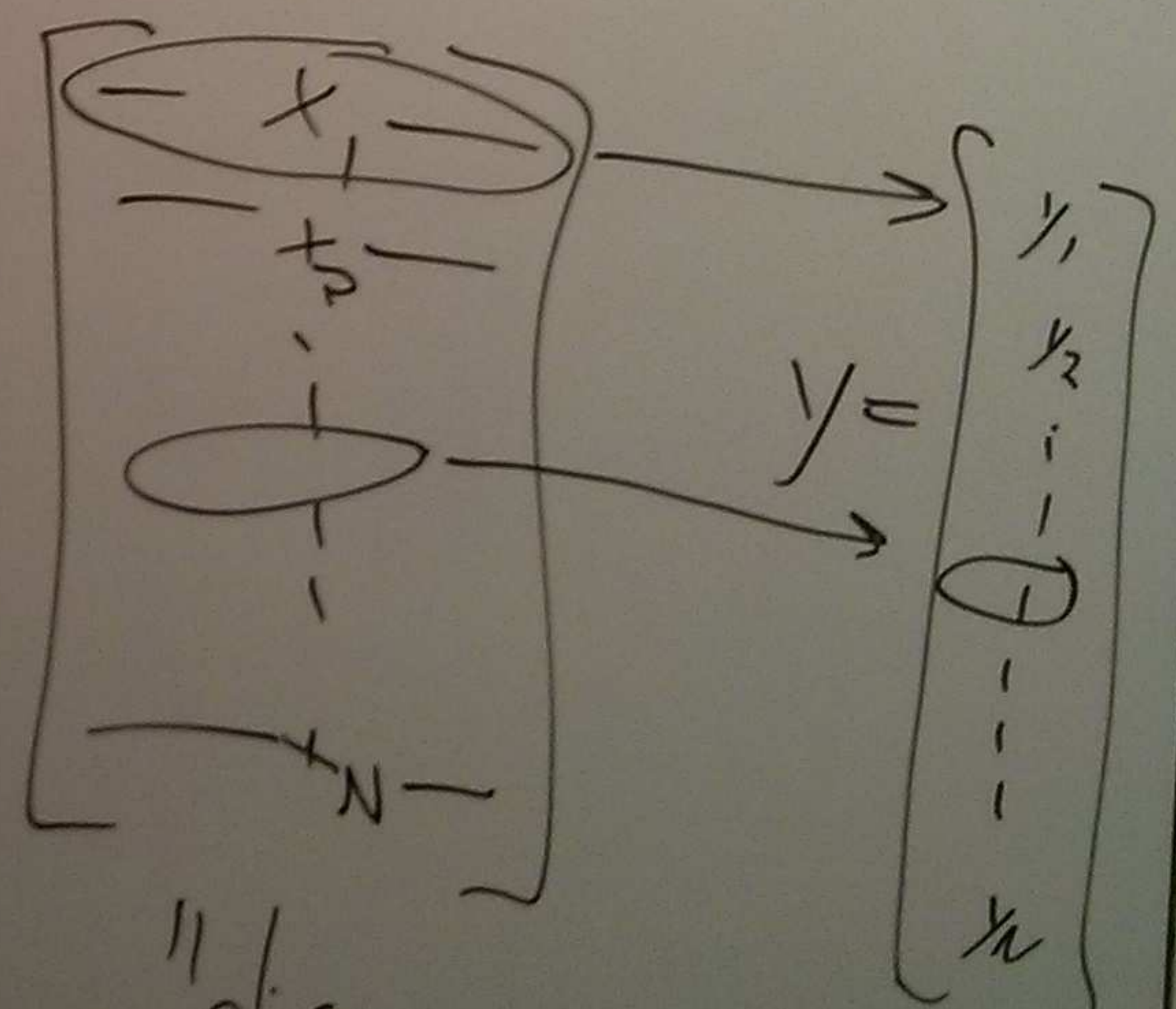
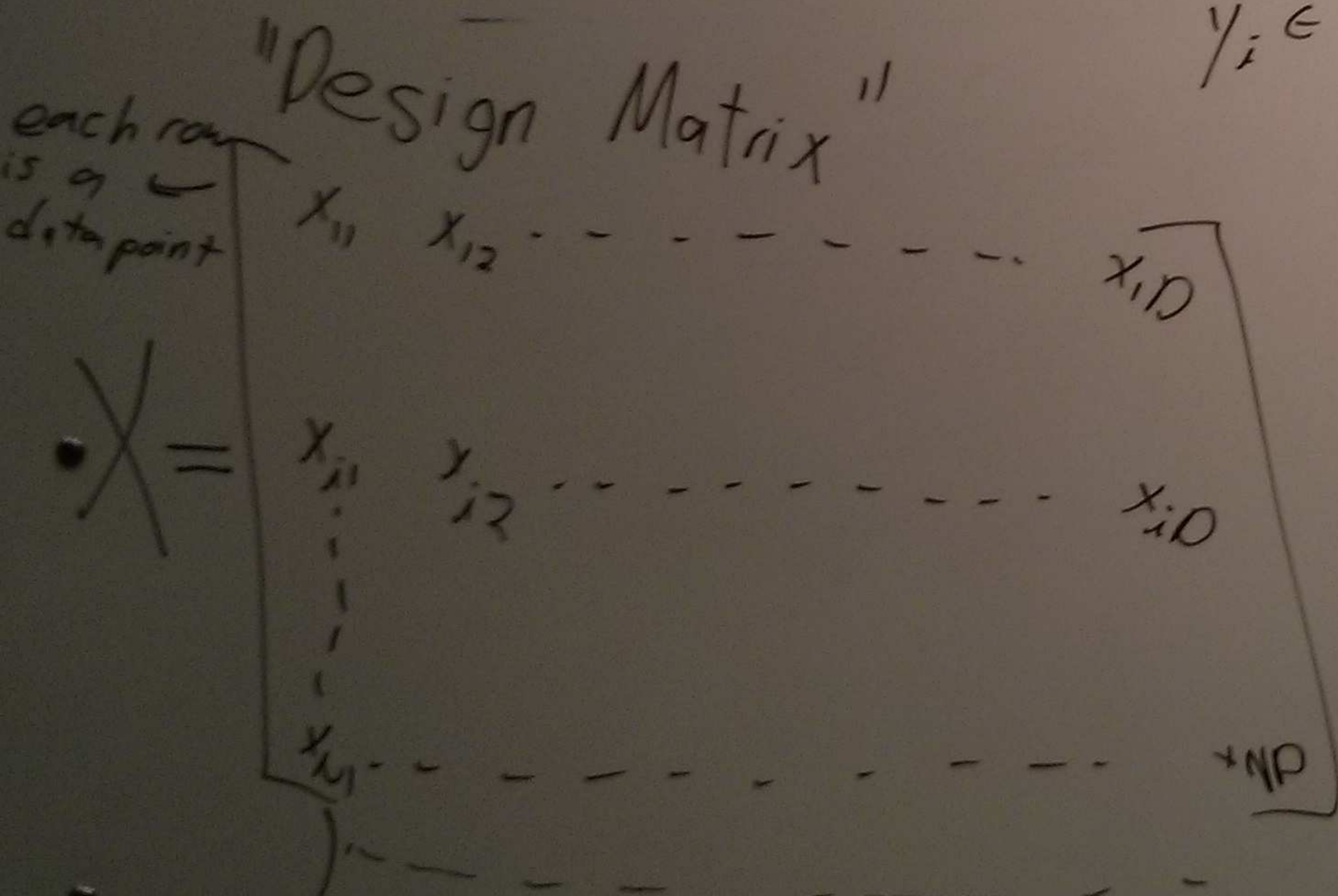


Supervised binary classification

Data Set $D = \left\{ (x_i, y_i) \right\}_{i=1}^N$

$x_i \in \mathbb{R}^D$ (.5, .6, 0, 1, -18, ...)

$y_i \in \{-1, 1\}$



Training

column is feature values across examples

Supervised learning is about finding f

$f(x_i) \mapsto y_i$

Testing phase

$$D_{\text{test}} = \sum_{i=1}^T x_i$$

$$X_{\text{test}} = \begin{bmatrix} \\ \\ \\ \\ \\ \\ \end{bmatrix}$$

$f(x_i)$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_T \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_T \end{bmatrix}$$

K-Nearest Neighbours

Training:

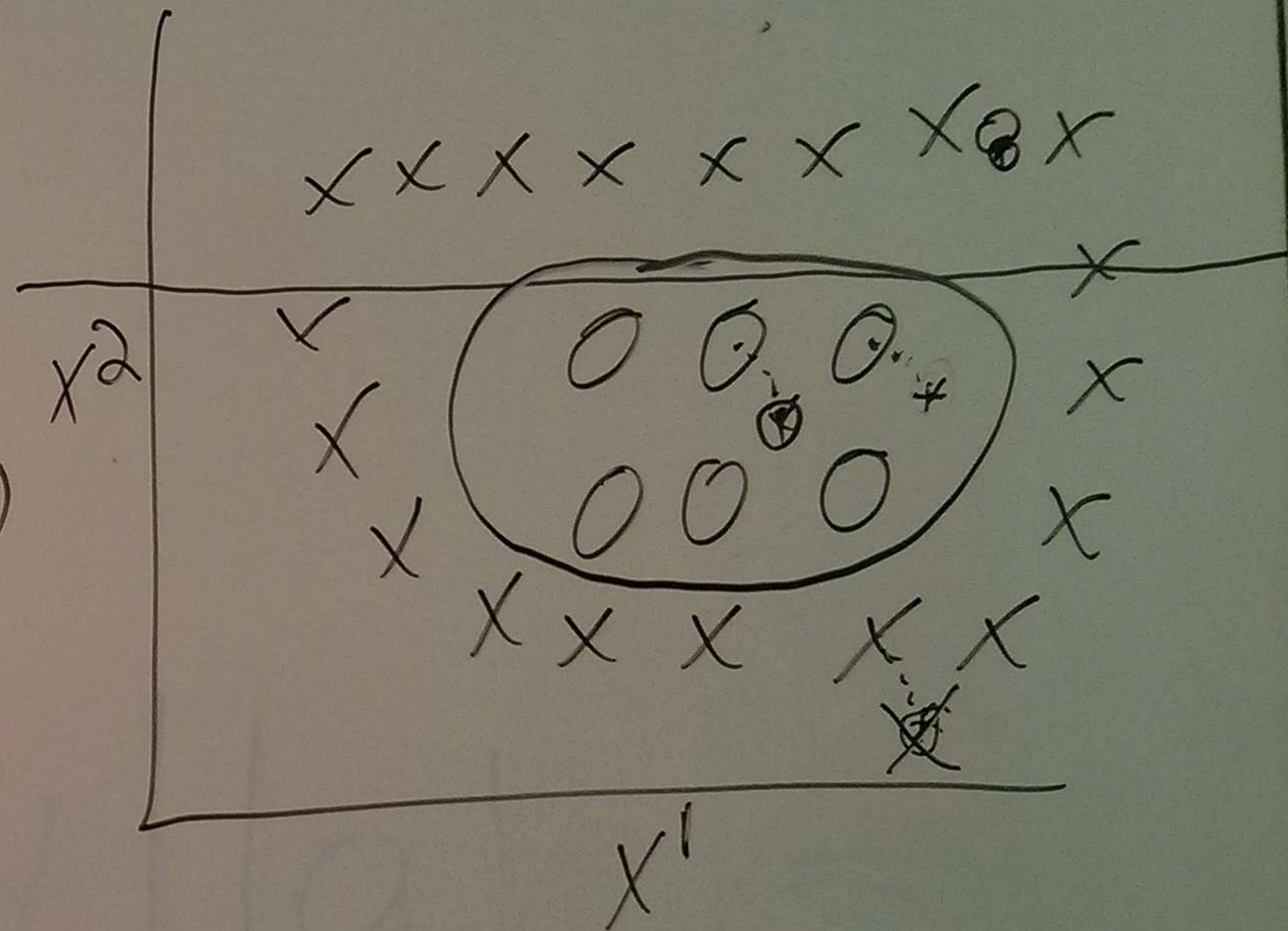
store X and y , choose K .

Testing:

for $x_i \in X_{\text{test}}$

find K 'closest' examples $\{x_1, x_2, \dots, x_K\}$ from training set.

look at corresponding $\{y_1, y_2, \dots, y_K\}$
choose majority



$O(ND)$

Handwritten scribbles at the top of the page.

Issues

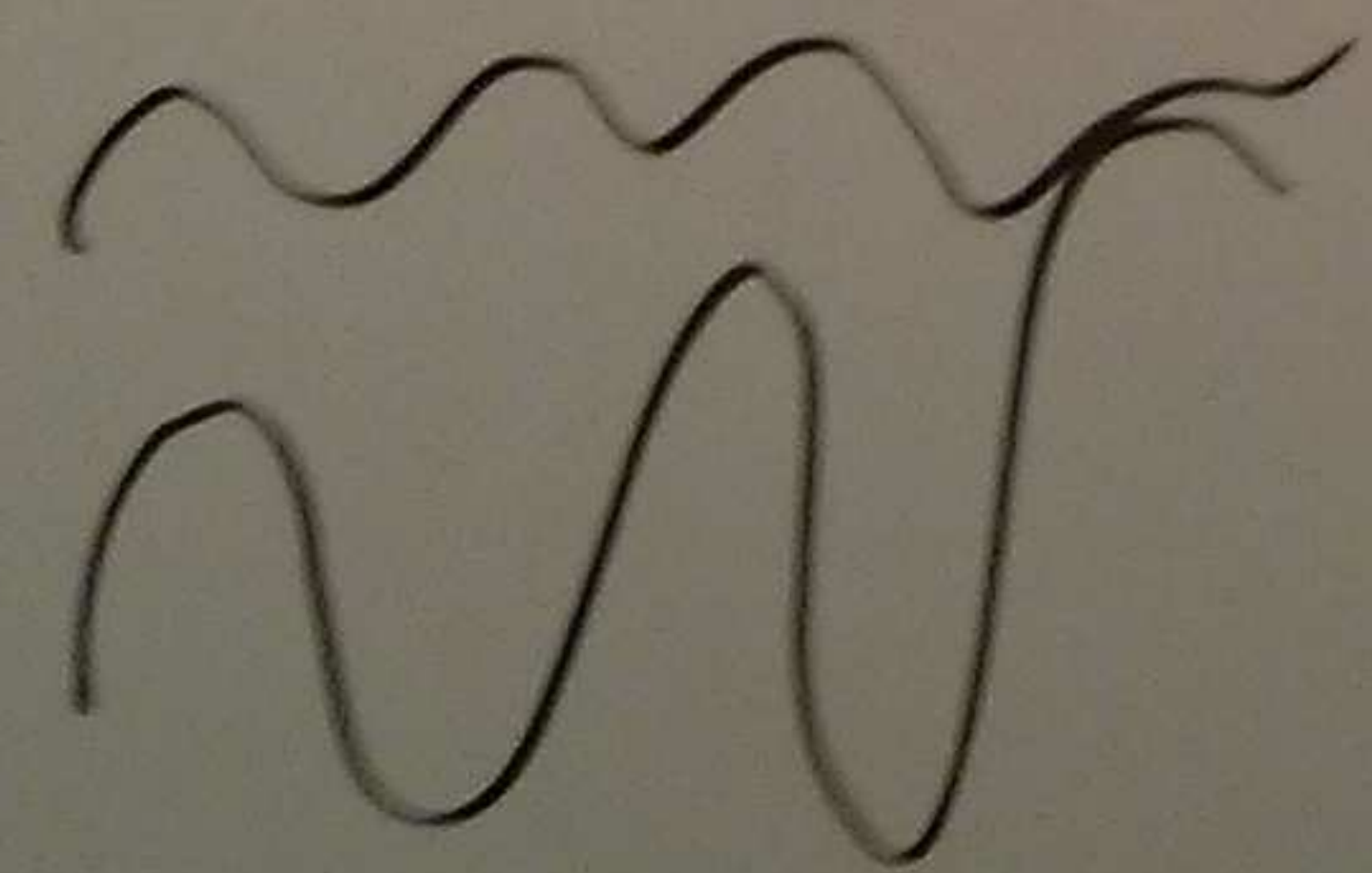
- Ties
- Weight neighbours 'closest'

Advantages

- only 1 choice
- simple, parallelizable
- training fast

xx	oo	xx	oo
oo	⊗	oo	xx
xx	oo	xx	oo
oo	xx	oo	xx

- fast updates
- highly flexible/non-linear



Disadvantages

- huge memory
- choose k
- class imbalance
- need to cover the space } 'curse of dimensionality'
- testing slow
- how to define 'distance'?
- Uniform distance measure?