

Plan

Recap

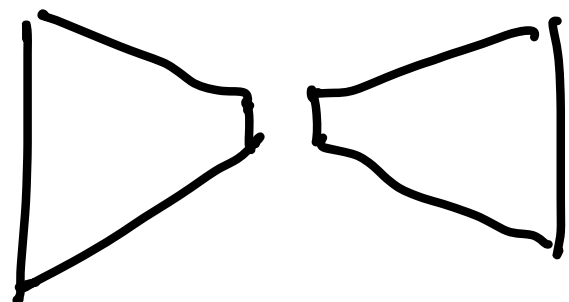
Logistics

Regularization

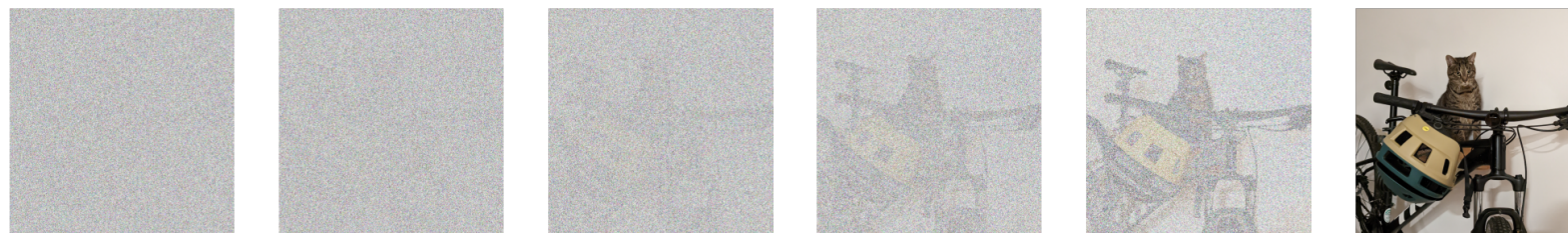
Implicit Regularization

Parameterization

CRF



Diffusion



x_t

x_{t-1}

$$x_t \in \mathbb{R}^{m \times m}$$

$$\epsilon_t \in \mathbb{R}^{m \times m}$$

$$x_t = x_{t-1} + \epsilon_t$$

$$\epsilon_t \approx f_{\theta}(x_t, t)$$

$$x_{t-1} = x_t - \epsilon_t$$

Stable Diffusion

↳ do diffusion in latent space

↳ text conditioning

Logistics

↳ Last day of new content

- last form tonight
- CRF

↳ Homework tomorrow

↳ Tomorrow

- no lecture/demo
- I'll be around
- games at 4

↳ presentation

Check Canvas Assignment

Supervised

(x, y) , n samples

$$f_w(x) \approx y$$

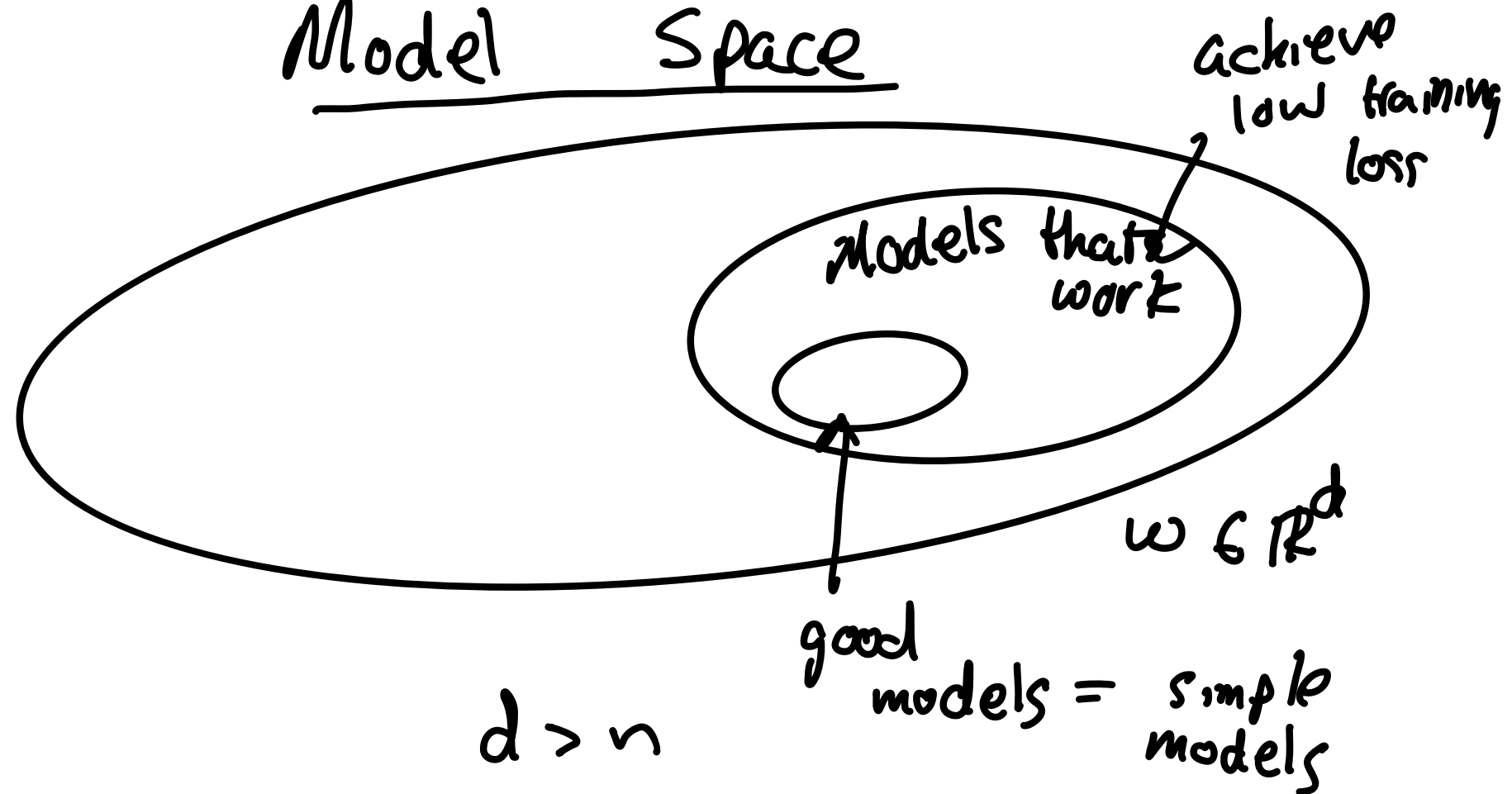
$$\mathcal{L}(w) = \|f_w(x) - y\|_2^2$$

training
 (x, y)

test
 (x, y)

Conceivably f_w could
memorize the training data

Model Space



Frumpf

Apple
Mouse

Computer
Board
phone

Clon

Pear
Kiwi
Tea

Pen
Ball
Horse

Strategies

- last letter
- first letter
- length
- vowels

What is a simple model?

↳ few parameters

What is a simple model
in the space of fixed
number of parameters?

↳ $\|w\|_2^2 \sim \text{complexity}$

Large $\|w\|_2^2 \rightarrow$

big changes when
we perturb

Regularization

$$\mathcal{L}(w) + \lambda \|w\|_2^2$$

↑
how
well we
solve
task

↑
complexity
of model

↳ "lambda" tradeoff

we do this implicitly

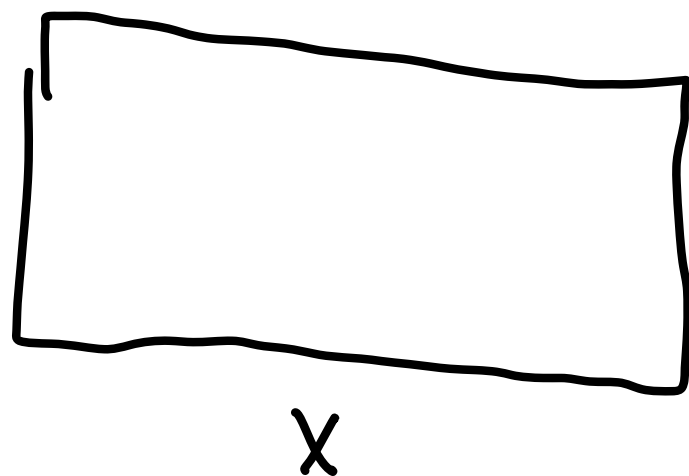
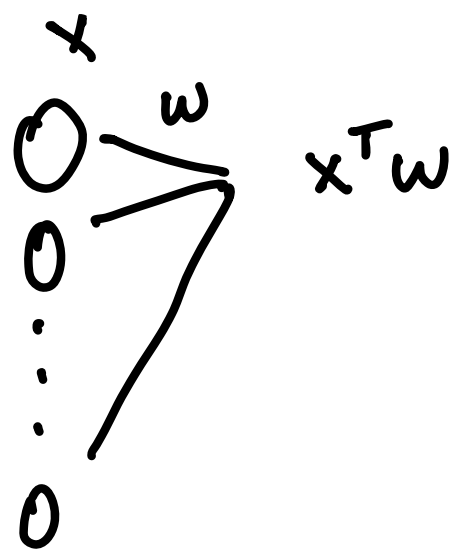
Implicit Regularization

Linear Regression

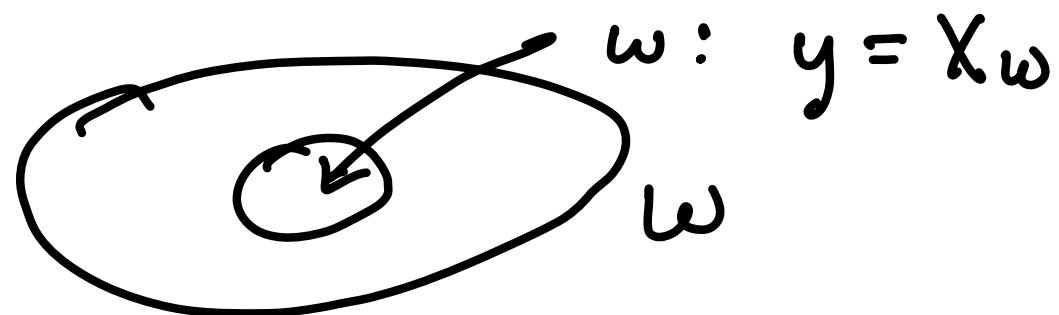
(y, X)

$$y \in \mathbb{R}^n \quad X \in \mathbb{R}^{n \times d}$$

$$y \approx f_{\theta}(x) = Xw \quad w \in \mathbb{R}^d$$



$d > n \Rightarrow$ lots of optimal w



$$\mathcal{L}(w) = \frac{1}{2} \|y - Xw\|_2^2$$

$$- \frac{\partial \mathcal{L}}{\partial w_i} = x_i^T (y - Xw)$$

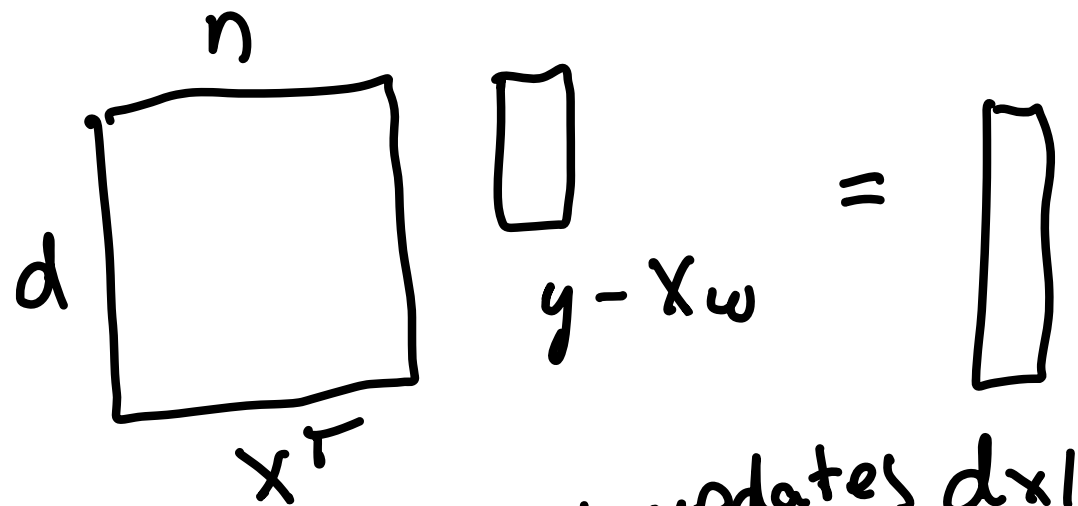
$$- \nabla_w \mathcal{L} = X^T (y - Xw)$$

\uparrow $\begin{matrix} d \times n & n \times 1 & n \times d \times 1 \end{matrix}$

direction of update

$$-\nabla_{\omega} \mathcal{L} = X^T (y - X\omega) = \sum_{i=1}^n x_i \cdot \text{something}$$

↑
span of $\{x_i\}$



$$\omega_0 \xrightarrow{\uparrow} \sum_{t=1}^T \nabla_{\omega} \mathcal{L}(\omega_t) = \hat{\omega}$$

↑ # updates $d \times 1$

Span of $\{x_i\}$

↑
span of $\{x_i\}$

$$\hat{\omega} = X^T \alpha$$

$y = X\hat{\omega}$ ← b/c GD gives an optimal solution

$$X\hat{\omega} = y$$

$$X X^T \alpha = y$$

$n \times d \times n$

assuming $X X^T$ is full rank ($d > n$)

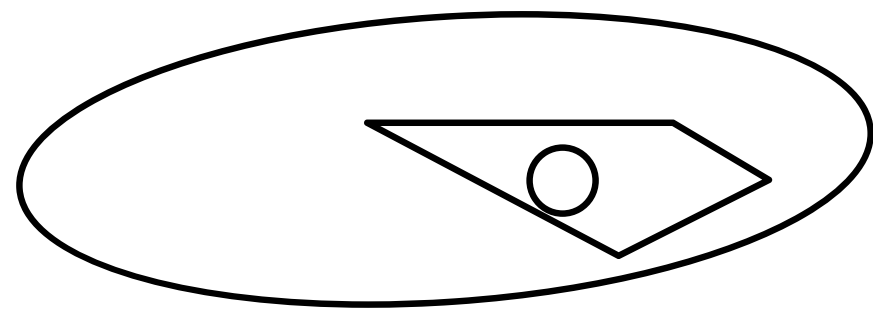
$$(X X^T)^{-1} X X^T \alpha = (X X^T)^{-1} y$$

$$\alpha = (X X^T)^{-1} y$$

$$\hat{\omega} = X^T \alpha = X^T (X X^T)^{-1} y$$

$$\min \|w\|_2^2$$

$$w: Xw = y$$



$$\min_{w \in \mathbb{R}^d} \max_{\beta \in \mathbb{R}^n} \underbrace{\|w\|_2^2 + \beta^T (y - Xw)}_{\mathcal{L}}$$

Solution is saddle point

Optimal $w, b \dots$

$$\nabla_w \mathcal{L} = 0 = 2w - X^T \beta$$

$d \times n \quad n \times 1$

$$\Rightarrow w = X^T \beta \cdot \frac{1}{2}$$

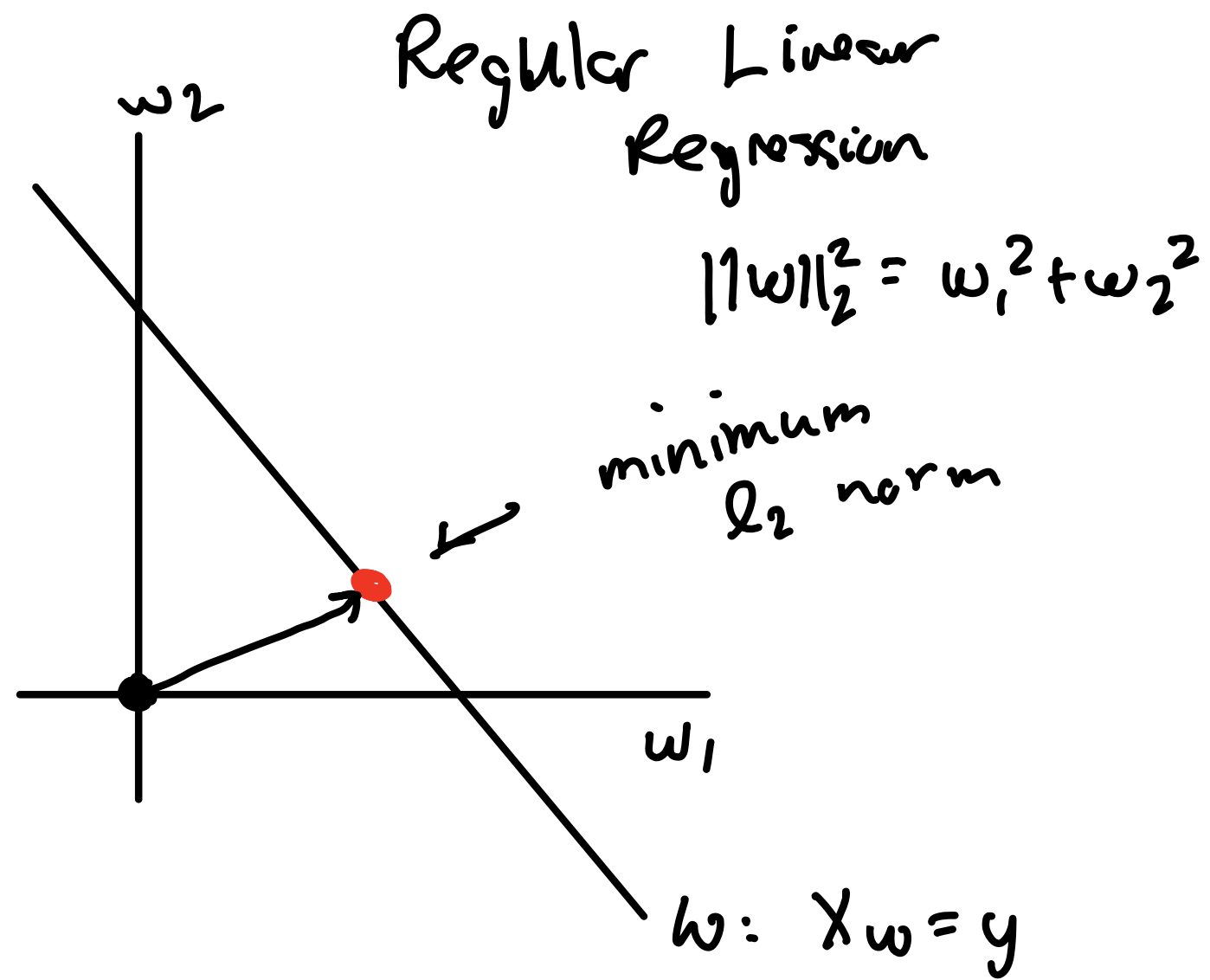
$$\nabla_{\beta} \mathcal{L} = 0 = y - Xw$$

$$0 = y - X X^T \beta \cdot \frac{1}{2}$$

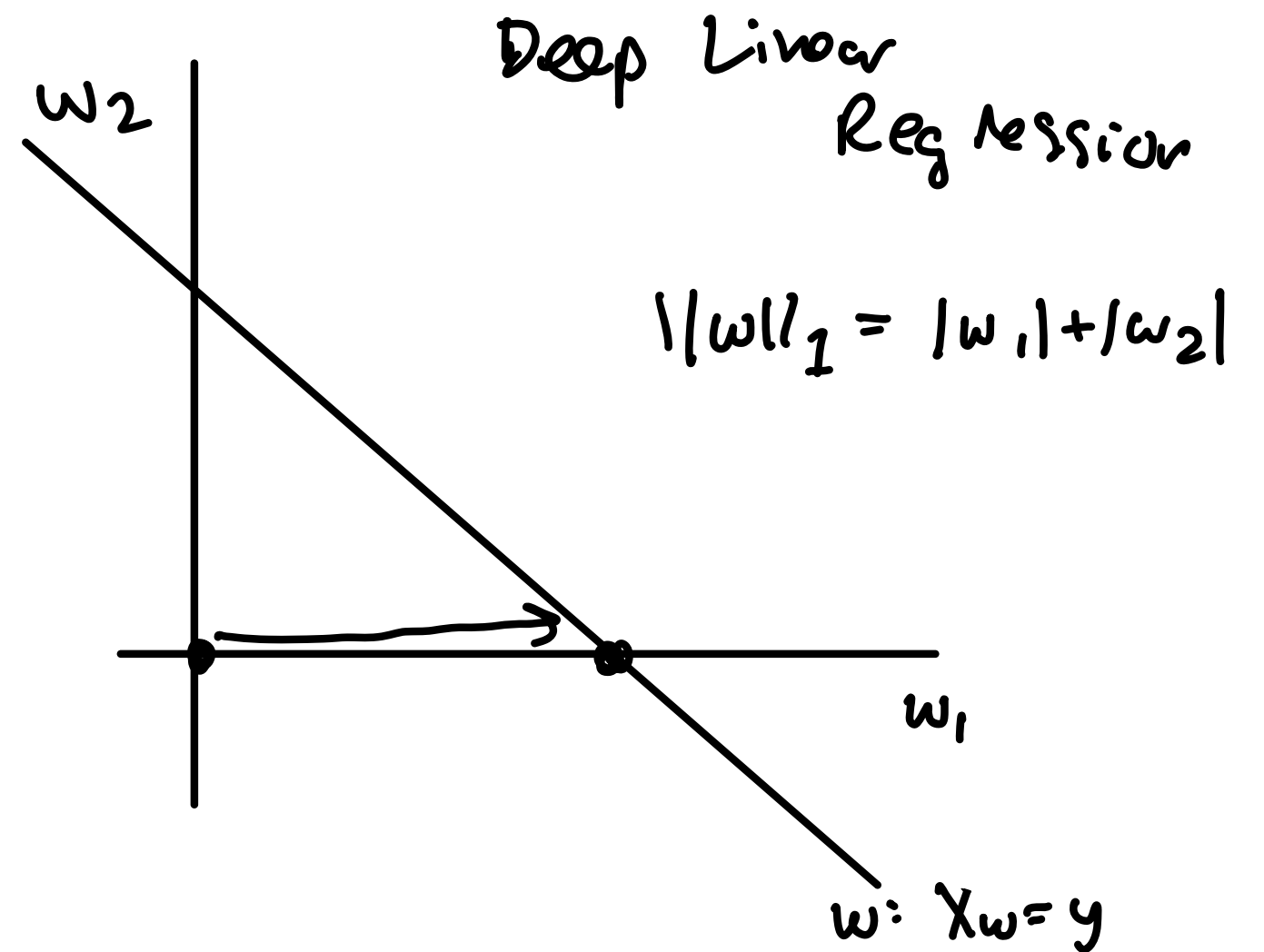
$$2y = X X^T \beta$$

$$\beta = 2 \cdot (X X^T)^{-1} y$$

$$w = X^T \beta \cdot \frac{1}{2} = X^T (X X^T)^{-1} y$$

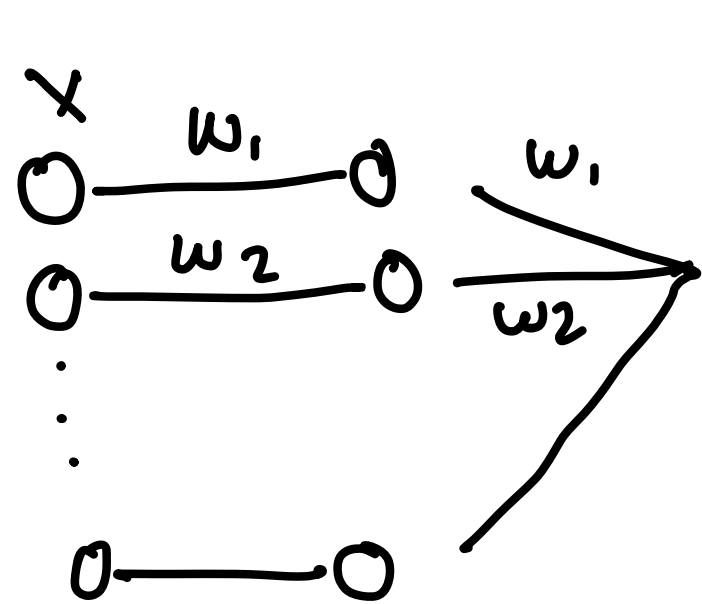
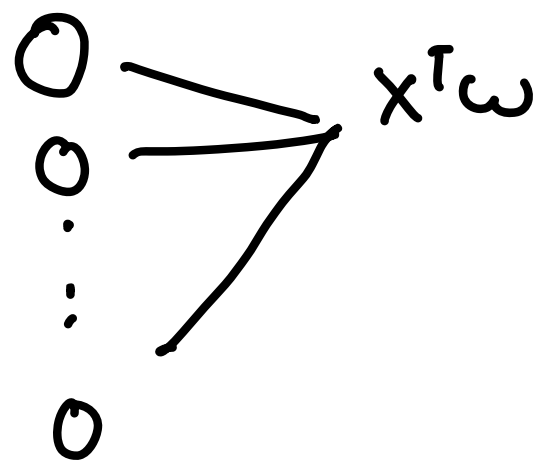


$\hookrightarrow D$ from $w_0 = 0$
 bring us to optimal
 with minimum l_2 norm



$\hookrightarrow D$ from $w_0 = \mathbb{I}$ · small constant
 bring us to optimal
 with minimum l_1 norm

Deeper Linear Regression



$$\sum_{i=1}^d x_i \cdot w_i^2 = x^T (w \circ w)$$

entrywise multiplication

$$\mathcal{L}(w) = \frac{1}{2} \|y - X(w \circ w)\|_2^2$$

$$\frac{\partial \mathcal{L}}{\partial w_i} = x_i^T (y - X(w \circ w)) \cdot \frac{\partial w_i^2}{\partial w_i}$$

$$\nabla_w \mathcal{L} = X^T (y - X(w \circ w)) \circ 2w$$