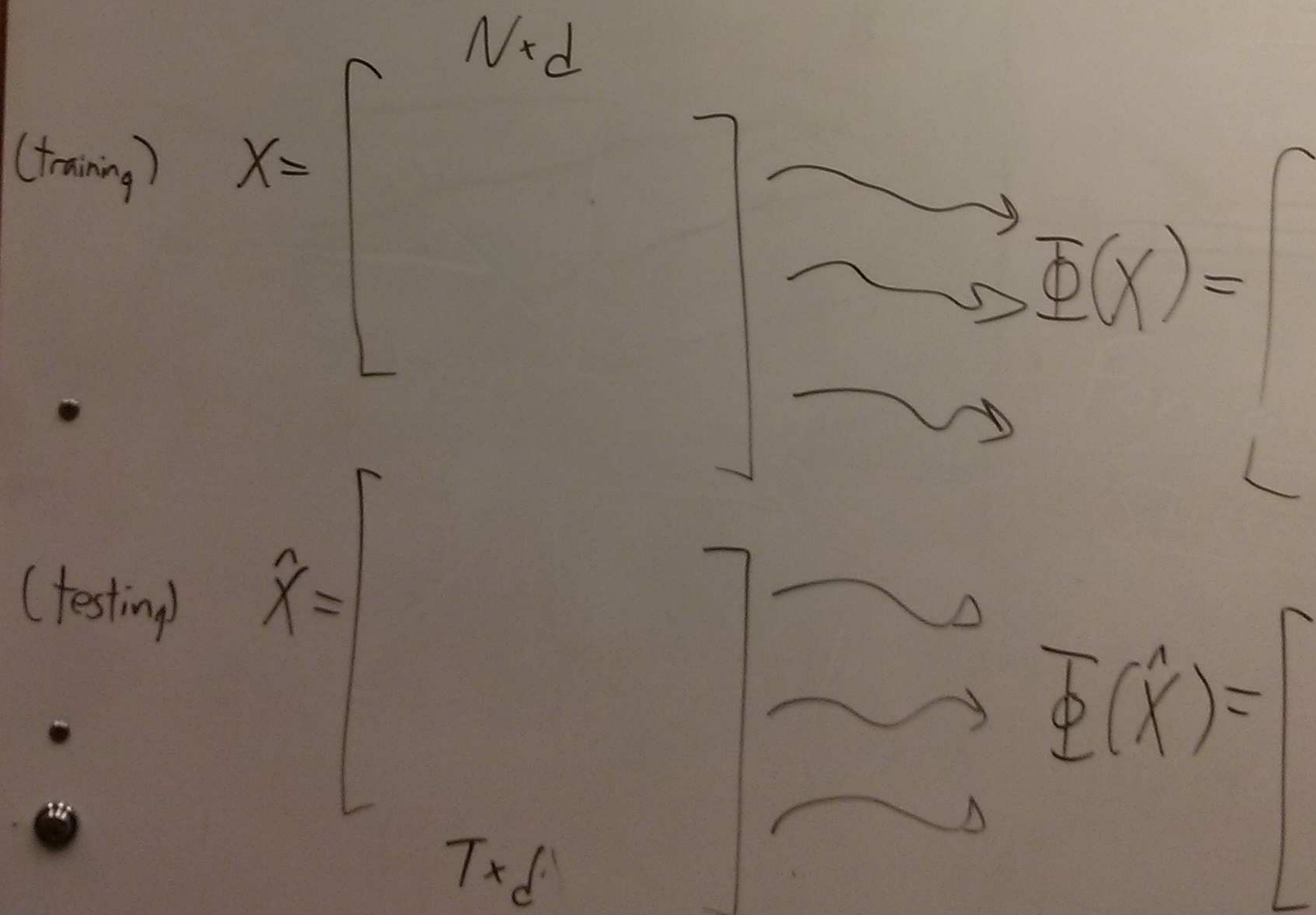


Valid kernels
 L_1 -Regularization
 Logistic Regression

Ass 2 due Wed
 Marked Ass 1 due now.

Review of kernel "trick"



$N \times d^{100}$

$T \times d^{100}$

led

now.

$N \times d^{100}$

$T \times d^{100}$

$$\hat{Y} = \overbrace{\Phi(\hat{X})}^{T \times d^{100}} \overbrace{\bar{W}_{MAP}}^{d^{100} \times 1}$$

$$= \underbrace{\Phi(\hat{X})}_{T \times d^{100}} \underbrace{(\Phi(X)^T \Phi(X) + \lambda I)^{-1}}_{d^{100} \times d^{100}} \underbrace{\Phi(X)^T Y}_{\substack{d^{100} \times N \\ N \times 1}}$$

$$= \underbrace{\Phi(\hat{X}) \Phi(X)^T}_{T \times N} \underbrace{(\Phi(X) \Phi(X)^T + \lambda I)^{-1}}_{N \times N} Y$$

$$\stackrel{(*)}{=} \underbrace{K(\hat{X}, X)}_{T \times N} \underbrace{(K(X, X) + \lambda I)^{-1}}_{N \times N} Y$$

If you can compute $k(x, x)$ and $k(\tilde{x}, x)$,
no need for $\Phi(x), \Phi(\tilde{x})$.

$$\Phi(x) \Phi(x)^T = \begin{bmatrix} \phi(\tilde{x}_1)^T \phi(\tilde{x}_1) & \phi(\tilde{x}_1)^T \phi(x_2) & \dots \\ \phi(\tilde{x}_2)^T \phi(\tilde{x}_1) & & \\ \vdots & & \\ \vdots & & \end{bmatrix}$$

$$k(x, x) = \begin{bmatrix} k(\tilde{x}_1, \tilde{x}_1) & k(\tilde{x}_1, \tilde{x}_2) & \dots \\ k(\tilde{x}_2, \tilde{x}_1) & & \\ \vdots & & \\ \vdots & & \end{bmatrix}$$

Valid kernels

Q: When does there exist feature map ϕ such that for kernel k we have $\phi(\tilde{x}_i)^T \phi(\tilde{x}_j) = k(\tilde{x}_i, \tilde{x}_j)$???

A: If $k(x_1, x_2) \succeq 0$, for all x_1, x_2 .
 (may be hard to show)

Valid kernels

L_1 - Regularization

Logistic Regression

Assume $k_{x,x'}$ is valid:

- $ck, c > 0$

- $k + k'$

- $k(\phi(x), \phi(x'))$

polynomial
- $\exp(k)$

- $\varphi(k)$, non-negative coefficient

- $f(x)k(x, x')f(x')$

Course Project Ideas:

- Take your application, design a kernel.

- Scale up to large 'N'.

- Learning the kernel (MKL, etc)

Ideas:

Search over features combinations

- NP-hard under most criteria (BIC)
- (CV)

Greedy methods

- Forward Selection: start with empty, add best.
(sub-modular: best approximation)
- Backward select: start w all vars, remove worst.
- Stagewise

L₁ - Regular

$\arg \min_w \frac{1}{2} \| \dots \|$

- regularization

- encourages var

$\hat{X}_w = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$

combinations

criteria (BIC)
(CV)

start with empty, add best.

(best approximation)

with all vars, remove worst.

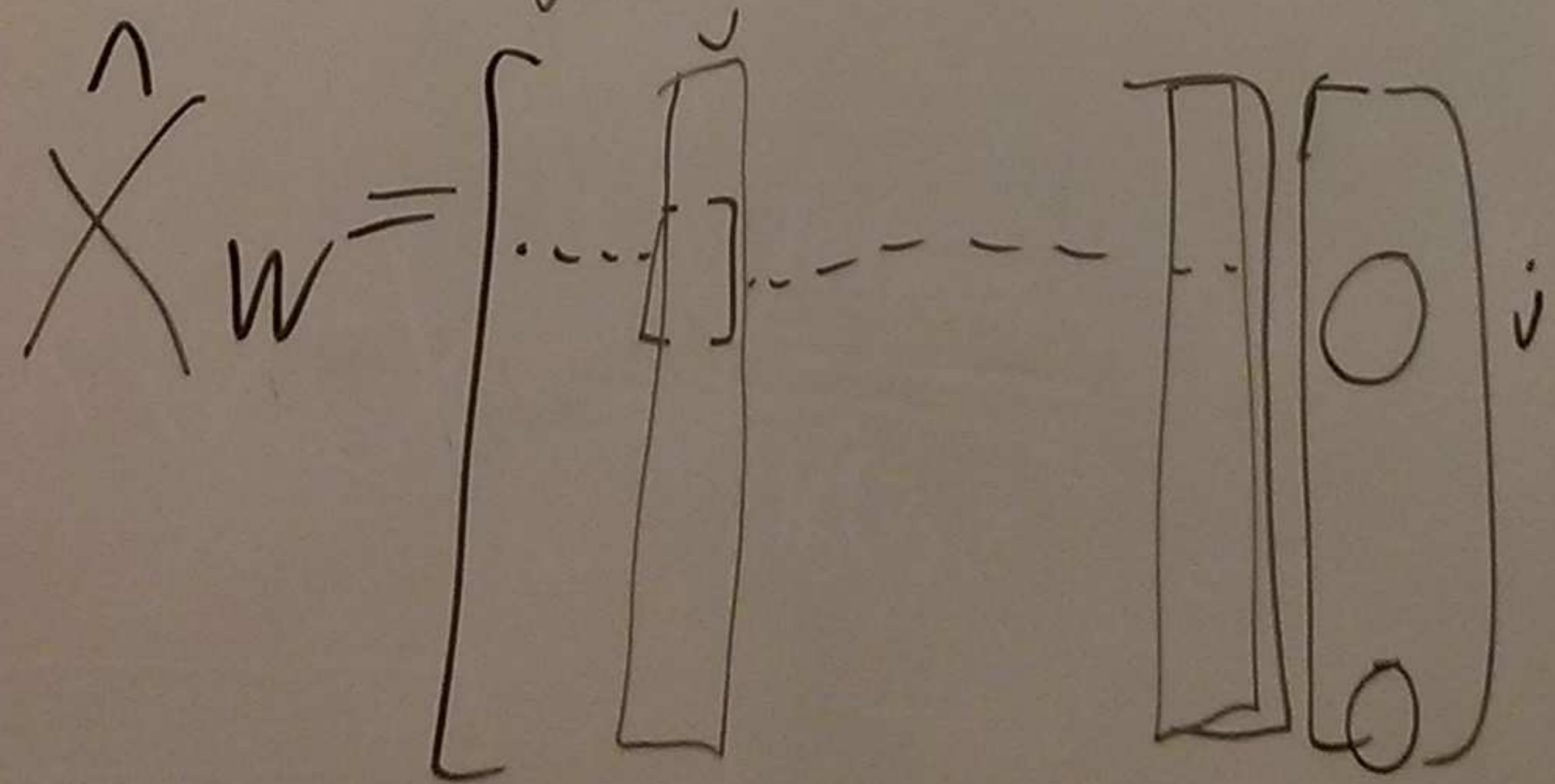
L₁ - Regularization

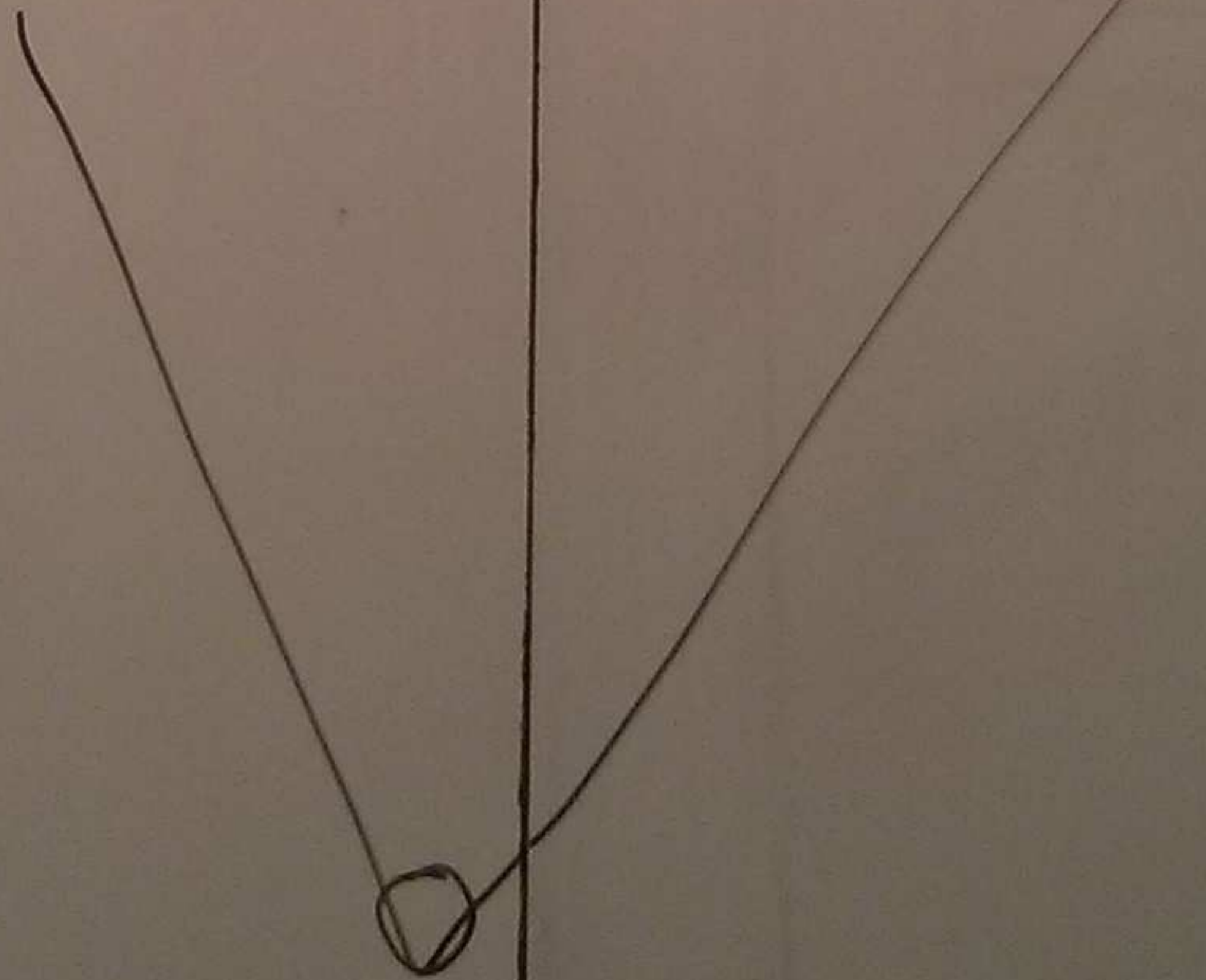
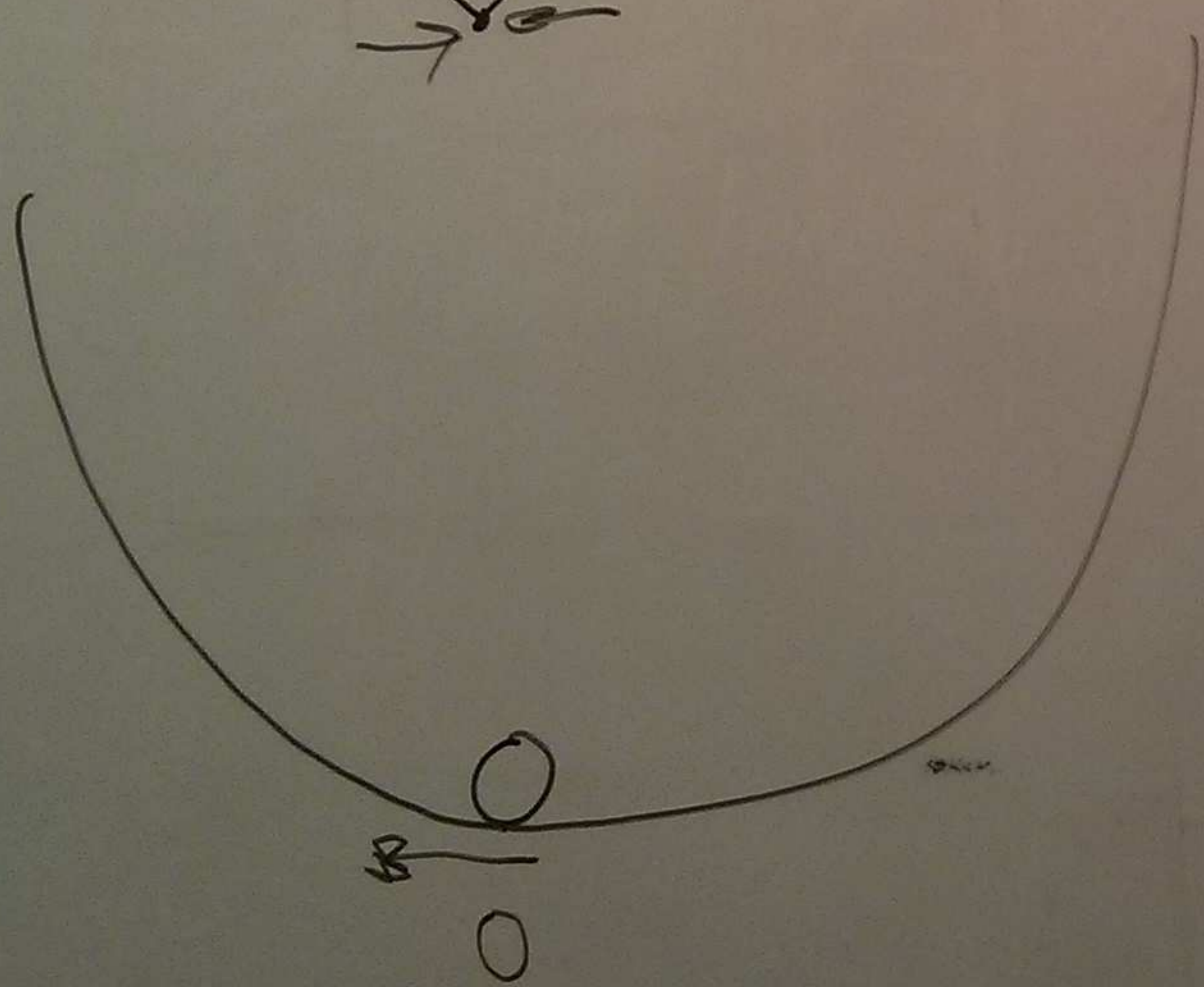
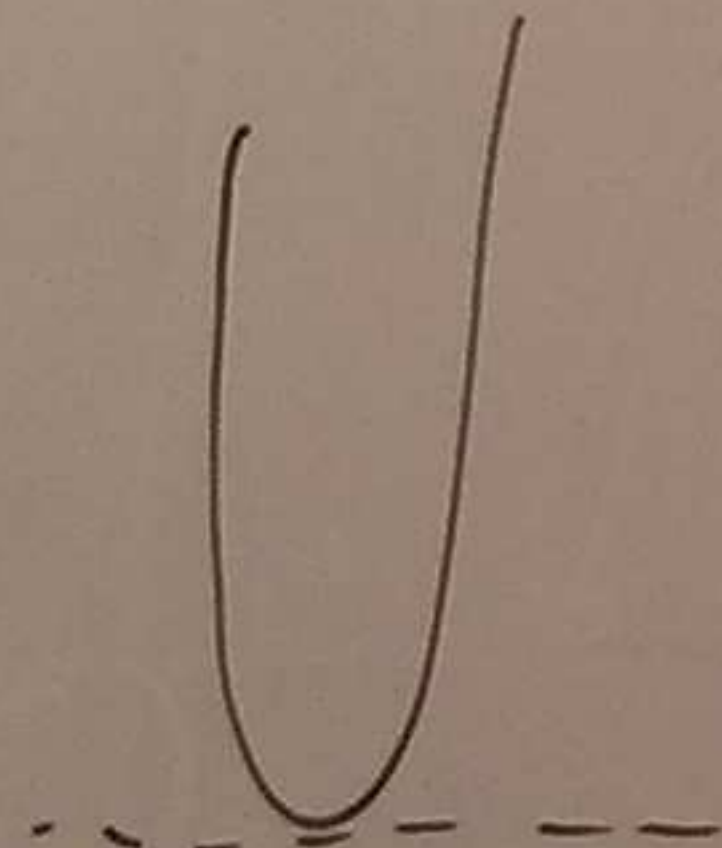
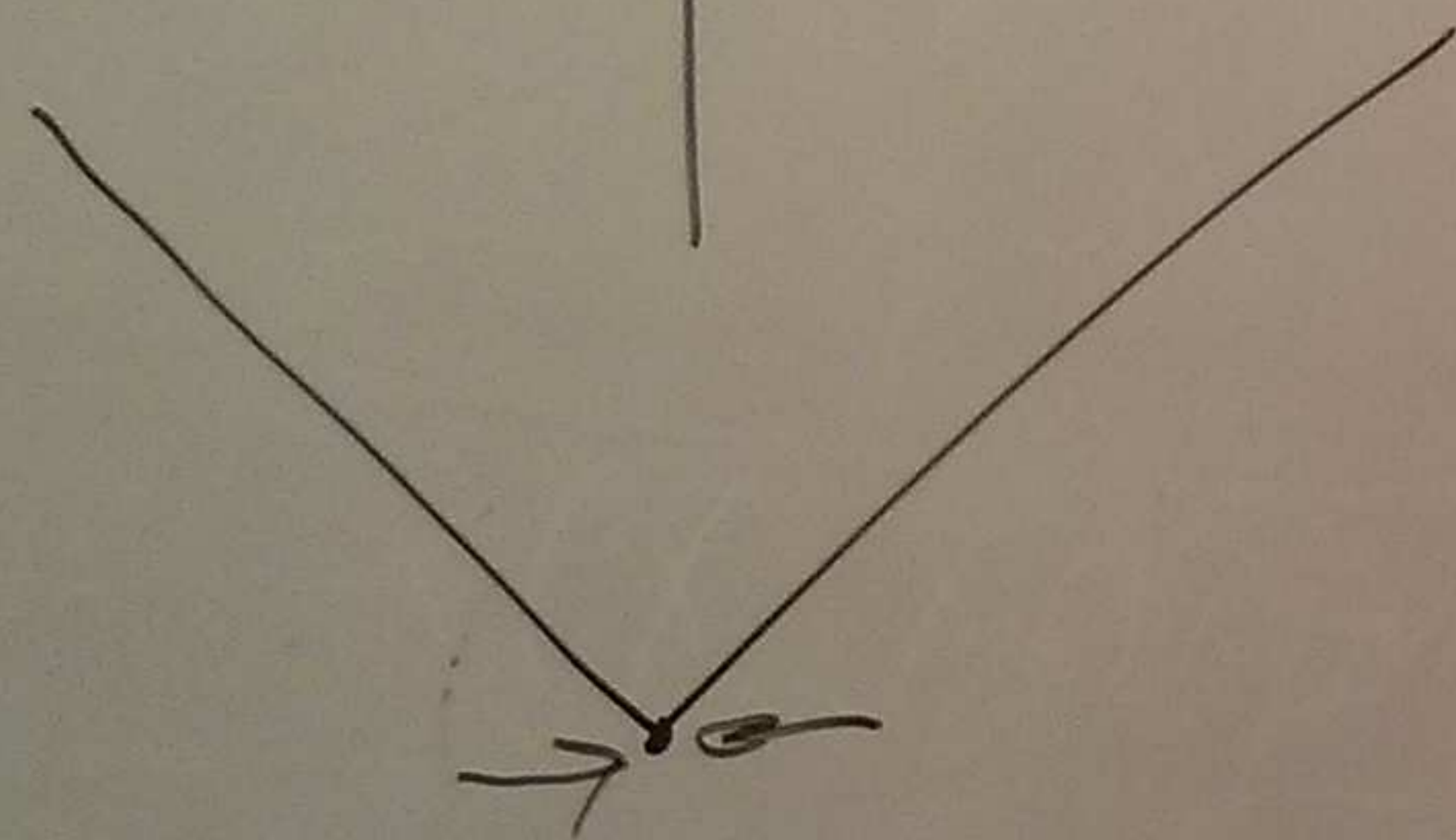
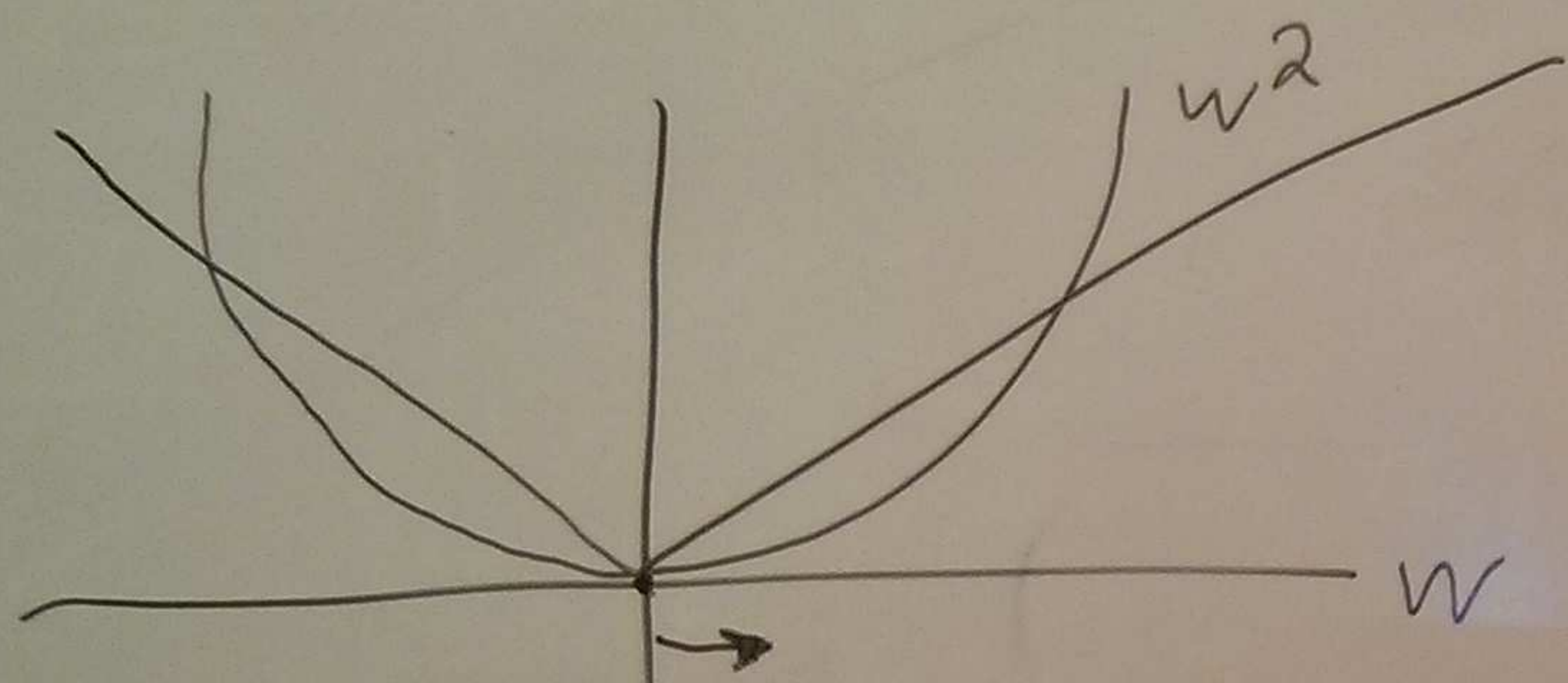
$$\arg \min_{\bar{w}} \frac{1}{2} \|X\bar{w} - Y\|^2 + \lambda \|\bar{w}\|_1$$

$$\|\bar{w}\|_1 = \sum_{i=1}^d |w_i|$$

— regularization: protect against over-fitting
(not unique)

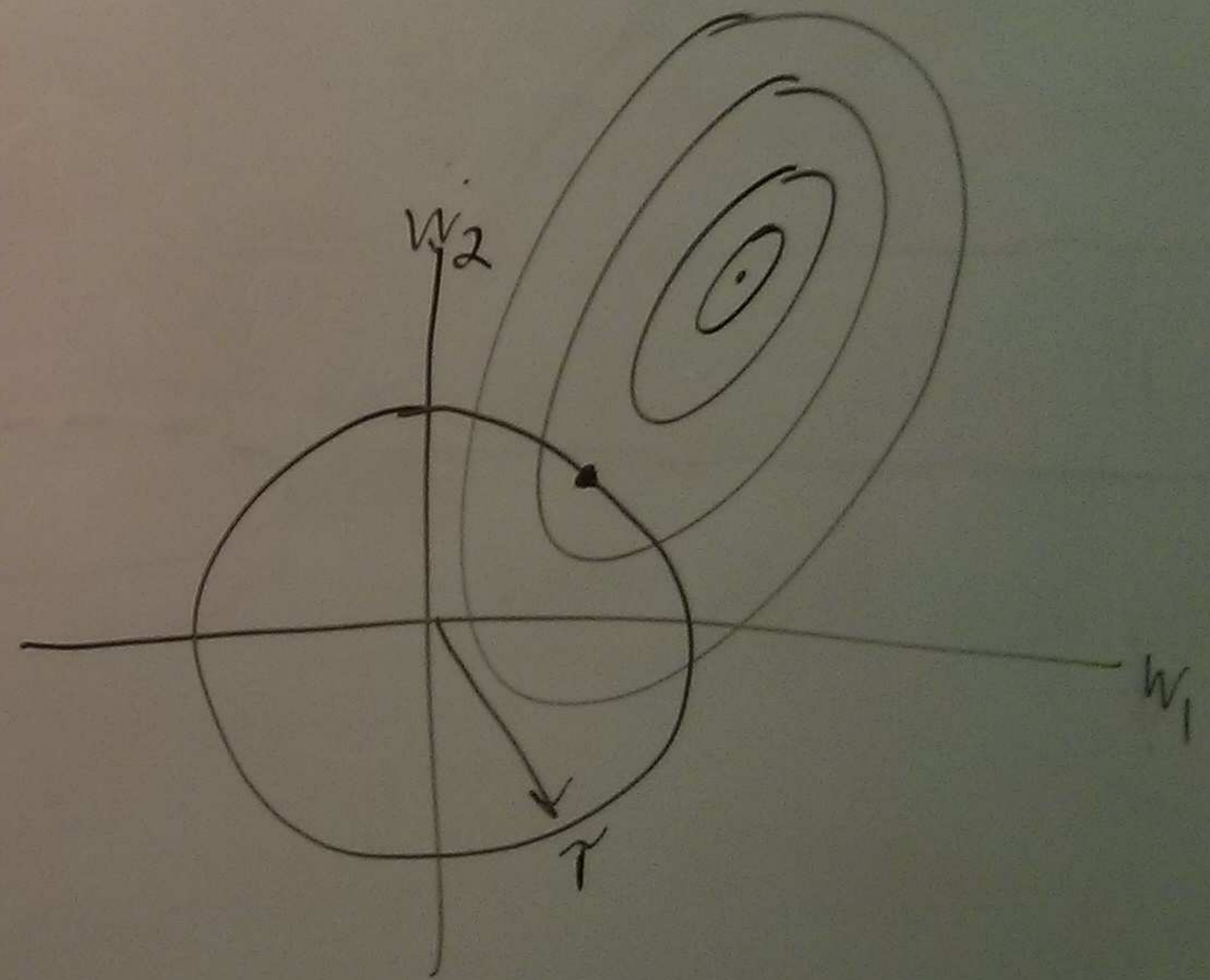
— encourages variable to be exactly 0





$$\operatorname{argmin}_{\bar{w}} \frac{1}{2} \|X\bar{w} - Y\|^2$$

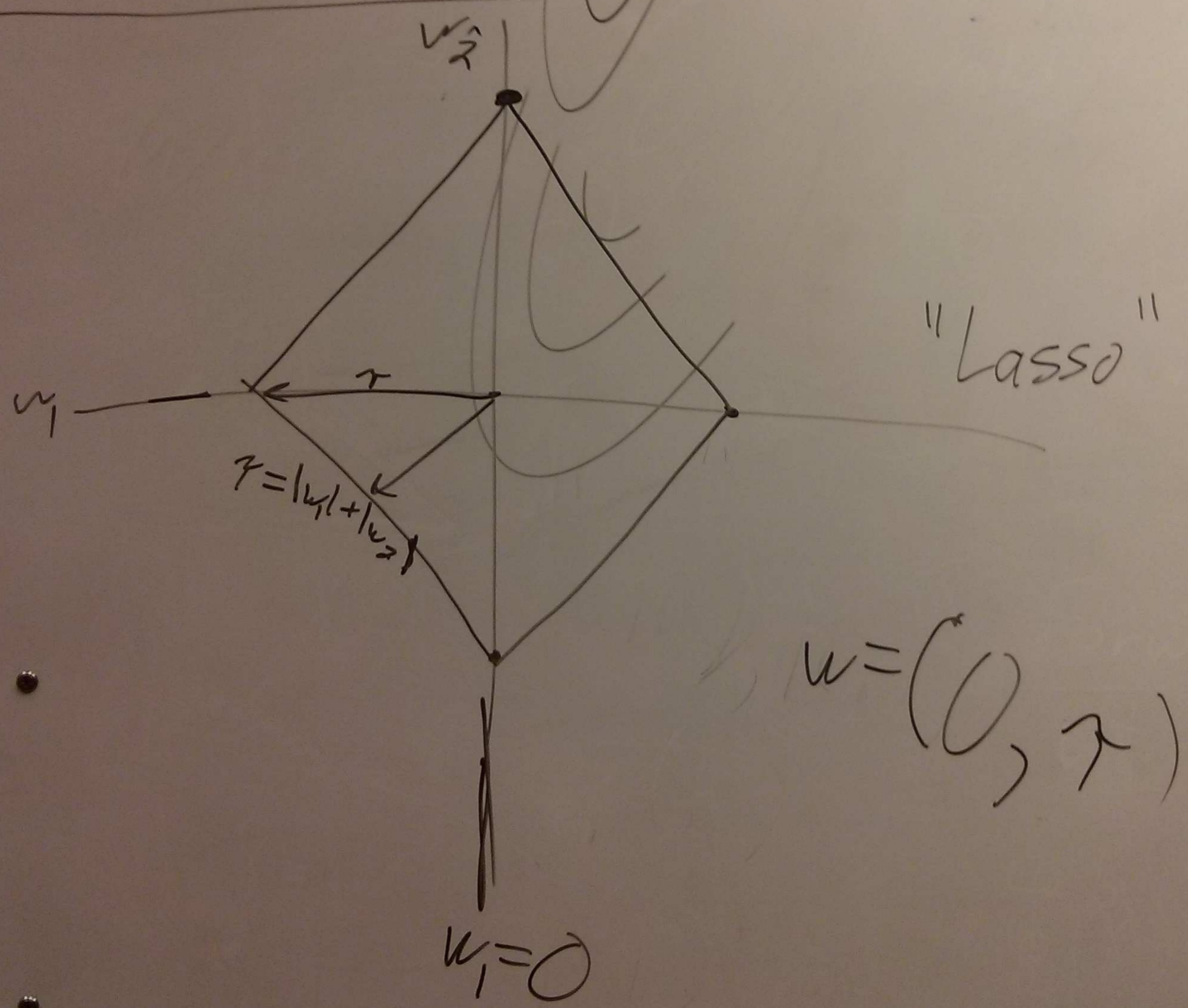
$$\text{subject to } \|\bar{w}\| \leq \gamma$$



Valid kernels

L_1 - Regularization

Logistic Regression

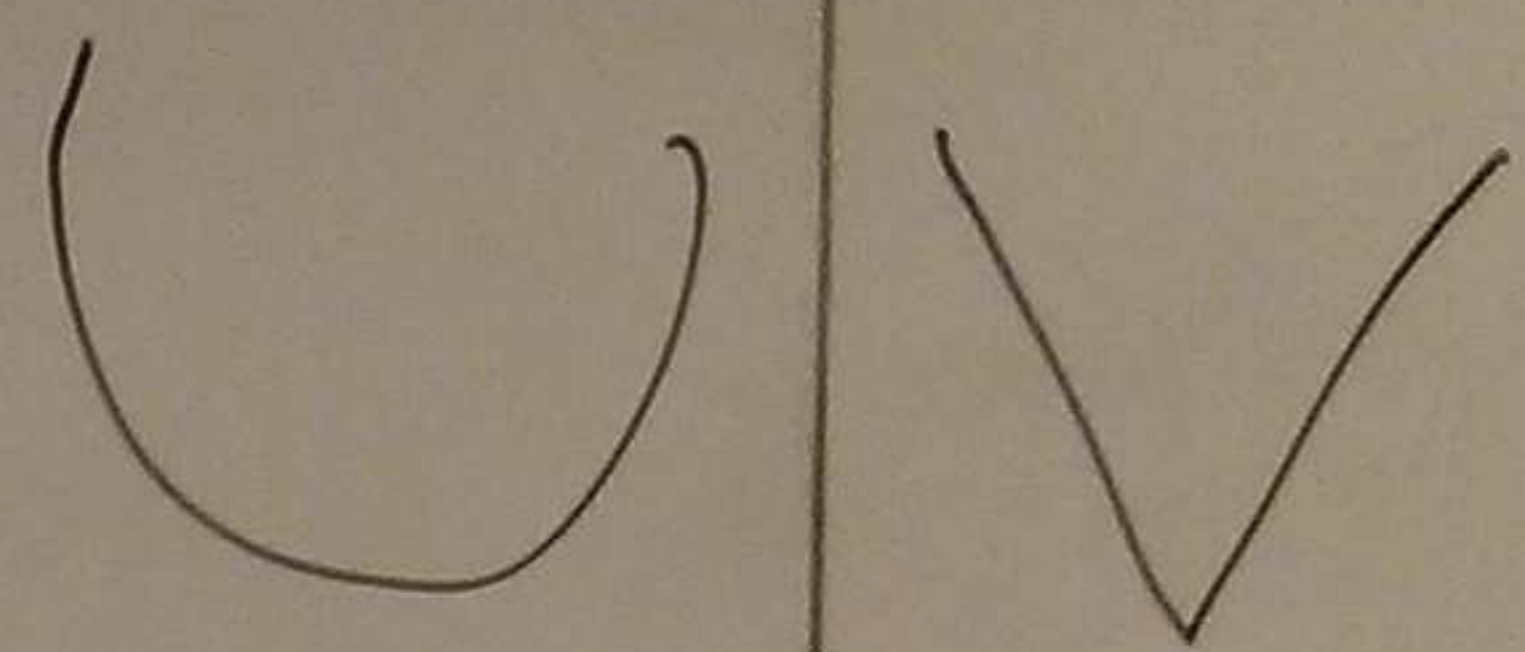


(smooth) + (non-smooth)
but separable

$$f(x) = \sum_{i=1}^d f_i(x_i)$$

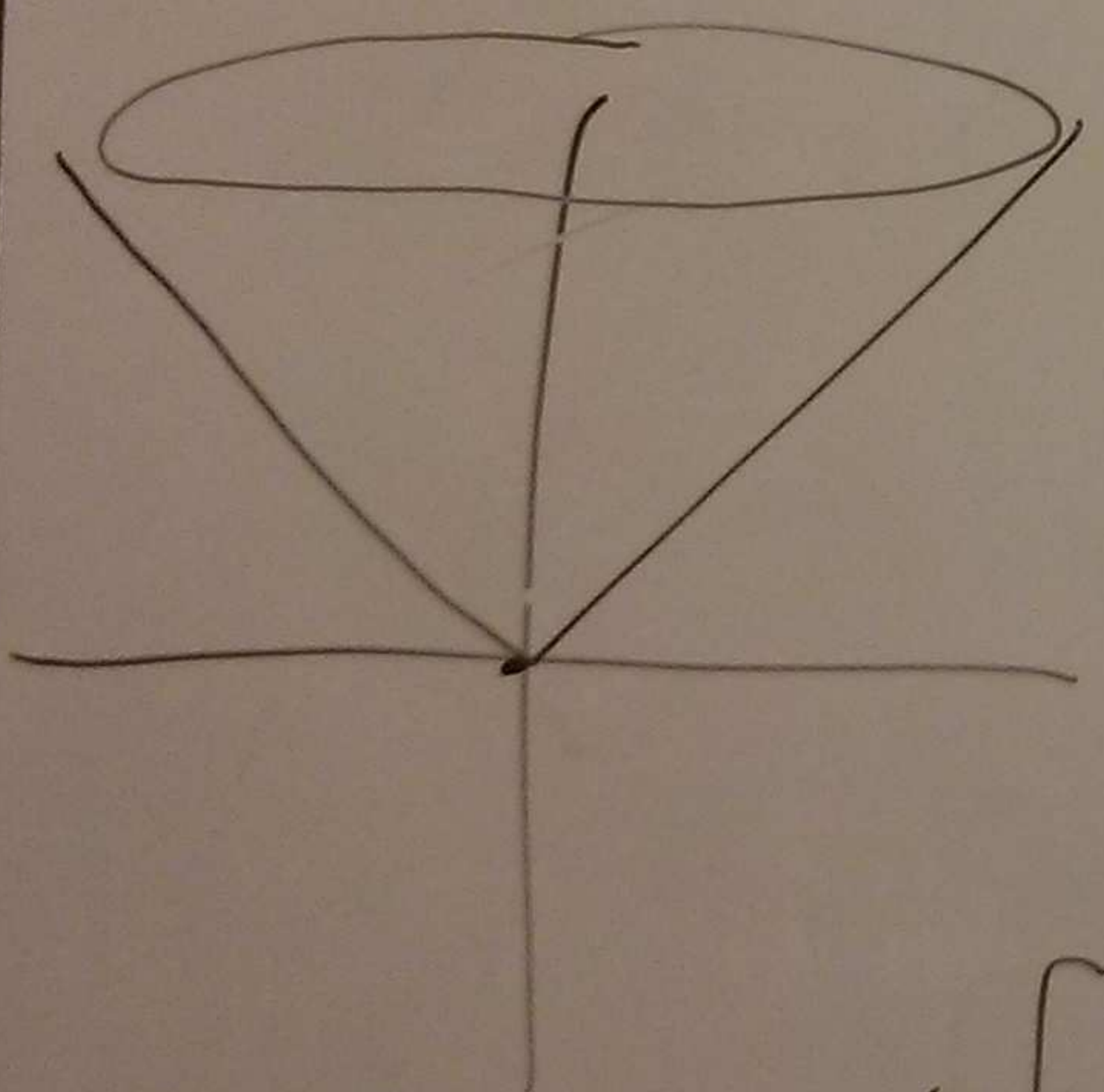
Elastic net: $\frac{\lambda_2}{2} \|w\|^2 + \lambda_1 \|w\|_1$

unique Sparse

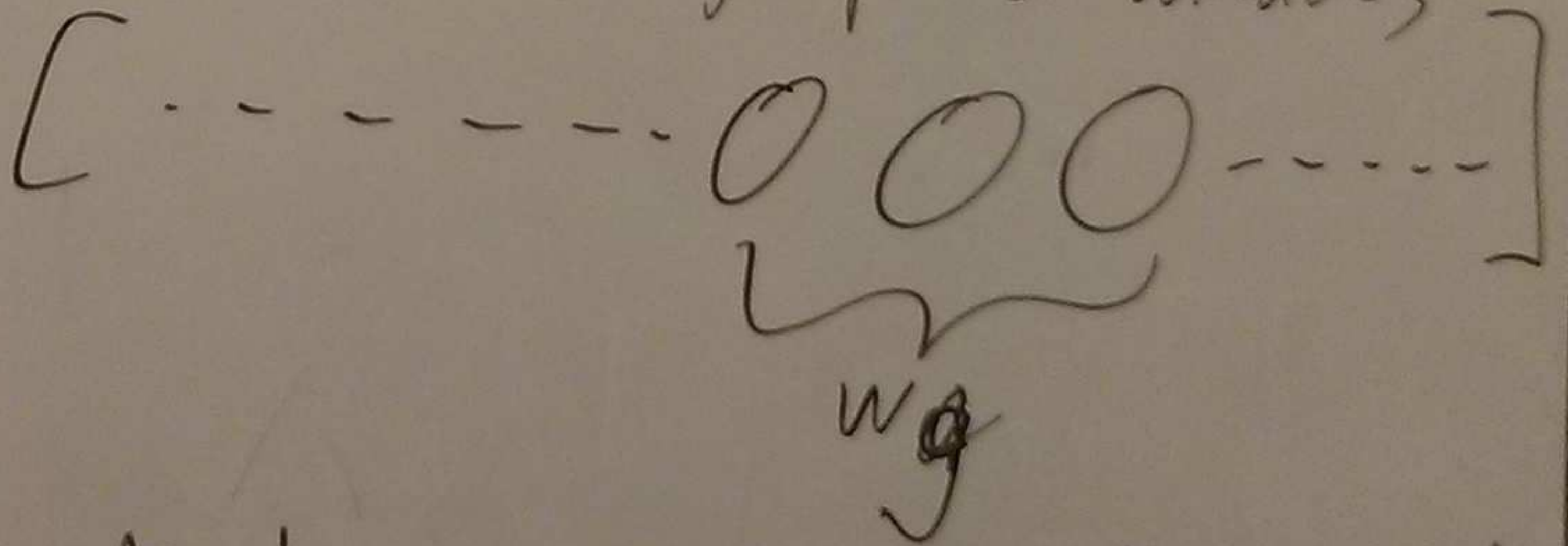


Group L1-Regularization: $\lambda \sum_g \|w_g\|$

w_g : sub-vector



Selects "groups" of variables

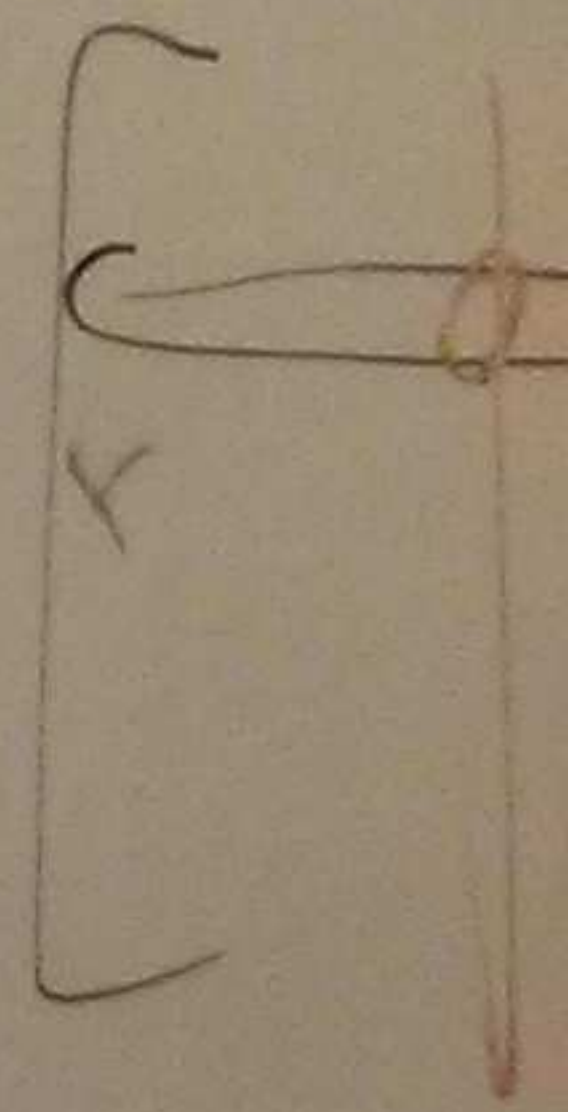


"multiple regression"

$X = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}$ $N \times d$

$Y = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}$ $N \times c$

XW

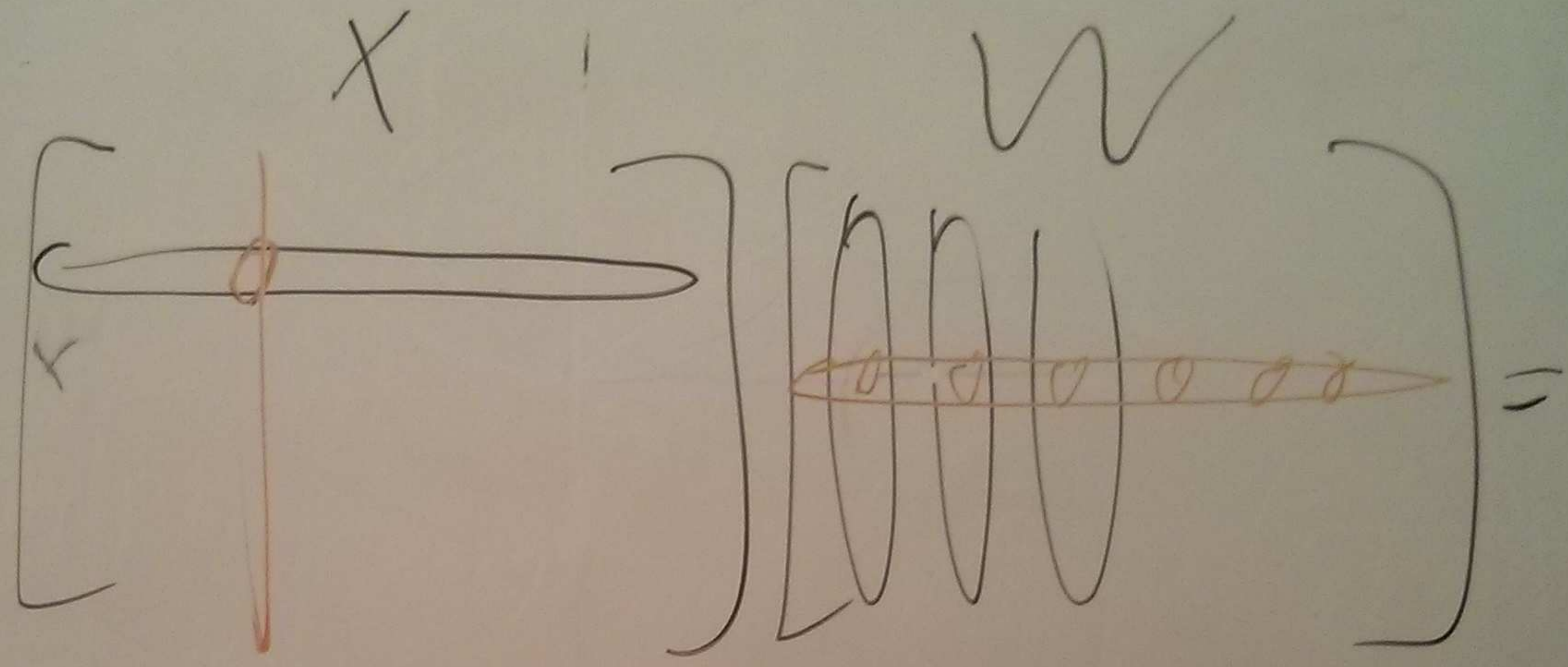


Nuclea

-vec for

5]

psion



x_1
 x_1^2 $x_1 x_2$

Nuclear-norm regularization:

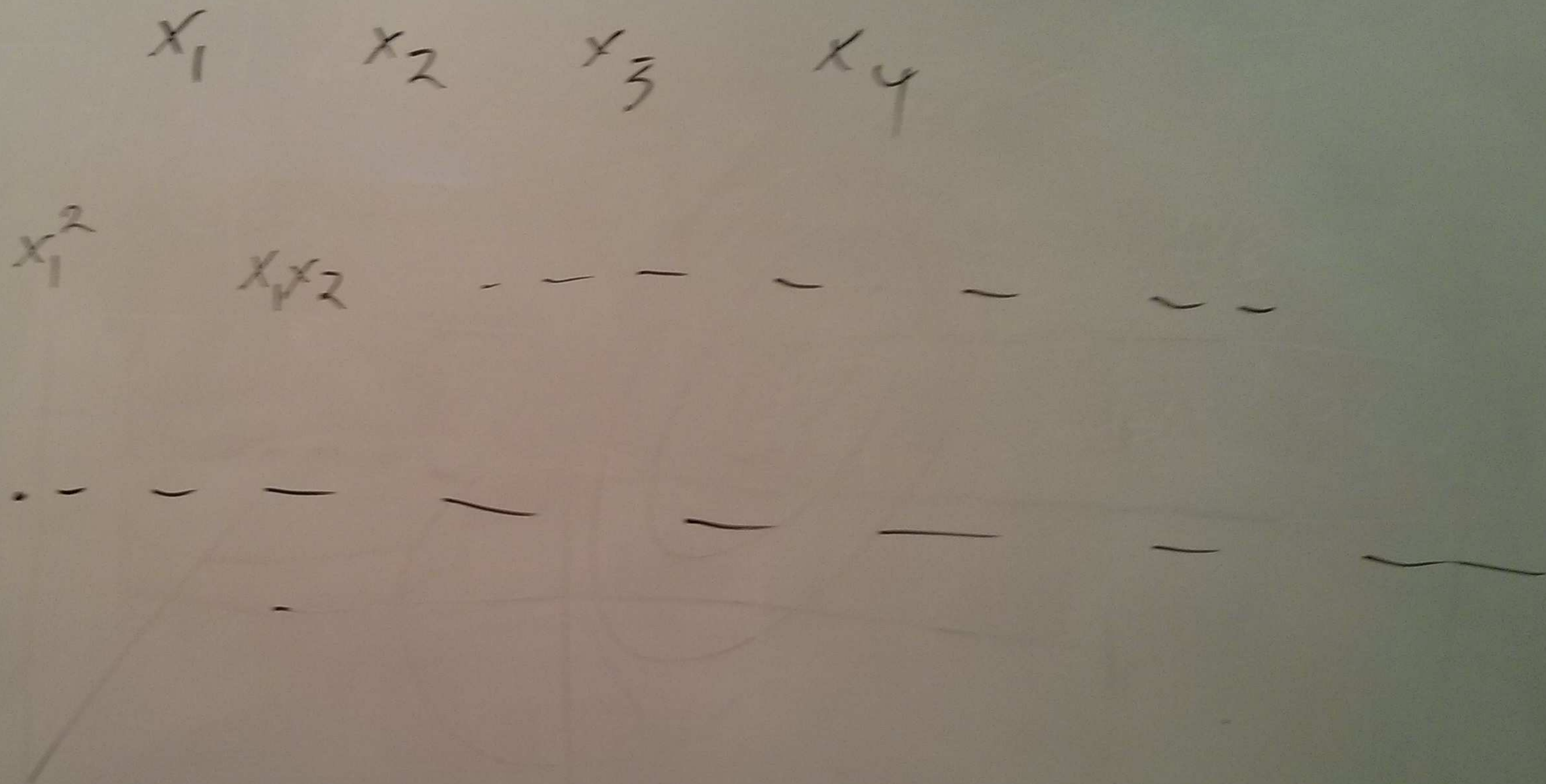
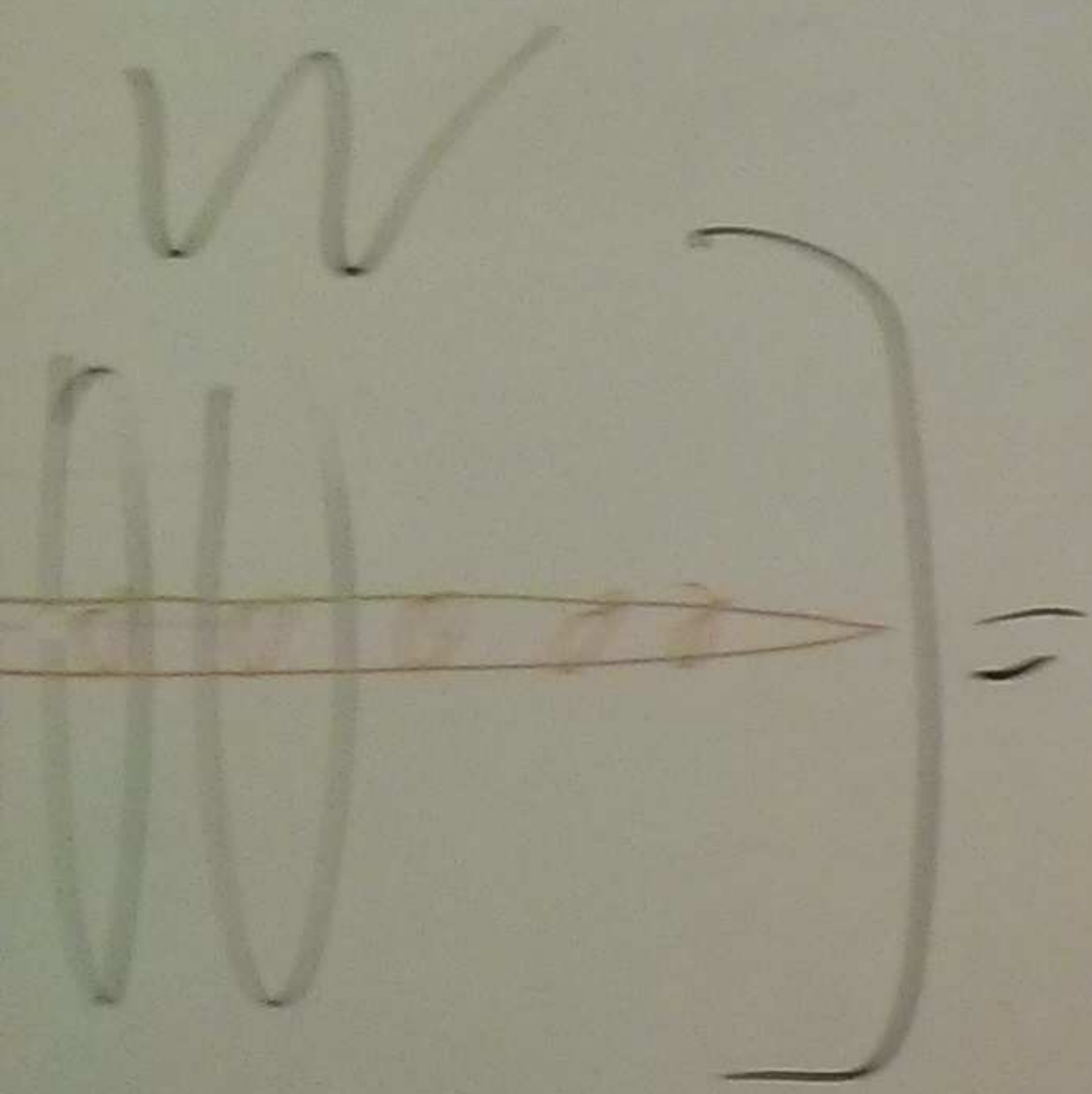
- sum of singular values
- \Rightarrow low-rank matrices

Structured Sparsity: " $w_2 \neq 0$ only
 - if $w_1 \neq 0$ "

(good project)

\Rightarrow publication

$x_1 x_2 x_5$



ization
of singular values
low-rank matrices

$\neq 0$ only
if $w_1 \neq 0$

$$x_1 x_2 x_5 x_6$$

Supervised Learning

Generative $p(y_i, x_i)$

Naive Bayes
GDA

Discriminative $p(y_i | x_i)$

Least squares

[Logistic Regression]

Discriminant function

KNN

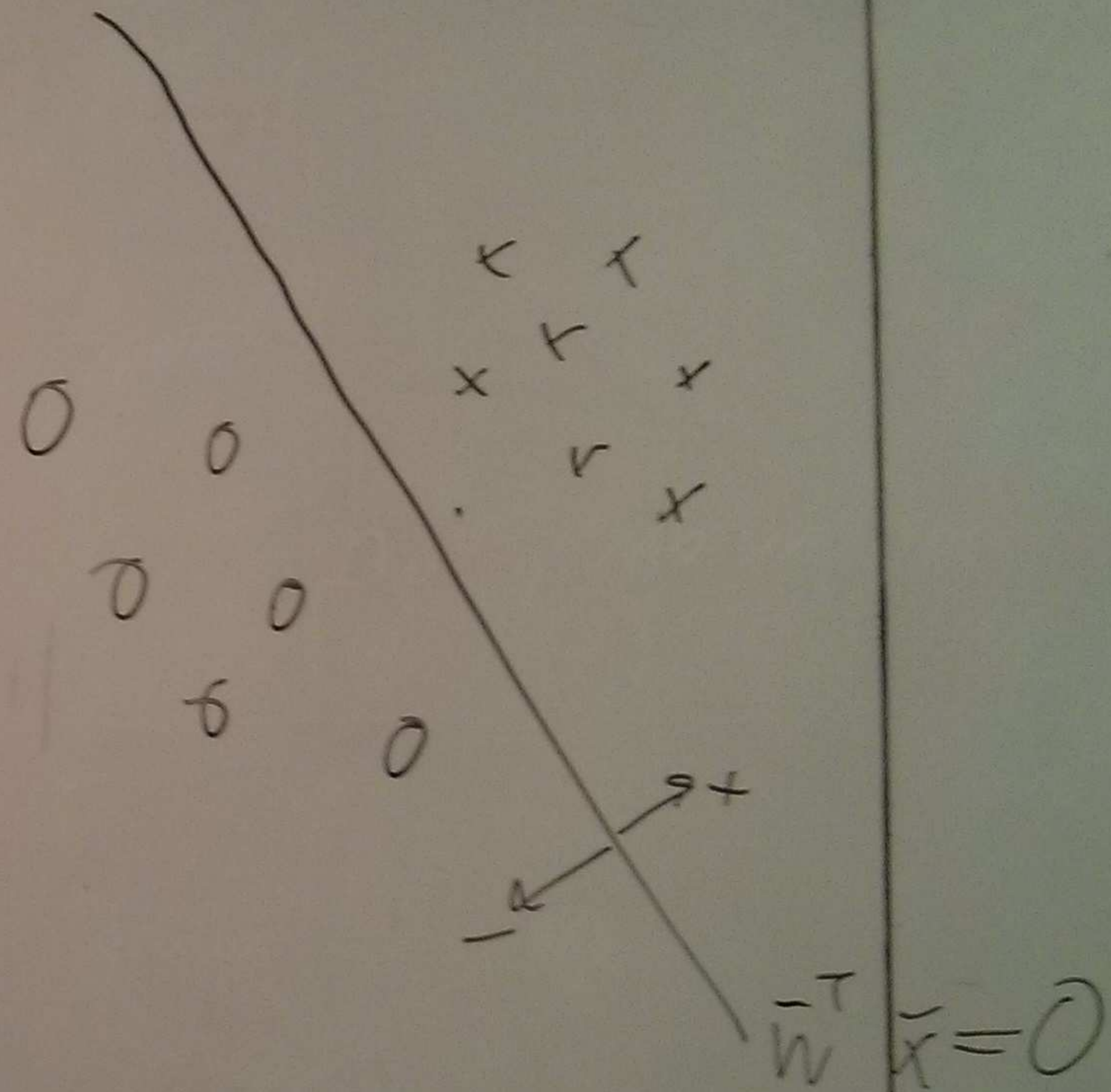
[Support vector machine]

features: $\bar{x}_i \in \mathbb{R}^d$

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

Linear classifiers:

$$\hat{y}_i = \text{sign}(w^T \bar{x}_i)$$

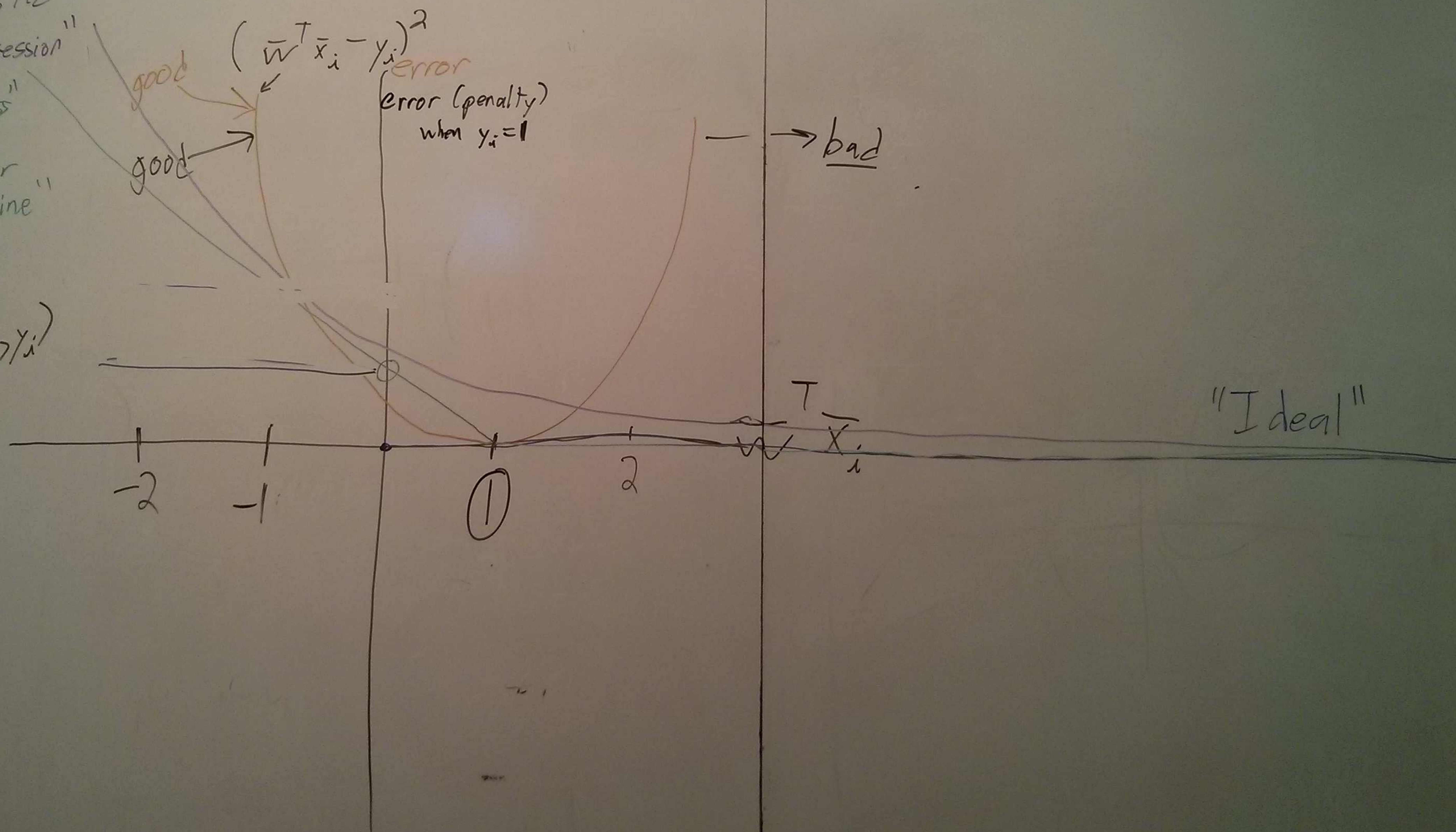


"logistic regression"

"hinge loss"

"support vector machine"

$$\sum_{i=1}^N f_i(\bar{w}^T \bar{x}_i / y_i)$$



$$(\bar{w}^T \bar{x}_i - y_i)^2$$

error (penalty) when $y_i = 1$

bad

"Ideal"

kernels

regularization
Regression

log - Odds ratio

$$\log \frac{p(y_i = 1 | \bar{x}_i, \bar{w})}{p(y_i = -1 | \bar{x}_i, \bar{w})} = \bar{w}^T \bar{x}_i$$

$$\frac{p(y_i = 1 | \bar{x}_i, \bar{w})}{p(y_i = -1 | \bar{x}_i, \bar{w})} = \exp(\bar{w}^T \bar{x}_i)$$

$$1 - p(y_i = 1 | \bar{x}_i, \bar{w})$$

$$p_1 = (1 - p_1) \exp(\bar{w}^T \bar{x}_i)$$

$$p_1 (1 + \exp(\bar{w}^T \bar{x}_i)) = \exp(\bar{w}^T \bar{x}_i)$$

$$p_1 = \frac{\exp(\bar{w}^T \bar{x}_i)}{1 + \exp(\bar{w}^T \bar{x}_i)} = \frac{1}{1 + \exp(-\bar{w}^T \bar{x}_i)}$$

"sigmoid function"

$$p_2 = \frac{1}{1 + \exp(\bar{w}^T \bar{x}_i)}$$