* **Sample, Batch, Epoch**

→ A sample is a single row of data.

→ The batch size defines the number of samples to work through, before updating the internal model parameters.

   • Depending on the batch size different learning algos are defined.

→ The number of epochs is a hyperparameter that defines the number of times the learning algo will work through the entire training dataset.

   • one epoch

* **An epoch**

one epoch :

1. Randomly divide training set into $m = N/k$ batches.
2. Use a batch of training samples to compute $J(w)$.
3. Update $w$ as: $w^{t+1} = w^t - \eta \dfrac{\partial I}{\partial w}$

4. Repeat 2 & 3 using different subsets all samples are used once.

Q. what is the size of the batch?

   - 1 ...... k ...... N

   - May depend on hardware.

   - we repeat epochs until convergence.

1 .... k
k+1 .... 2k
2k+1

## * Batch Gradient Descent

→ Training set: Ex: ImageNet has 14M images

→ Approach:

- Compute the loss $J(w)$ on the entire training, update the parameters $w$.
- At the next epoch, shuffle the training data, and repeat above process.

→ Typical batch size: Size of the training set.

## * Mini-batch Gradient Descent

→ It is wasteful to compute the loss over the entire set to perform a single parameter update for large datasets.

- Ex: ImageNet has 14M images
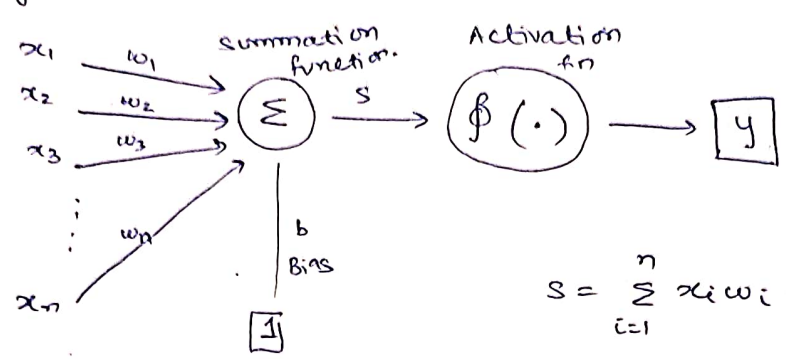- GD is replaced with mini-batch GD

→ Mini-batch GD

- Approach:
  - Compute the loss $L(w)$ on a batch of images, update the parameters $w$

How to learn w & b? How will our model learn?

learn   loss fn            How many neurons
         gradient descent           in one layer.

$x_1$   $w_1$    summation function.      Activation fn

$x_2$   $w_2$    $\Sigma$   $s$    $\phi(\cdot)$  →  $y$

$x_3$   $w_3$

    $w_n$    b

      Bias

$x_n$

$$s = \sum_{i=1}^{n} x_i w_i + b \qquad y = \phi(s)$$

## error
## /loss/ cost/ objective function

→ supervised learning.

→ The loss function provides the cost of being wrong, by measuring the quality of a particular set of parameters based on how well the output of the network agrees with the ground ~~labe~~ truth labels in the training data.

                Input, features    → label, ground truth

$$L(\theta) = \text{distance} \left( (f_\theta(x), y) \right)$$

       error             → label (true)       $y_{pred} = f_\theta(x)$

                                            learned by the model

     difference b/w actual & predicted

            $\theta$ = parameters (weights, biases)

learning process



$y = w_0 + w_1 x_1 + w_2 x_2$

depending upon
feedback.

start random value (for slope).

1. Start with random values of $w_i$ (training data)
2. Evaluate the goodness of the line, determined with a loss function, $J(w)$.
3. The weights $w_i$ are changed acc moving the line to a better position.
   → Note: $J(w)$ should be min when the training samples are correctly classified.
4. Repeat 2 & 3 until $J(w) < \gamma$

Gradient Descent → linear regr$^n$ / NN / CNN / RN

learns parameters.

low loss function: Gradient Descent in Action.



$w_1$
$w_2$

Force you to go to dir$^n$ in which max descent is less.