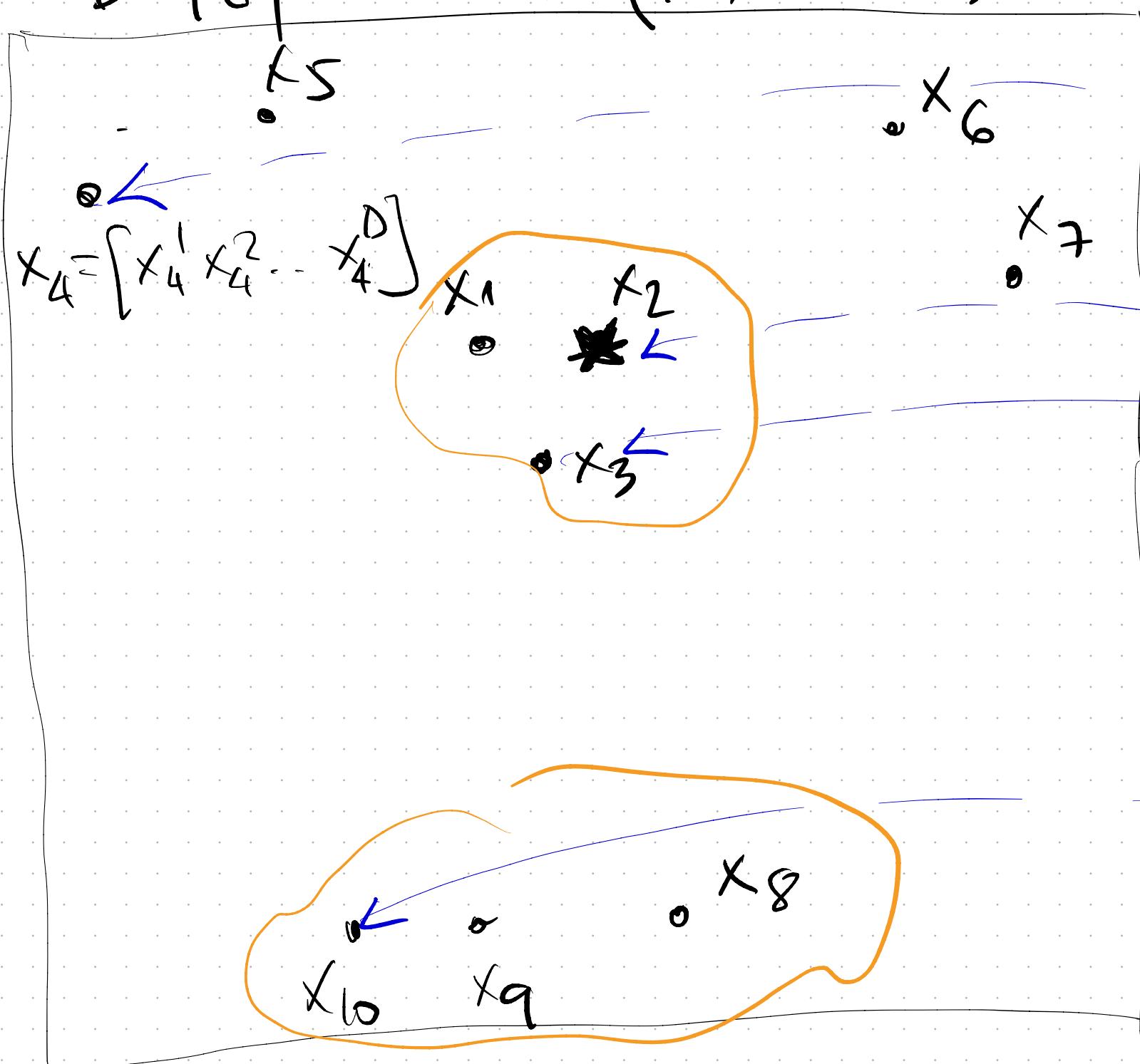
$D\text{-dim } X\text{-space (original)}$
 $D=784$
 $X = [x^1 \ x^2 \ \dots \ x^D]$

$$X_i \Rightarrow Y_i$$

$G\text{-dim } Y\text{-space (new)}$

$$G=3$$

$$Y = [Y^1 \ Y^2 \ \dots \ Y^G]$$



- do not worry about exact distances / positioning (in algebra)

$$\text{dist}(x_1, x_5) \quad ?? \quad \text{dist}(y_1, y_5)$$

$$\text{angle}(x_1, x_9, x_7) \quad ?? \quad \text{angle}(y_1, y_9, y_7)$$

= care about: preserve "close" vs "far" neighbors.

x_2 much closer to
 x_3 than to x_{10} \iff y_2 much closer to y_3
than to y_{10}

one time
Fixed

P = distribution of distances (x)

$$P_{ij} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma^2)}{\text{Normalized}}$$

Q = distribution of distances (y)
not gaussian. Use t-distib.

$$Q_{ij} = \frac{1}{[1 + \|y_i - y_j\|^2]}$$

$$Q_{ij} = \frac{1}{N \times N}$$

not fixed, updated with y

- Given X_{NxD} compute $P_{N \times N}$ (Search for "width" t^2 by perplexity)

- Want Y_{NxD} such that $\text{dist}_{\text{dist}}(Q(Y)) \approx P(X)$

$$\text{MIN OBJ} = \text{KL}(P || Q) = \sum_{ij} P_{ij} \log \frac{P_{ij}}{Q_{ij}}$$

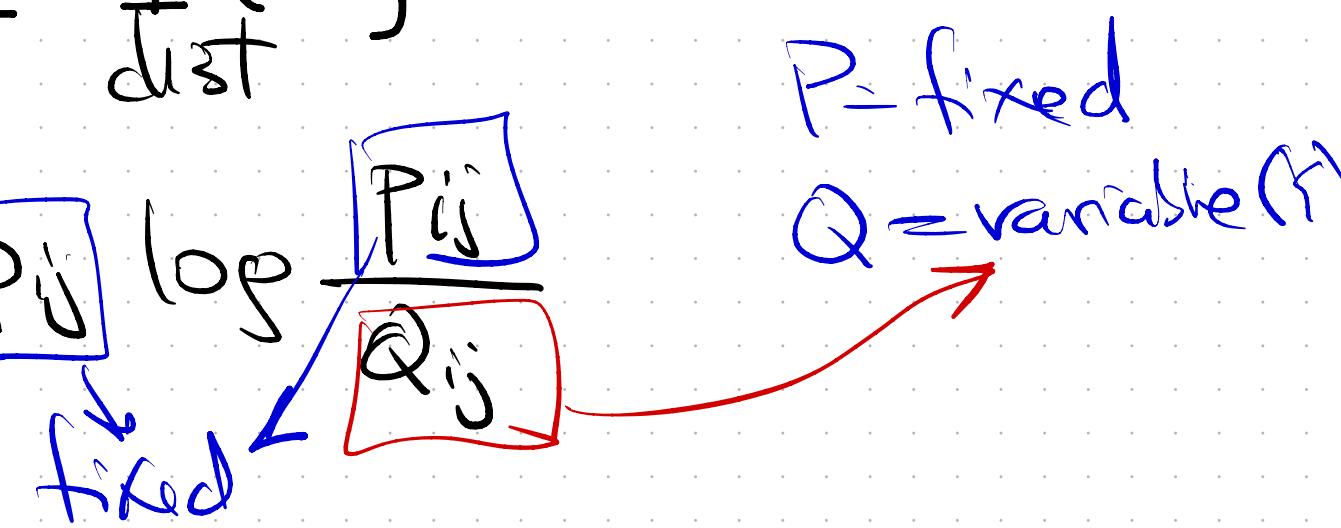
- init $Y = \text{random}$ (better cluster)

- iterate (Grad descent)

- compute $Q = t \cdot \text{dist}(Y \text{ positions})$

- update Y based Grad Descent Rule

$$\frac{\partial \text{OBJ}}{\partial Y_i} = 4 \sum_j (P_{ij} - Q_{ij}) (y_i - y_j)$$



$$\frac{1}{1 + \|y_i - y_j\|^2}$$

Q_{ij} numerator

easy to compute (fast using place lacan - kind of comp)
 $y_i^{\text{new}} = y_i + \gamma \frac{\partial \text{OBJ}}{\partial y_i} + \beta(y_i - y_i^{\text{prev}})$ on graph degrees