

"SGD | GD"
Basics of optimization
for ML | AI

→ for-loop
→ Num-Methods

- ✓ (i) ~15 code-walkthroughs
- ✓ (ii) pseudo-code + interactive
- ✓ (iii) Logistic Reg + L₂ reg (SGD)

L₁-SVM

Linear

$$\checkmark D_{tr} = \{(x_i, y_i)\}_{i=1}^m$$

$$\checkmark D = \{(x_i, y_i)\}_{i=1}^n$$

$$x_i \in \mathbb{R}^d$$

GD

Regression:

$$\min_{w, b} \sum_{i=1}^n \left[y_i - (\underline{w}^\top \underline{x}_i + \underline{b}) \right]^2 + \lambda \sum_{j=1}^d w_j^2$$

squared-loss

D_{tr}, w, b

$L(w, b)$

$$\frac{\partial L}{\partial w}, \frac{\partial L}{\partial b}$$

GD / SGD

$$\lambda \underline{w}^\top \underline{w}$$

L_2 -reg

Derivative of $\mathcal{L}(w, b)$

param: w_i 's b

$\partial \mathcal{L}$,
 w
 b

$$\nabla_w \mathcal{L} = \boxed{\frac{\partial \mathcal{L}}{\partial w}}$$

$\sum_{i=1}^n 2(y_i - (w^T x_i + b)) (-x_i) + \lambda(2w)$

Vector

def $\frac{\partial \mathcal{L}}{\partial w}$



$\partial \mathcal{L}, w, b$

$$\boxed{\frac{\partial \mathcal{L}}{\partial b}} =$$

$\sum_{i=1}^n 2(y_i - (w^T x_i + b)) (-1)$

Scalar

Update - eqns:

$$w_{t+1} = w_t - \eta \frac{\partial L}{\partial w} \Big|_{w_t}$$

learning-rate

$$b_{t+1} = b_t - \eta \frac{\partial L}{\partial b} \Big|_{b_t}$$

init: $t=0(w_0, b_0)$

$t=1, 2, 3, \dots$



$$w_1, b_1$$

Gradient-descent:

$$\{(x_i, y_i)\}_{i=1}^n = \underbrace{D_{tr}}_{\text{Training Data}} \underbrace{D_{tst}}_{\text{Test Data}}$$

SGD

1. Pre-processing:
Normalize data \mathcal{D}_{tr}

Sklearn \rightarrow Dim-redu

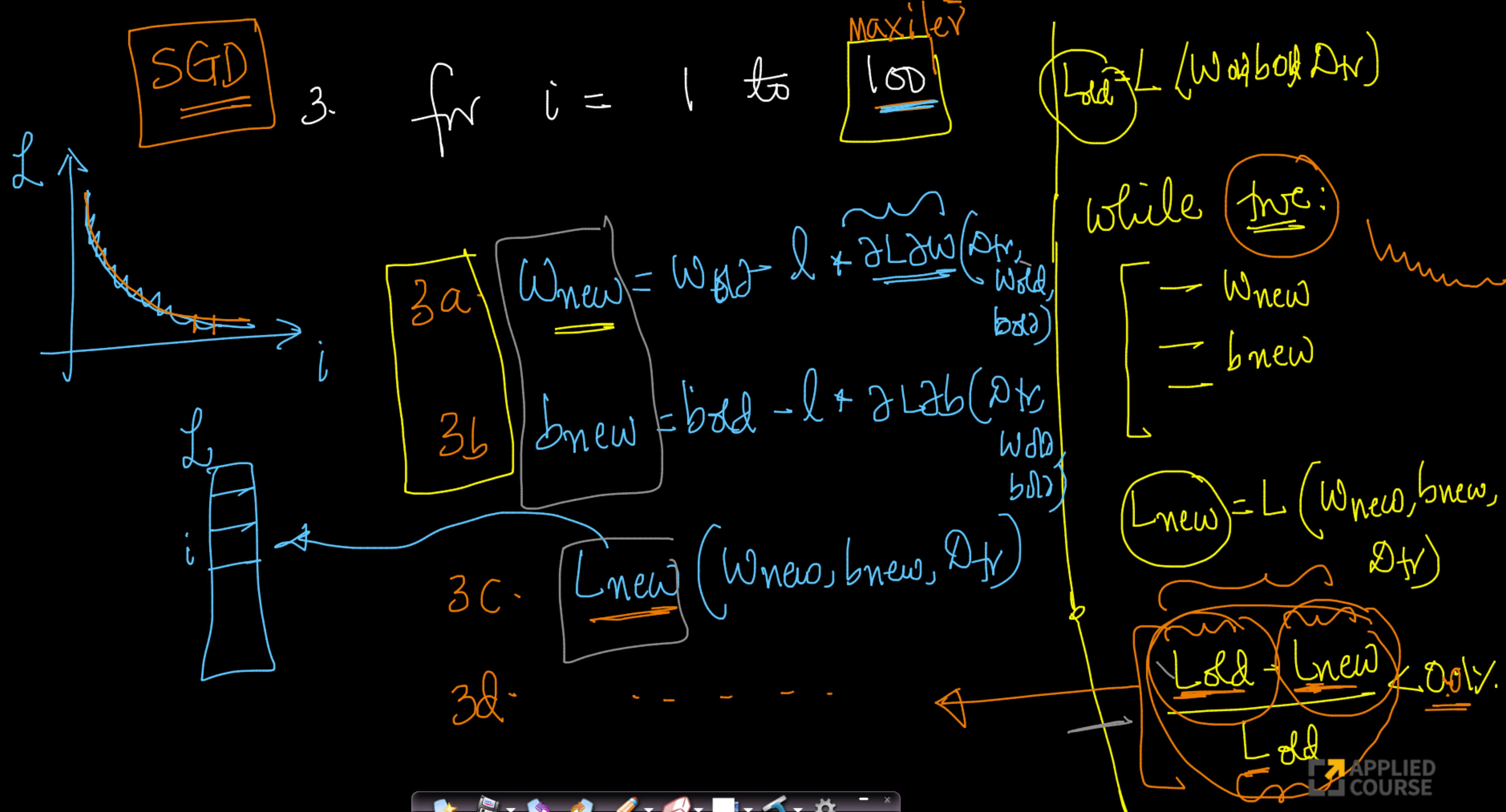
- probability
- boundary
- intervals

2. init: hyper-params

$$\lambda = 1 \quad l = 0.5$$

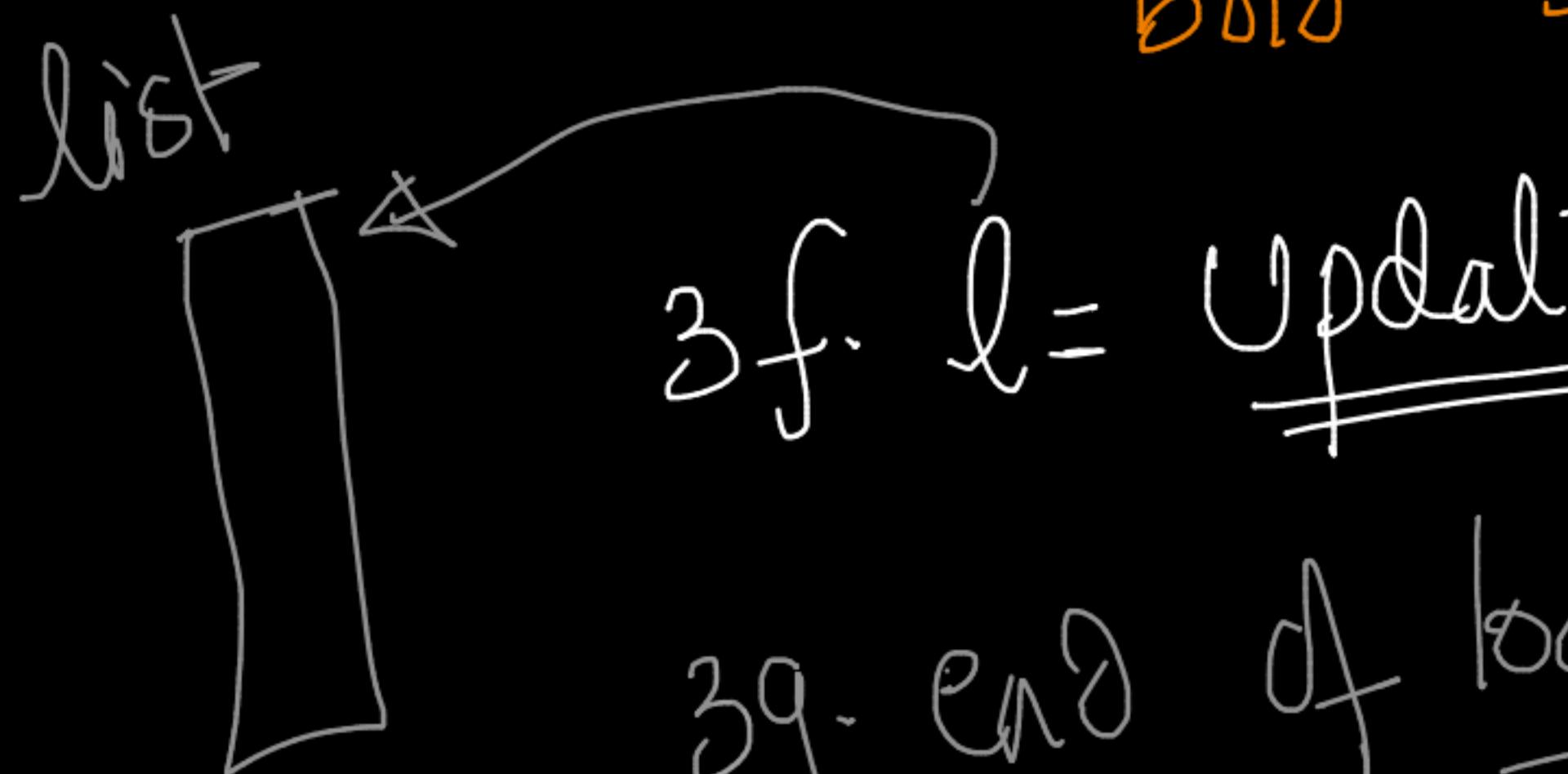
$\checkmark n = \checkmark$
 $\checkmark d = \checkmark$

$\left\{ \begin{array}{l} w_{old} = \text{random-vector} \\ b_{old} = \text{scalar - randomly} \\ \quad \text{random.seed(100)} \end{array} \right.$



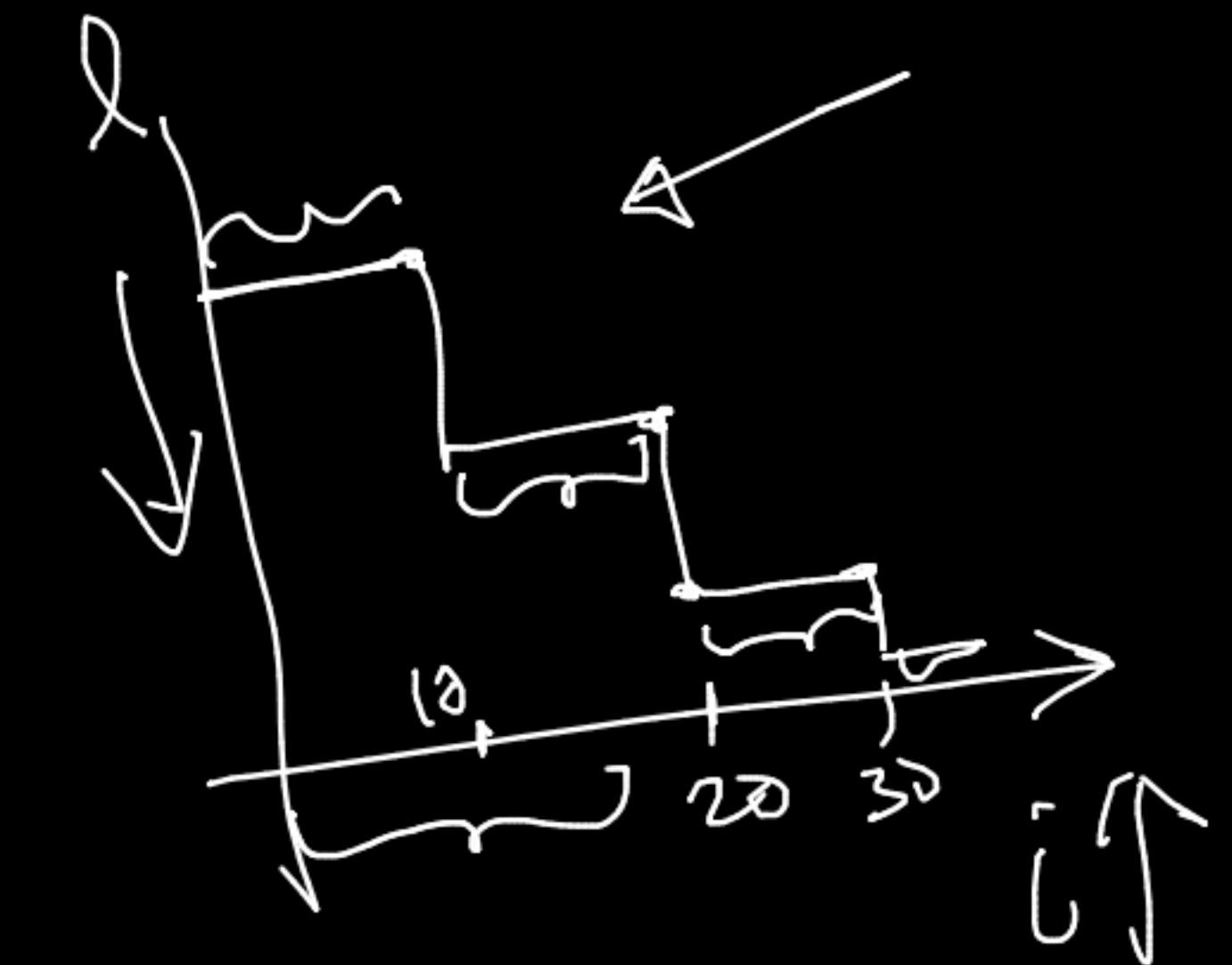
3e- $w_{\text{old}} = w_{\text{new}}$

$b_{\text{old}} = b_{\text{new}}$



3g. end of loop

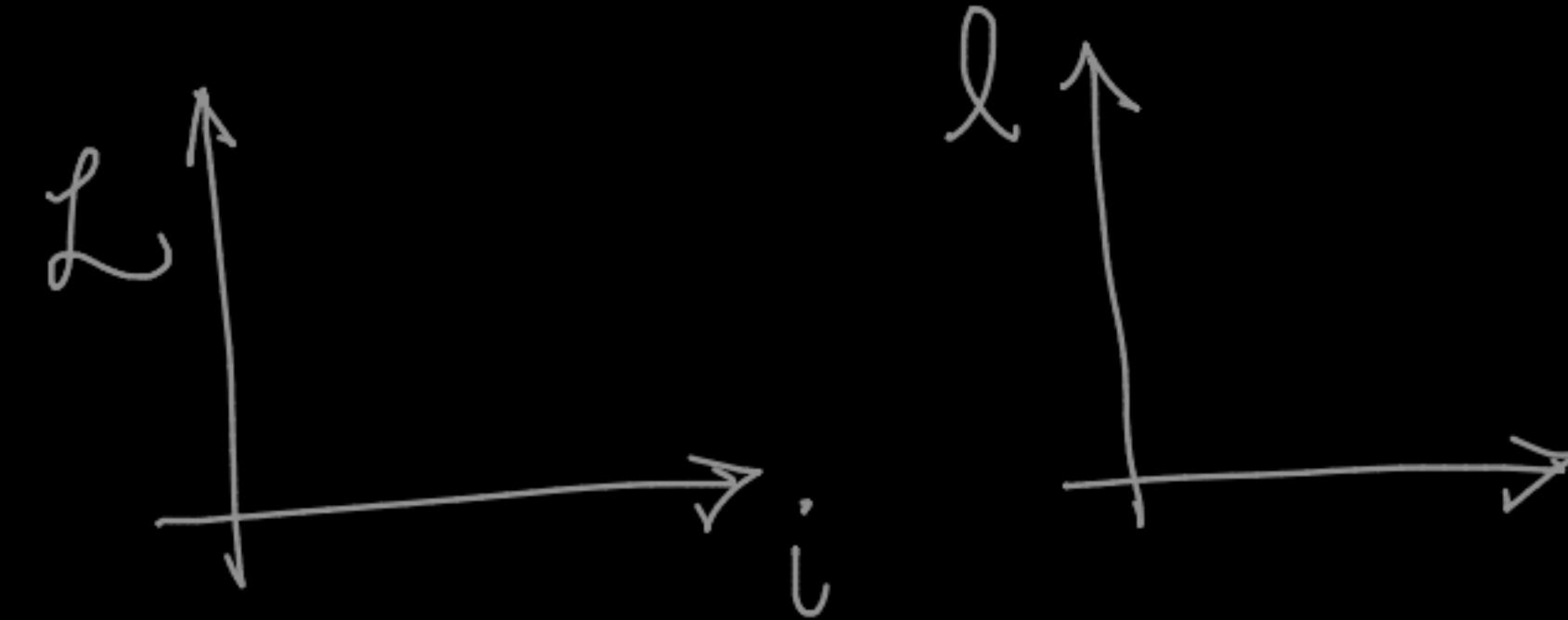
$w_{\text{new}}, b_{\text{new}}$



✓ $l = 1/(i)$

✓

4. plotting



5. return $w_{\text{new}}, b_{\text{new}}$

2 hrs

Eval Model $(w, b, \underline{D_{cv}})$

SQ. loss

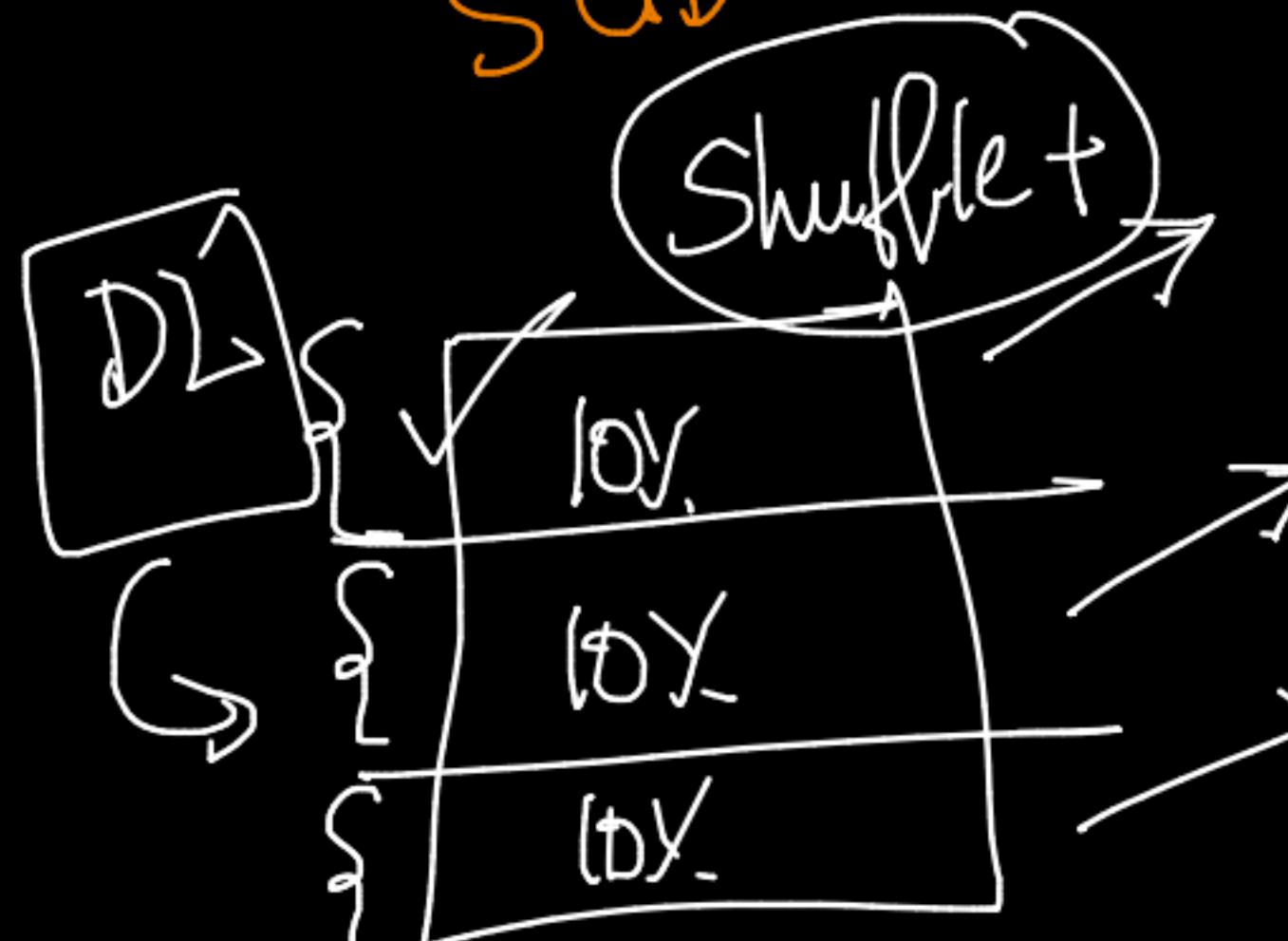
x_i, y_i SGD:

$$3a. \quad w_{\text{new}} = w_{\text{old}} - l * \boxed{dLdw(w_{\text{old}}, b_{\text{old}}, \cancel{\text{DK}})}$$

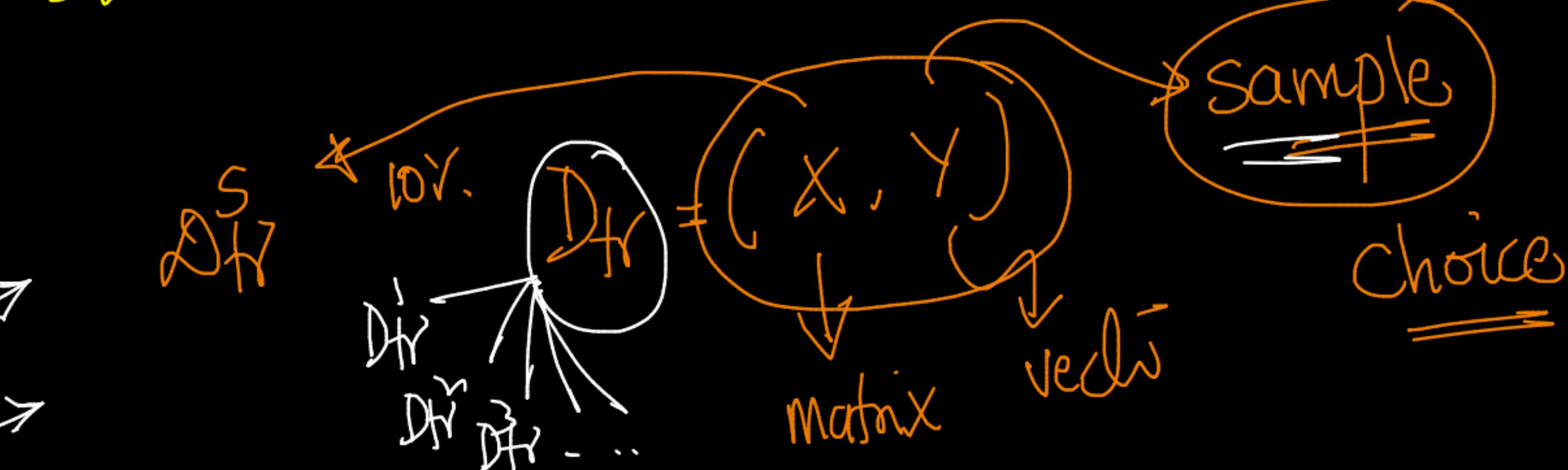
$$3b. \quad b_{\text{new}} = b_{\text{old}} - l * \cancel{dLdb(w_{\text{old}}, b_{\text{old}}, \cancel{DN})}$$

GD

SGD



Ioy.



SGD

$$\text{each iter} \rightarrow \frac{\partial L}{\partial w} \Big|_{w_t}$$

faster



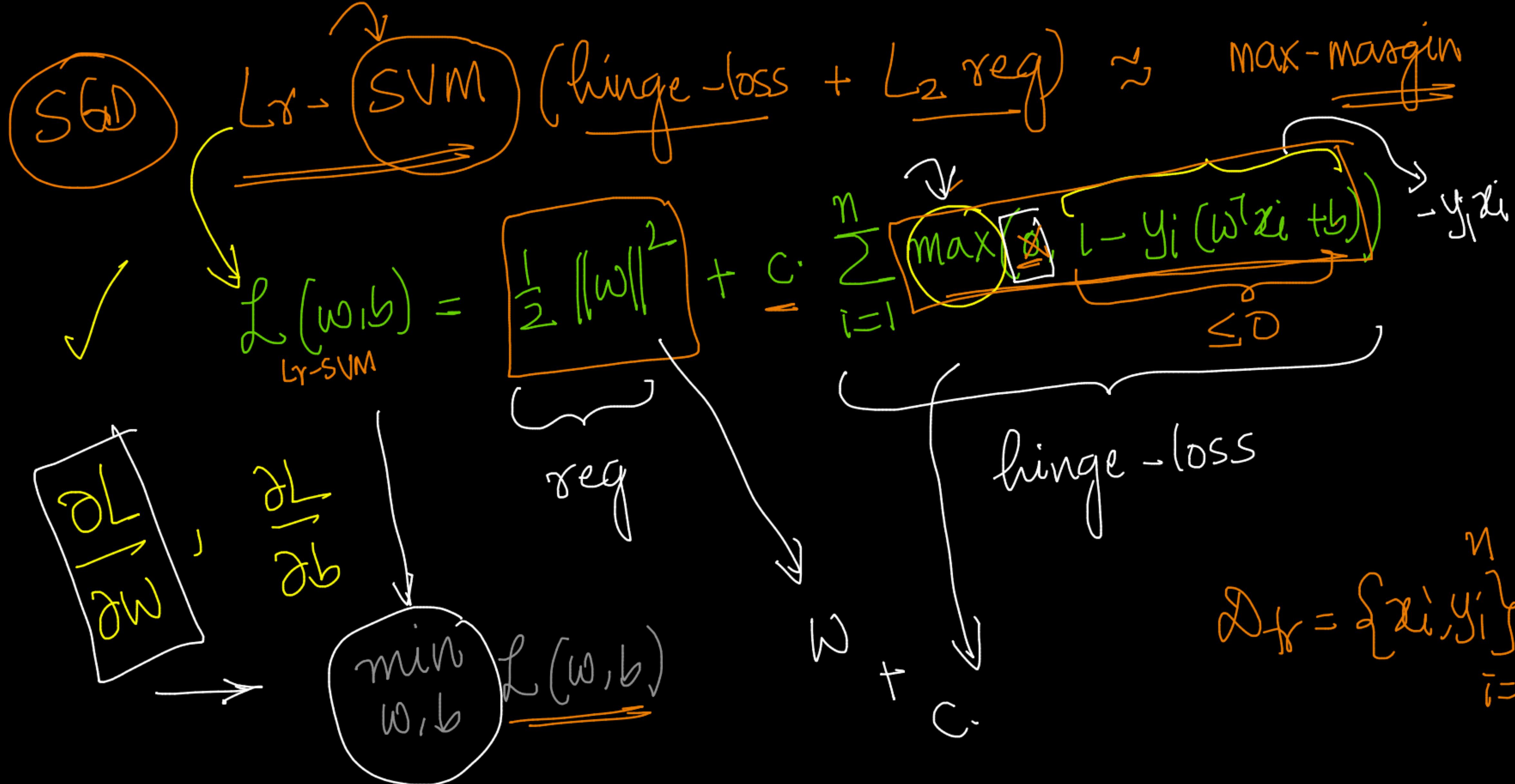
more - iters



GD

slower / Memory issues

fewer iter



Derivative of hinge-loss:

$$\frac{\partial L}{\partial w}$$

Vector \rightarrow

$$\frac{\partial L}{\partial w} =$$

word

$$w +$$

$$c \quad \checkmark$$

$$\sum_{i=1}^n$$

$$\begin{cases} 0 & \text{if } \\ -y_i x_i & \end{cases}$$

$$1 - y_i (w^T x_i + b) \leq 0$$

otherwise

$$\frac{\partial L}{\partial b}$$

$$\frac{\partial L}{\partial b} =$$

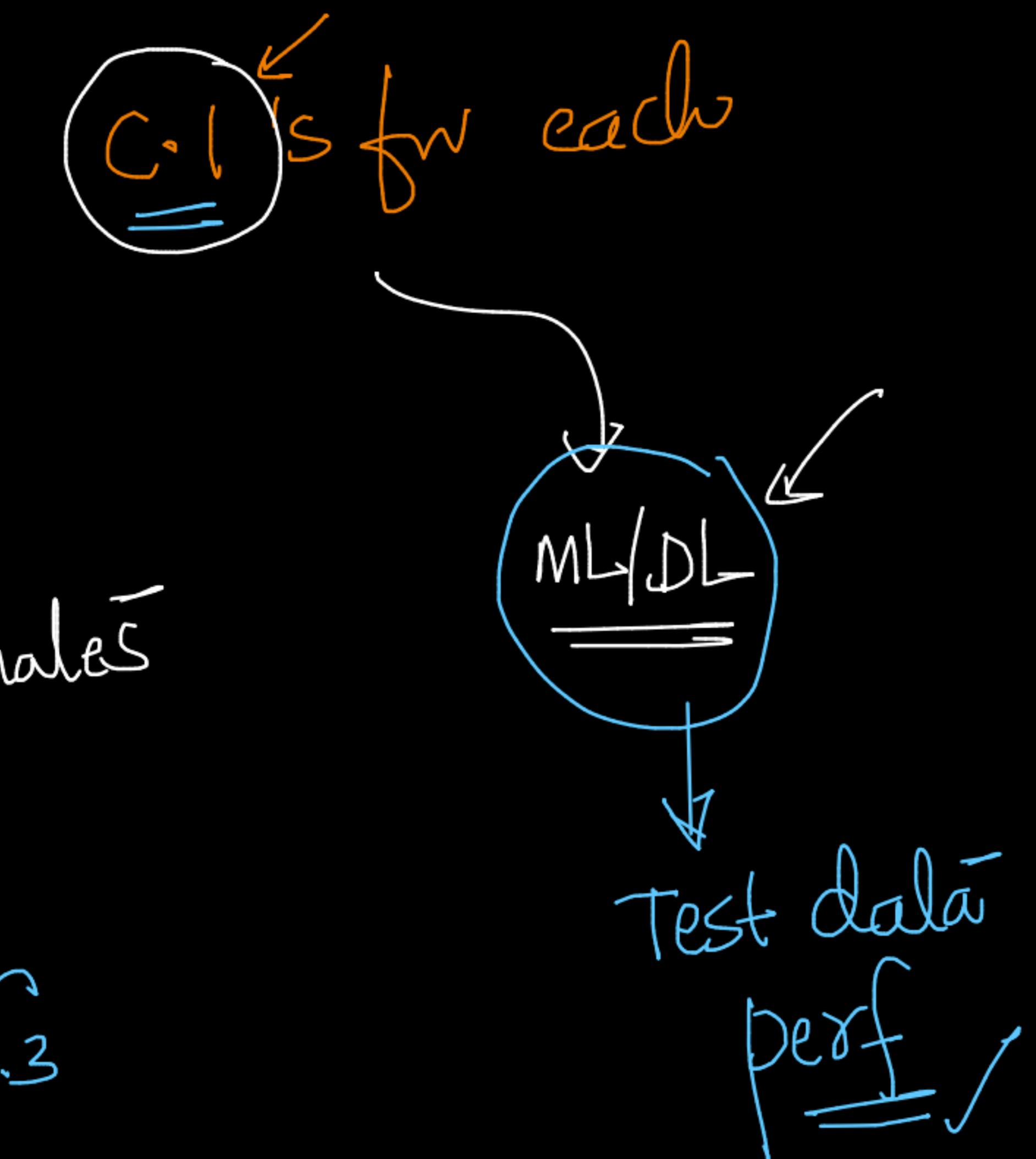
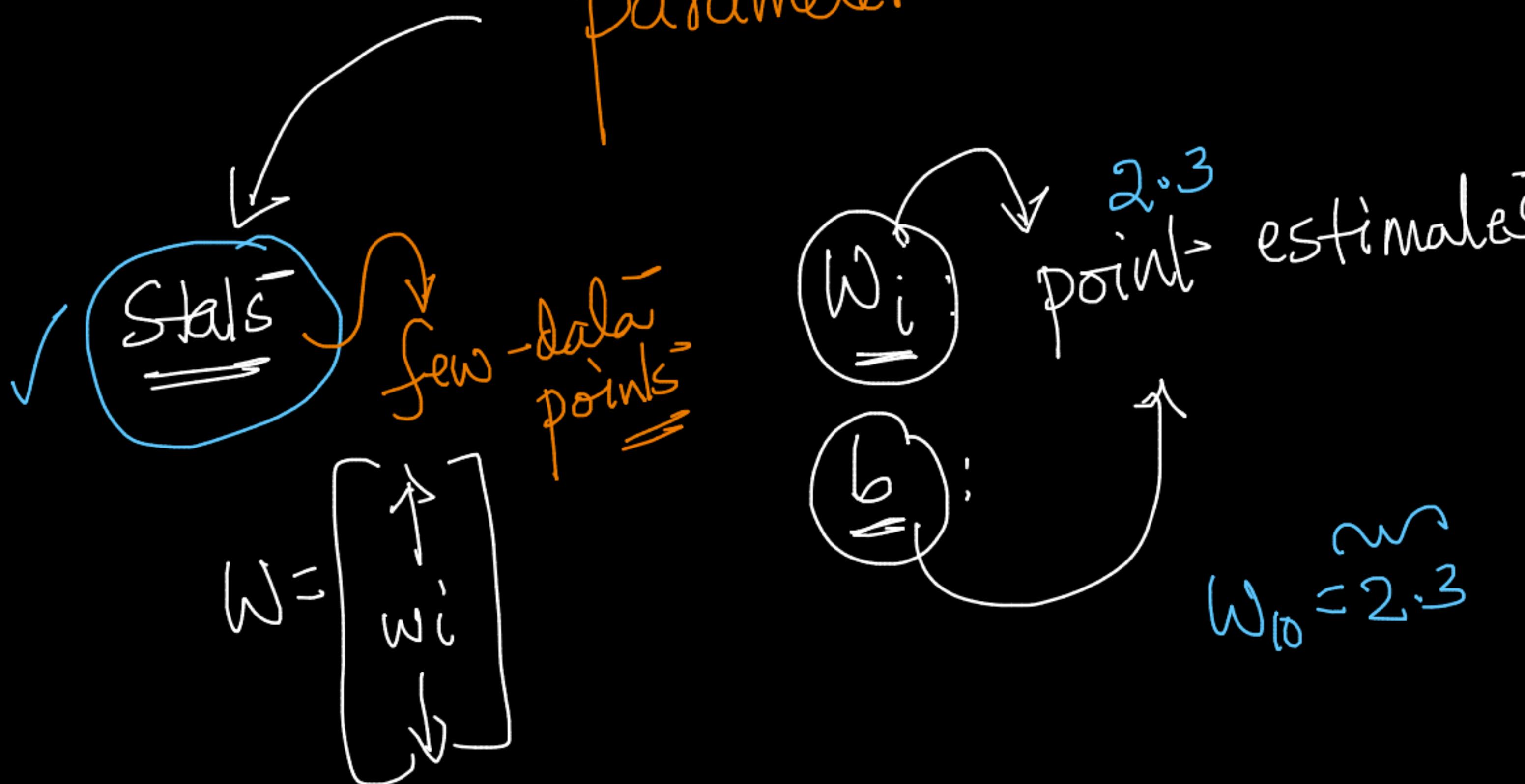
$$0 + c \sum_{i=1}^n \begin{cases} 0 & \text{if } \\ -y_i & \end{cases}$$

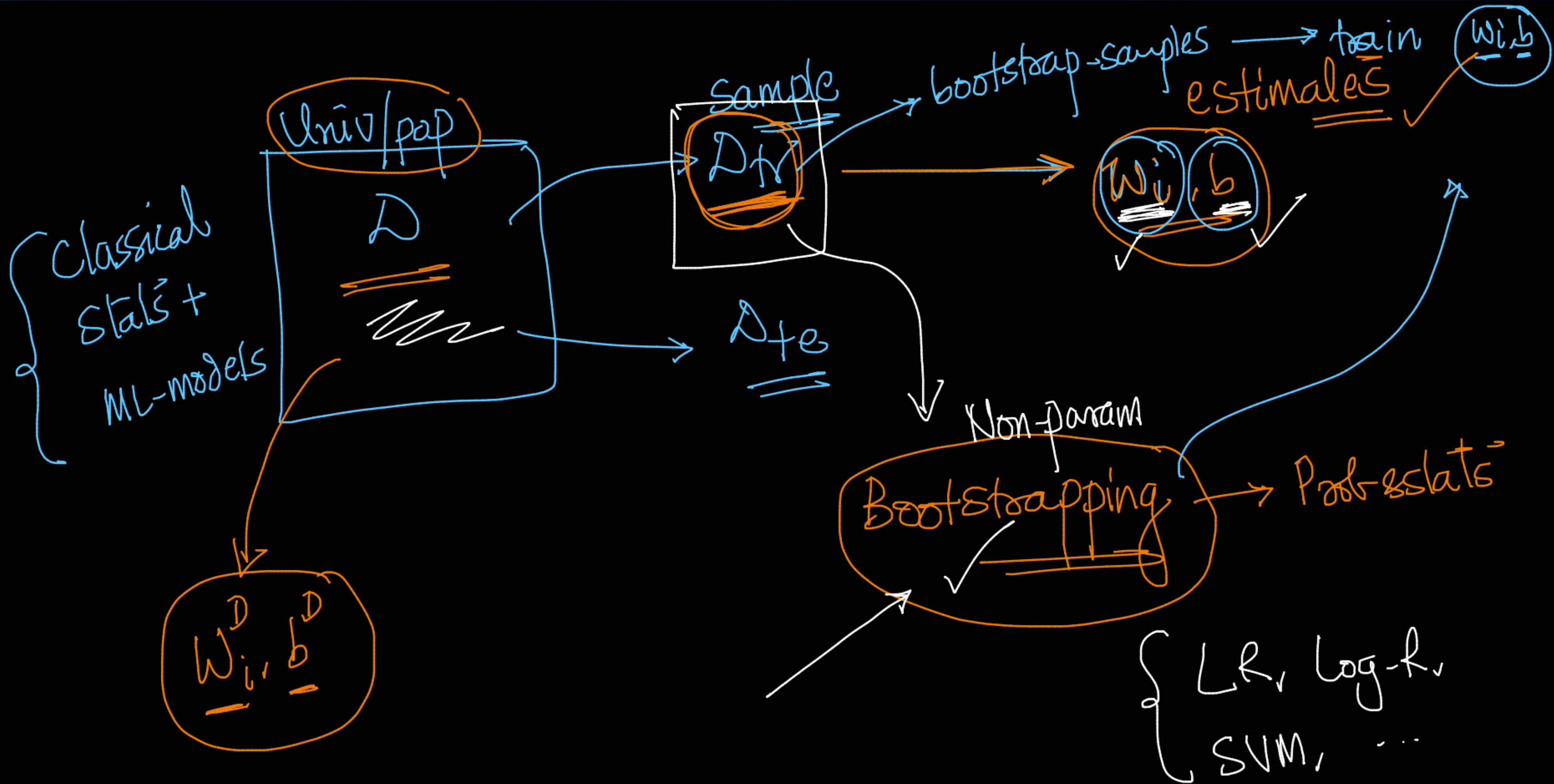
$$1 - y_i (w^T x_i + b) \leq 0$$

otherwise

(Q) how can we obtain $\underline{C_L}$'s for each

parameters: w_i or b





(Q) Can we obtain p-values for each param?

