

Sampling Techniques - Continuation 05/02/2025

MALA - Metropolis Adjusted Langevin Algorithm

Markov Chain Monte Carlo Method (MCMC) that combines Langevin Dynamics (a gradient-based proposal mechanism) with the Metropolis-Hastings algorithm (an acceptance-rejection step).

This combination ensures that the Markov chain converges to the correct target dist. $p(x)$, even when the Langevin dynamics are discretized.

MALA uses Langevin dynamics to propose new samples :-

$$\left(dx_t = - \nabla U(x_t) dt + \sqrt{2} dB_t \right)$$

gradient of the
potential
function

potential function

increment of
standard Brownian
motion.

Discretized using Euler - Metropolis Scheme.

$$\underline{x_{k+1}} = \underline{x_k} - \underbrace{\gamma \nabla U(x_k)}_{\text{step size}} + \sqrt{2\gamma} z_k.$$

Langevin Dynamics.
Standard normal R.V.

$$z_k \sim N(0, 1).$$

The proposal is biased towards regions of higher probability (lower $U(x)$) due to the gradient term $-\gamma \nabla U(x_k)$.

The proposal dist $Q(x'|x)$ describes the probability of proposing a new state x' given the current state x .

For MALA - the dist is Gaussian:-

$$Q(x'|x) \propto \exp\left(-\frac{1}{4\gamma} \|x' - x + \gamma \nabla U(x)\|^2\right)$$

$x - x + \gamma \nabla U(x) \rightarrow$ diff between x'
(proposed state)

and the predicted state (using Langevin dynamics)
 $(x - \gamma \nabla U(x))$

The Gaussian dist ensures that the proposals are more likely to be close to the state predicted by the Langevin Dynamics.

To correct the discretization error in Langevin Dynamics, MALA uses Metropolis Hastings

acceptance probability:

without this the discretization error could cause the Markov chain to converge to wrong prob. dist.

$$\alpha = \min \left\{ 1, \frac{p(x_{k+1}) Q(x_k | x_{k+1})}{p(x_k) Q(x_{k+1} | x_k)} \right\}$$

$p(x) \propto e^{-U(x)} \rightarrow$ Target distribution

$Q(x'|x) \rightarrow$ proposal distribution \rightarrow Gaussian

Acceptance ratio: $\frac{p(x_{k+1}) Q(x_k | x_{k+1})}{p(x_k) Q(x_{k+1} | x_k)}$

✓

This ensures that the Markov chain satisfies the detailed balance equation which guarantees the convergence to the target distribution $p(x)$

Step by Step MALA Algorithm :-

→ x_0 → start with the initial state

→ for each step k :-

- Propose a new state x_{k+1} → using Langevin dynamics proposal.

$$x_{k+1} = x_k - \gamma \nabla U(x_k) + \sqrt{2\gamma} \epsilon_k$$

- Compute the acceptance probability α

$$\alpha = \min \left\{ 1, \frac{p(x_{k+1}) Q(x_k | x_{k+1})}{p(x_k) Q(x_{k+1} | x_k)} \right\}$$

- Accept x_{k+1} with prob. α , otherwise keep x_k

→ Burn-in → discard the first few samples to allow the chain to converge.

Metropolis Hastings Algorithm:-

MCMC method used to sample from a target prob. dist $P(x)$. It is particularly useful when $P(x)$ is difficult to sample from directly, but its unnormalized density can be evaluated.

Idea:-

Construct a Markov chain whose stationary dist is the target distribution $P(x)$.

Proposing a new state x_i from a proposal dist. $Q(x|x_t)$, accepting/rejecting the proposed state based on acceptance prob. A .

P.T.O

Metropolis Hastings Algorithm:-

Start \rightarrow initial state $x_t = x_0$

Propose a new state x_i from the proposal dist $Q(x|x_t)$.

Compute acceptance probability A :-

$$A = \min \left(\frac{P(x_i) Q(x_t|x_i)}{P(x_t) P(x_i|x_t)}, 1 \right)$$

- this ratio ensures that the Markov chain satisfies detailed balance, which guarantees convergence to the target dist $P(x)$.

- Accept / Reject:-

- Sample $k_i \sim \text{Uniform}[0,1]$

if $k_i < A$, then accept the proposed state $x_{t+1} = x_i$ otherwise, reject, i.e.,
 $x_{t+1} = x_t$.

Theory of Input/Outputs via Model.

Target = T (Ground Truth, Labels)
 Obtained = Y (model-forward()).
 \hookrightarrow Model's output.

10 digit - binary numbers.

----- = 2^{10} To model this.
#1's > #0's ? $\rightarrow 1/0 \rightarrow$ model.
1010110110 = 1
1111110000 = 1
0000000010 = 0
1111100000 = 0

Activation?
 $0/1$
(Sigmoid)

I/P \rightarrow [Model] \rightarrow o/p.

(1×10)

$(10 \times 1) \rightarrow$

~~1×10
 10×1~~

1×1

$(0 - 1)$

(-1 to 1)
 ↑ ↑
 0 1

$\left\{ \underline{0} \text{ or } \text{+ve.} \right\} = \text{valid.}$

$(\infty \rightarrow \text{NaN})$

corrupt.

#1 > #0's

6
 ↓
 4
 ↓
 1

7 3
 8 2
 9 1
 10 0

W

1	-	1
1	-	1
1	-	1
0	-	1
0	-	1
0	-	1
1	-	1
1	-	1
0	-	1

to = -5

12411

#1s → 5, #0s → 5

1000

6 - 5 ⇒ 1

10 - 5 ⇒ (5)

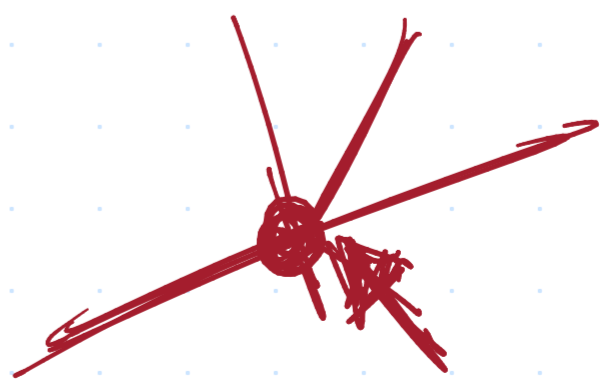
LR = 10^{-4}

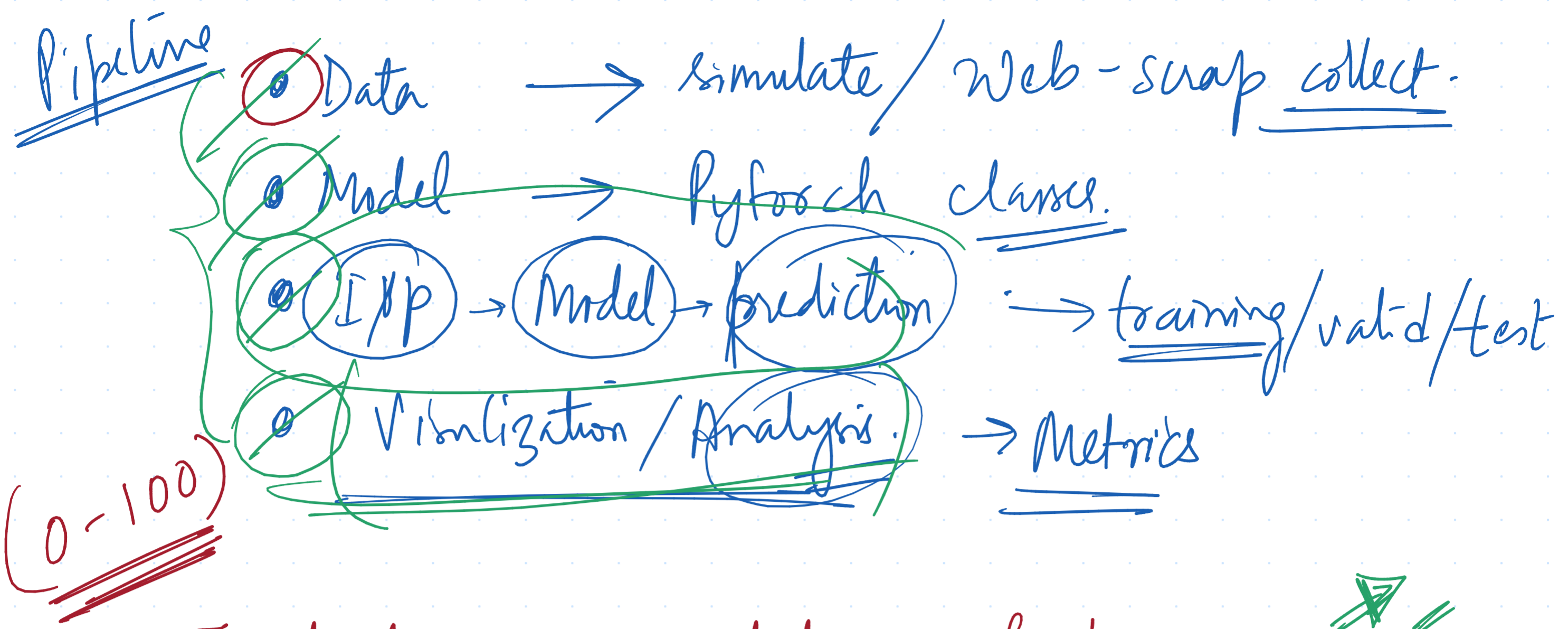
with = W - W

1000 10^{-4}

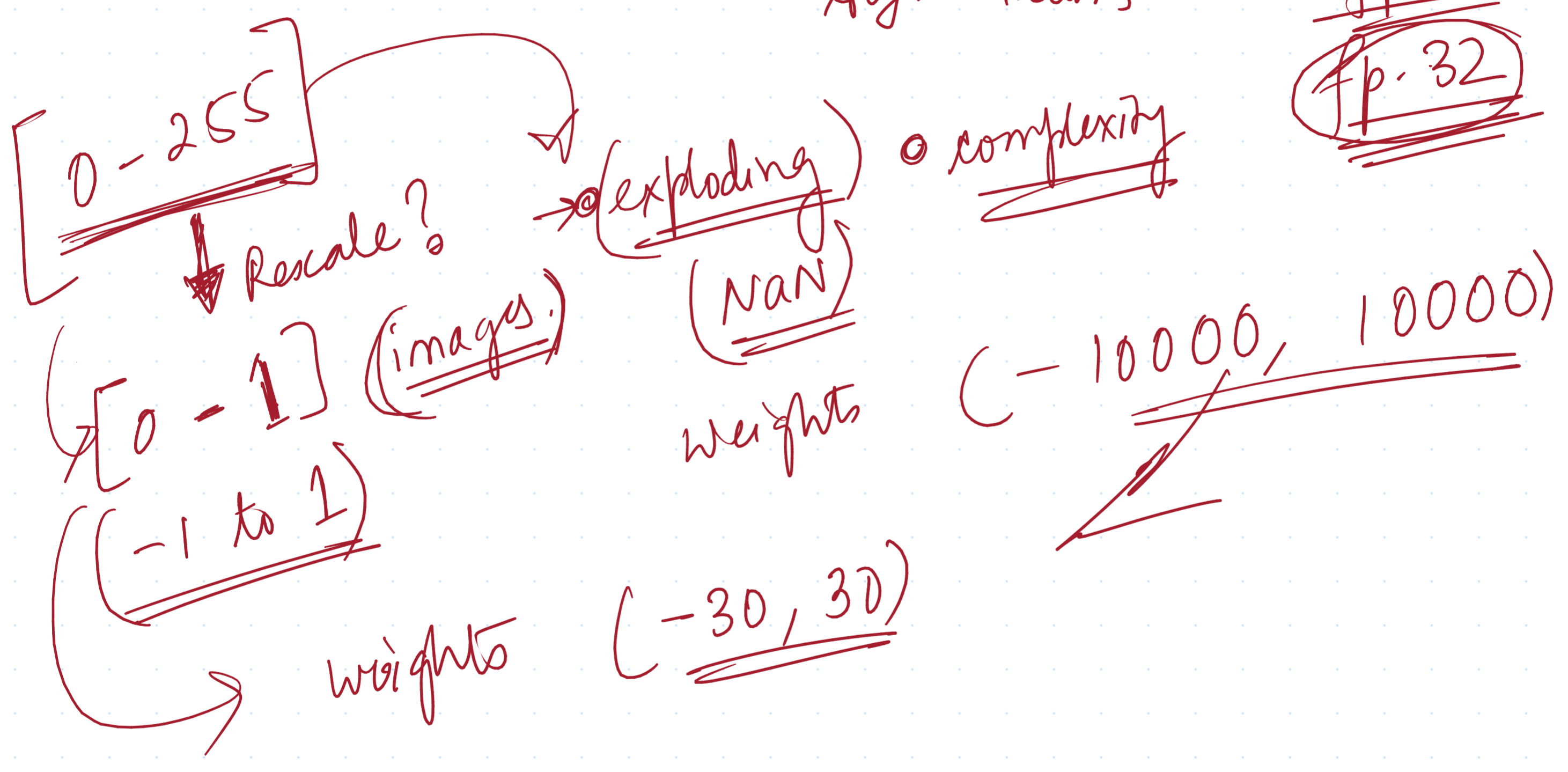
W = 1000

(0 to 1)
 ↳ Sigmoid.

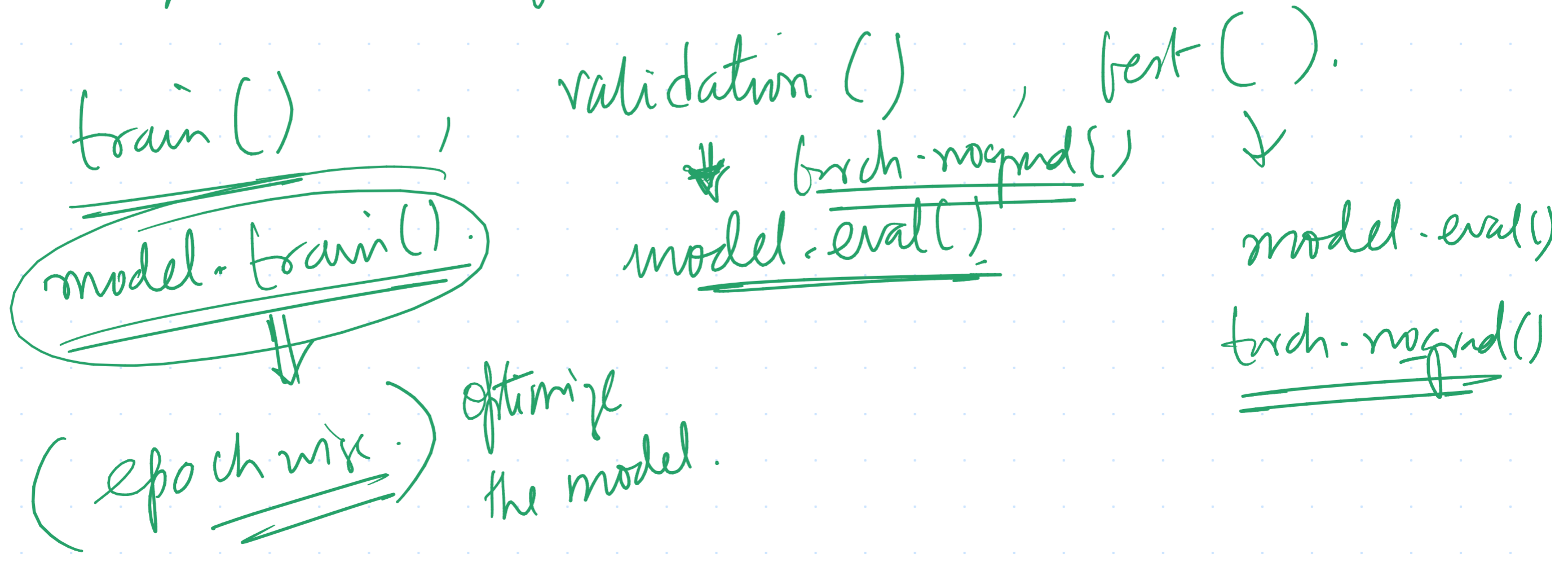




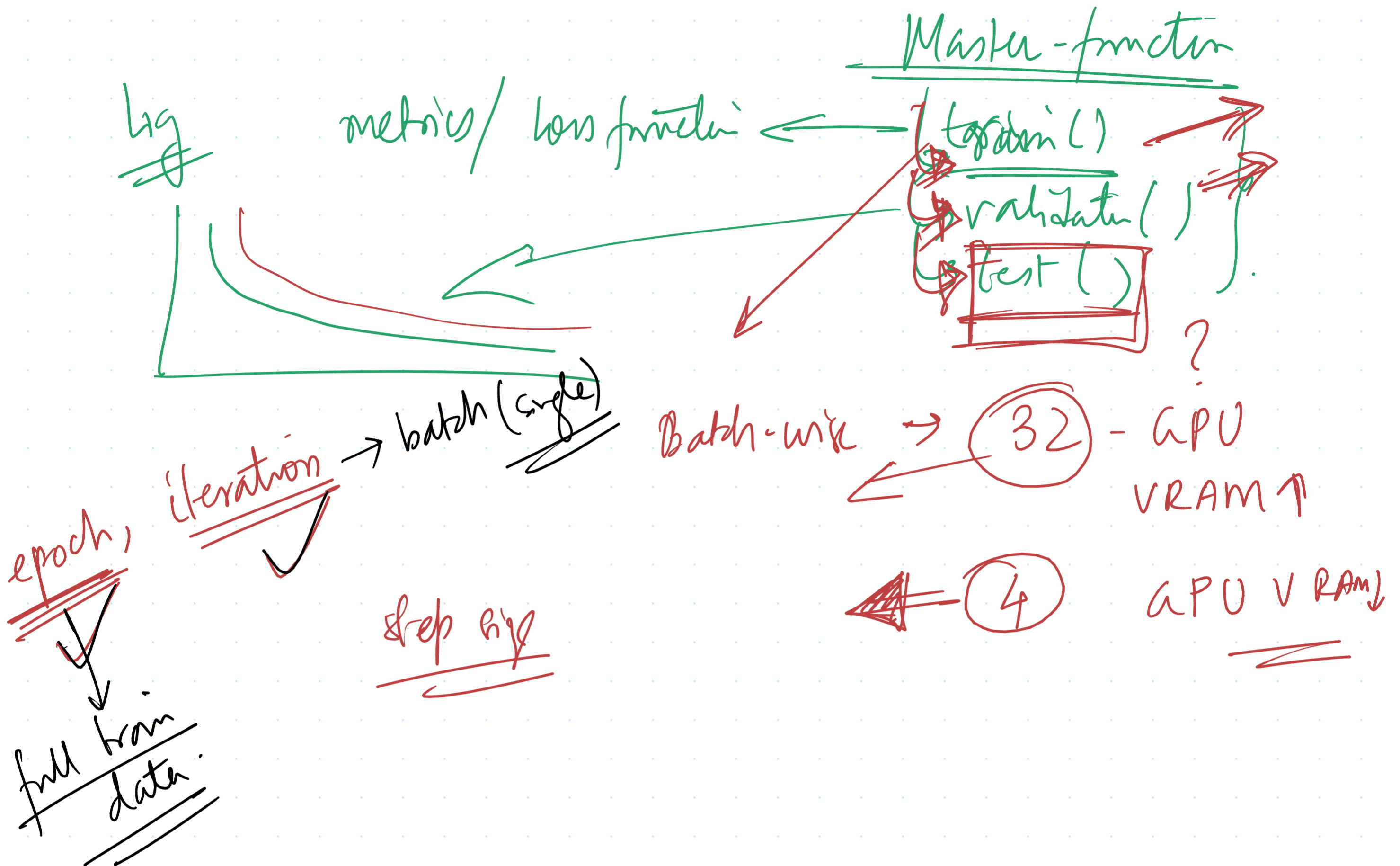
Dataloader → raw data → preprocessing
 augmentations



Model → Pytorch.



model checkpointing, model saving, model loading



Target $\rightarrow T$

obtained \Rightarrow

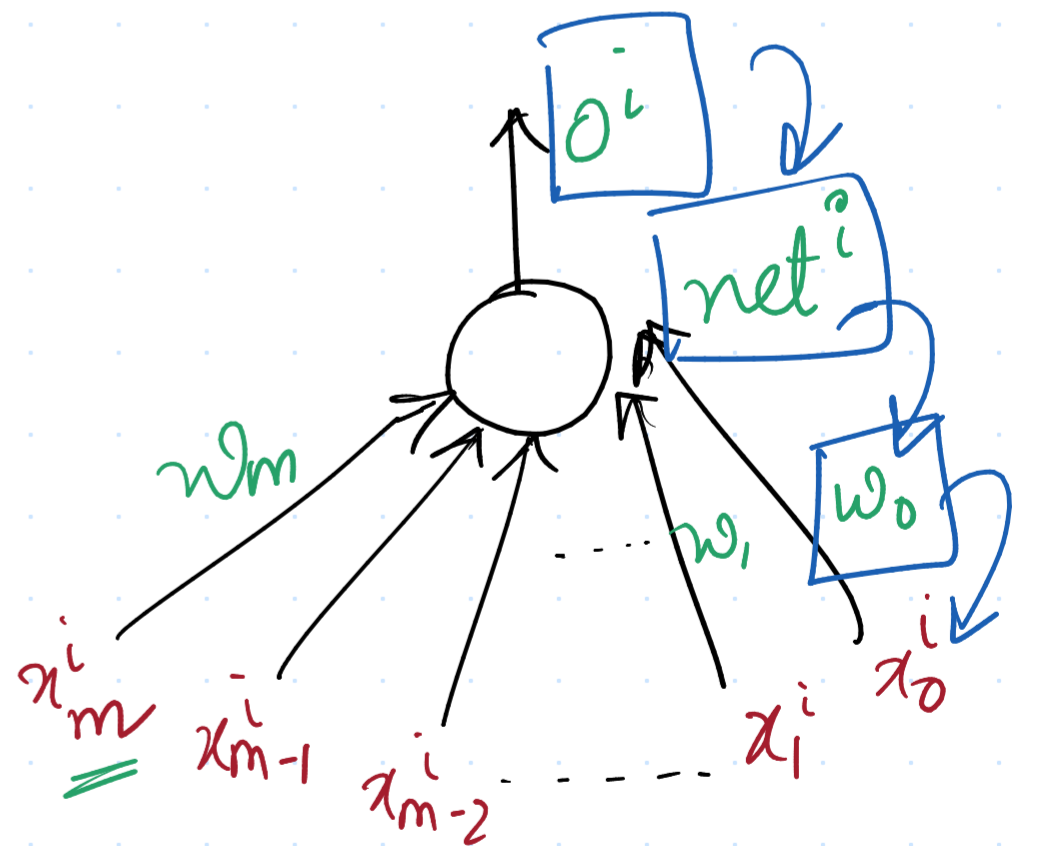
Sigmoid

$$o^i = \frac{1}{1 + e^{-x}} \rightarrow \text{net}^i$$

$$o^i = \frac{1}{1 + e^{-\text{net}^i}}$$

$x^i \rightarrow$ i^{th} sample.
 $x_j^i \rightarrow$ j^{th} scalar

$$\text{net}^i = W \cdot x^i = \sum_{j=0}^m w_j x_j^i$$



$$X = \begin{matrix} 0^{\text{th}} \\ 1^{\text{st}} \\ 2^{\text{nd}} \\ 3^{\text{rd}} \end{matrix} \left[\begin{array}{cccccccccc} \textcircled{1} & 0 & \textcircled{1} & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 1 & 1 & 0 & 1 & 1 & \textcircled{1} & 0 \\ 1 & 1 & 1 & 1 & \textcircled{1} & 1 & 1 & 1 & 1 & 1 \end{array} \right] \begin{matrix} \rightarrow \\ \rightarrow \\ \rightarrow \\ \rightarrow \end{matrix} \left. \begin{matrix} y \\ 0 \\ 1 \\ 1 \\ 1 \end{matrix} \right\} y^0$$

$x'_j \rightarrow x''_2$
 x''_8
 x^3_4

$$x^0 \leftarrow [x^0_0, x^0_1, x^0_2, \dots, x^0_m] \rightarrow [y^0]$$

Vectors \rightarrow ~~W~~ capital letters \leftarrow

$$x^1 \leftarrow [x^1_0, x^1_1, x^1_2, \dots, x^1_m] \rightarrow [y^1]$$

Scalar \rightarrow small letters \leftarrow

$$x^{n+1} \leftarrow [x^{n+1}_0, x^{n+1}_1, \dots, x^{n+1}_m] \rightarrow [y^{n+1}]$$

$(m+1)$

$x^i \rightarrow$ i^{th} input vector.

$o_i \rightarrow$ output scalar.

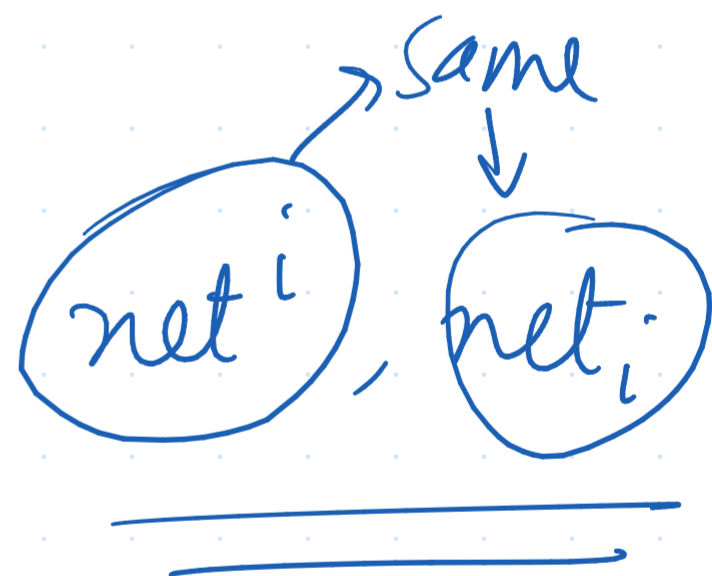
$$\text{net}_i \rightarrow W x^i \quad (\text{m-th } \underline{\text{i/p / o/p}})$$

$W \rightarrow$ weight vector.

W & x^i has $(m+1)$ components

$$W = \langle w_m, w_{m+1}, \dots, w_2, w_1, w_0 \rangle$$

$$x^i = \langle x^i_m, x^i_{m-1}, \dots, x^i_2, x^i_0 \rangle$$



Derivative of Sigmoid :-

$$o^i = \frac{1}{1 + e^{-net_i}}$$

$$\log(o^i) = -\ln(1 + e^{-net_i})$$

Derivative w.r.t. netⁱ :-

$$\frac{1}{o^i} \cdot \frac{\partial o^i}{\partial net_i} = \frac{1}{(1 + e^{-net_i})} (e^{-net_i})$$

$$\frac{1}{o^i} \frac{\partial o^i}{\partial net_i} = \frac{e^{-net_i}}{1 + e^{-net_i}} = 1 - \frac{1}{1 + e^{-net_i}}$$

↓
oⁱ

$$1 + \frac{x}{x+1} \rightarrow \left(1 - \frac{1}{1+x}\right) \frac{x}{1+x}$$

↓

$$\frac{1}{o^i} \frac{\partial o^i}{\partial net_i} = (1 - o^i)$$

$$\frac{\cancel{1+x} - \cancel{x}}{1+x} \Rightarrow \frac{x}{1+x}$$

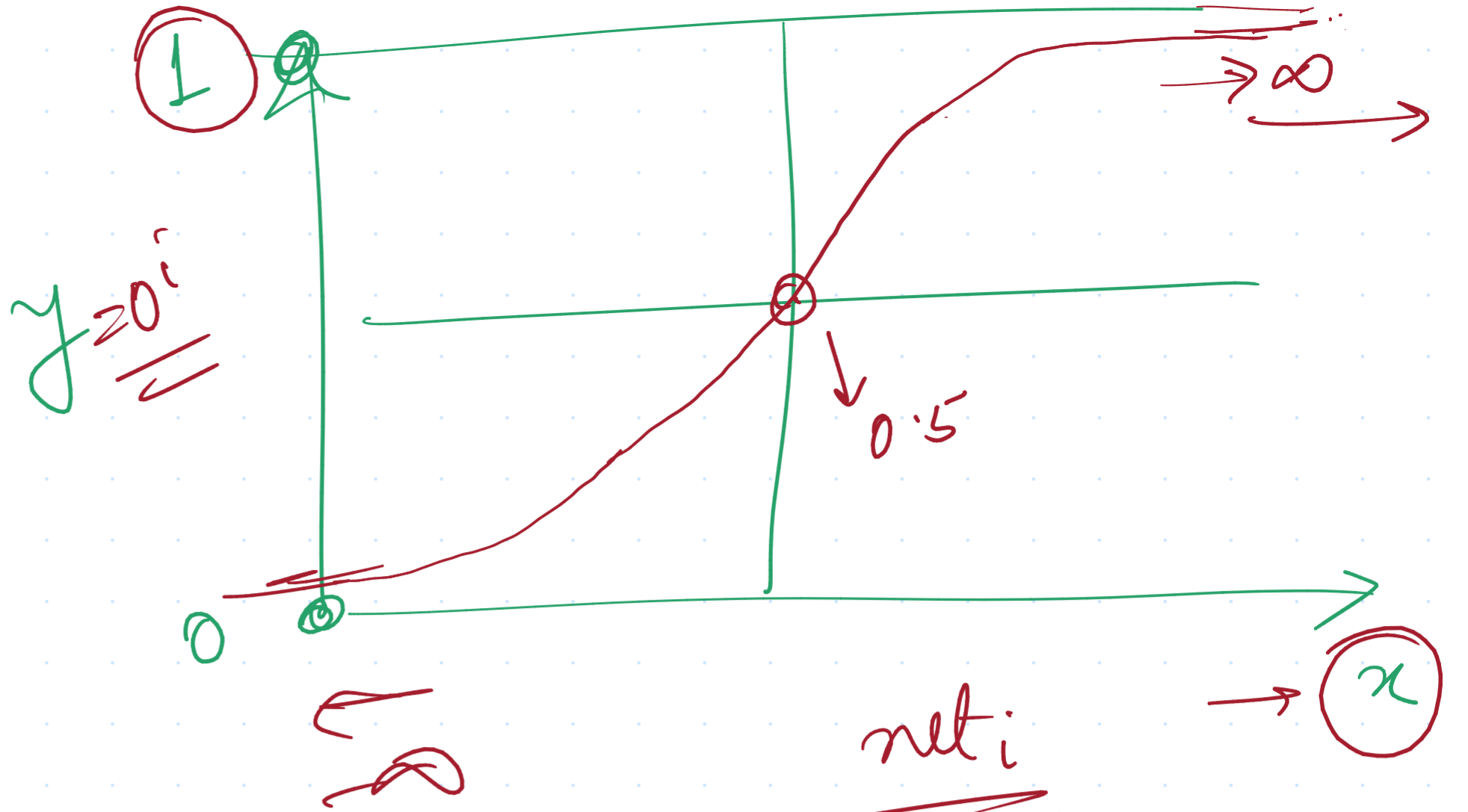
$$\frac{\partial o^i}{\partial net_i} = \underline{\underline{o^i(1 - o^i)}}$$

when $x = \infty$;

$y = 1$

when $x = -\infty$

$y = 0$



when $x = 0.5$,
 $y = ?$

$y = ?$ $x = 0.5$

where is the derivative of sigmoid maximum?

$y = \frac{1}{1 + e^{-0.5}}$

$y = 0.6224$

$\frac{\partial}{\partial x} (x(1-x)) \Rightarrow \frac{d}{dx} (x - x^2) \Rightarrow 1 - 2x$

Derivative is maximum at $1 - 2x = 0$.

$x = \frac{1}{2}$

$\therefore x = \frac{1}{2}$

But

maximum value of derivative = ?

maximum value of derivative

maximum value of derivative of sigmoid $\frac{1}{4}$

$x(1-x) \Rightarrow \frac{1}{2} (1 - \frac{1}{2})$

$\Rightarrow \frac{1}{2} \times \frac{1}{2} \Rightarrow \underline{\underline{0.25}}$

* When $x = \frac{1}{2}$ derivative of sigmoid is maximum.

value of sigmoid $\frac{1}{1+e^{-x}}$ $\Rightarrow \frac{1}{1+e^{-\frac{1}{2}}} \approx \underline{\underline{0.62}}$

Derivative of sigmoid $\Rightarrow o_i(1-o_i)$
function. $\Rightarrow \frac{1}{(1+e^{-x})} \left(1 - \frac{1}{1+e^{-x}}\right)$
 $\Rightarrow 0.62 \times (1 - 0.62)$

$\Rightarrow \underline{\underline{0.2356}}$

$net_i = \infty$

$o_i = \frac{1}{1+e^{-net_i}}$

$net_i = -\infty$

$o_i = \frac{1}{1+e^{\infty}} = \frac{1}{e^{\infty}} = 0$

$o_i = \frac{1}{1+e^{-\infty}}$
 $= \frac{1}{1+\frac{1}{e^{\infty}}}$

$o_i = \frac{1}{1+0} = 1$

$$\frac{\partial}{\partial a} \left(\frac{1}{1+e^{-a}} \left(1 - \frac{1}{1+e^{-a}} \right) \right) = 0.$$

$\boxed{x} \rightarrow$ different?

$\frac{1}{2}x$ different \boxed{x} .