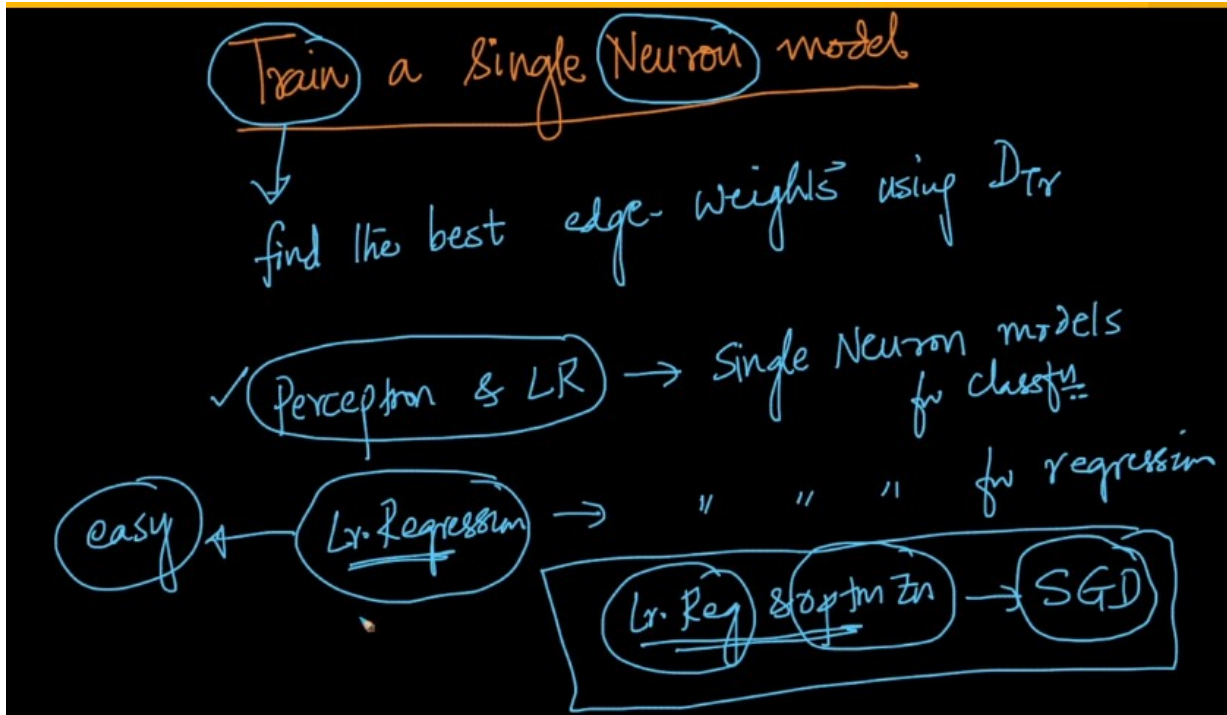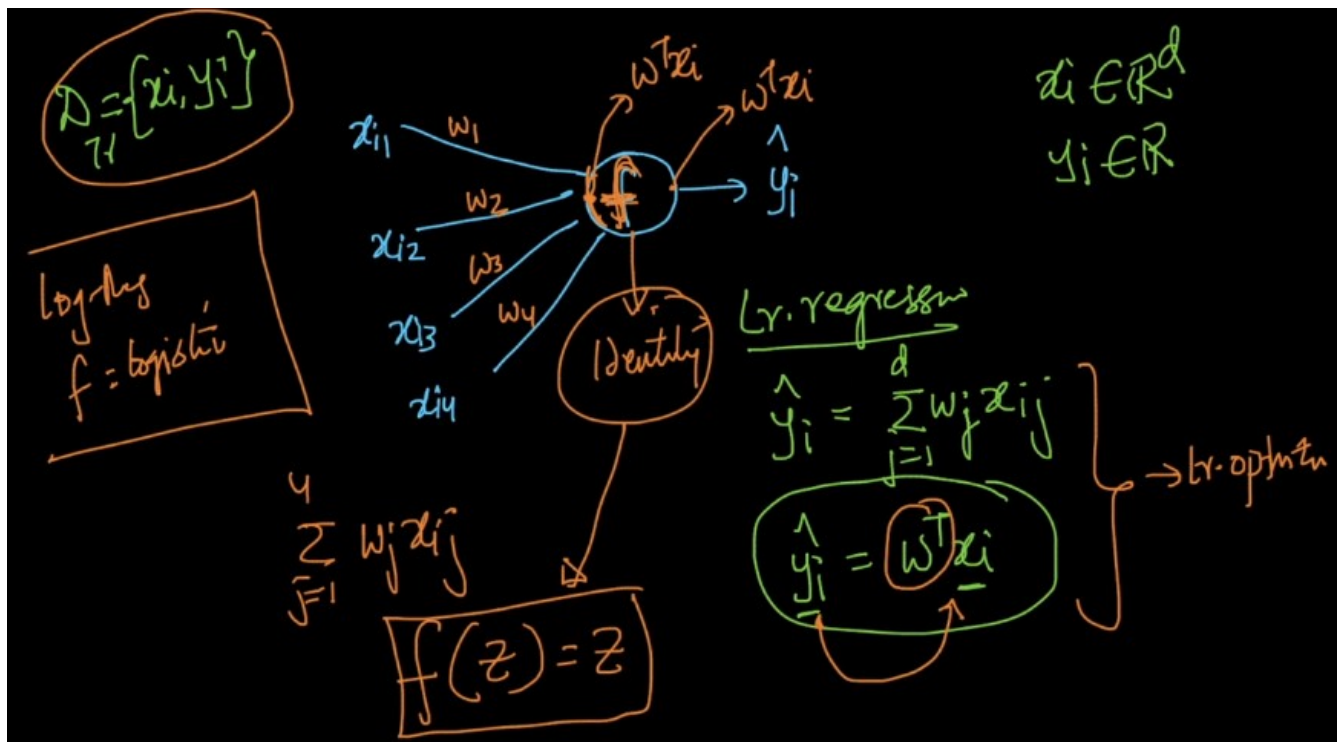# Neural Networks

Train a single Neuron Model:
To find the best edge weights using Dtrain data, Perception and Log reg are single neuron models for classification.



Linear regression problem: Linear regression and optimization. Here F is an identity function, $F(x) = x$

Logistic regression: But in case of logistic regression F is a logistic function.
The optimization problem for linear regression is as follows.

$$\text{Lr. regrs}^n:$$

$$\underset{w_i}{\min} \quad \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \text{reg} \qquad n \rightarrow \#\text{pts Tr.date}$$

$$\hat{y}_i = w^T x_i$$

$$\text{optmzn} \qquad \underset{w_i}{\min} \quad \sum_{i=1}^{n} (y_i - w^T x_i)^2 + \|w\|_2^2$$

Neural networks in Depth -

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

$$x_i \in \mathbb{R}^n, \quad y_i \in \mathbb{R}$$

① Define loss-fn

$$\mathcal{L} = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \text{reg}$$

$$\mathcal{L} = \sum_{i=1}^{n} \mathcal{L}_i$$

$$\mathcal{L}_i = (y_i - \hat{y}_i)^2$$

$$\hat{y}_i = w^T x_i$$

Optimization problem:
Find Wi that minimizes the optimization problem.
F is identity in case of Linear regression.
F is Threshold in case of perception.
F is Sigmoid in case of Logistic regression.



Solve the optimization problem:

Updating weights -

$$\boxed{C} \quad W_{new} = W_{old} - \eta \left[ \nabla_w L \right]_{w_{old}}$$

$$\left\{ (w_i)_{new} = (w_i)_{old} - \eta \left[ \frac{\partial L}{\partial w_i} \right]_{(w_i)_{old}} \right.$$
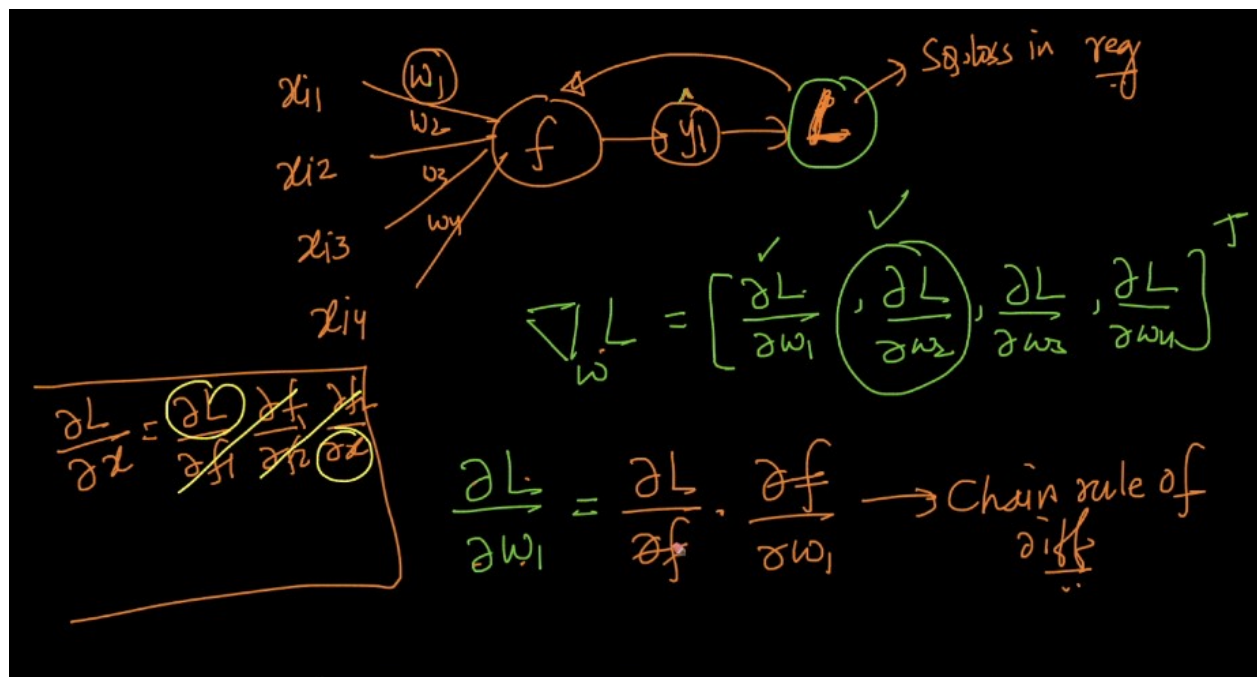
$$for \ iter = 1 \ to \ k$$

The key part is computing gradient descent we will use all the data for computing the new weights. In SGD we will use the small batch of points and approximate the weight vector.

$$GD: \quad \nabla_w L \quad \rightarrow \quad x_i's \ \& \ y_i's$$

$$SGD: \left\{ \left( \nabla_w L \right) \approx one \ pt \ \{x_i, y_i\} \right.$$
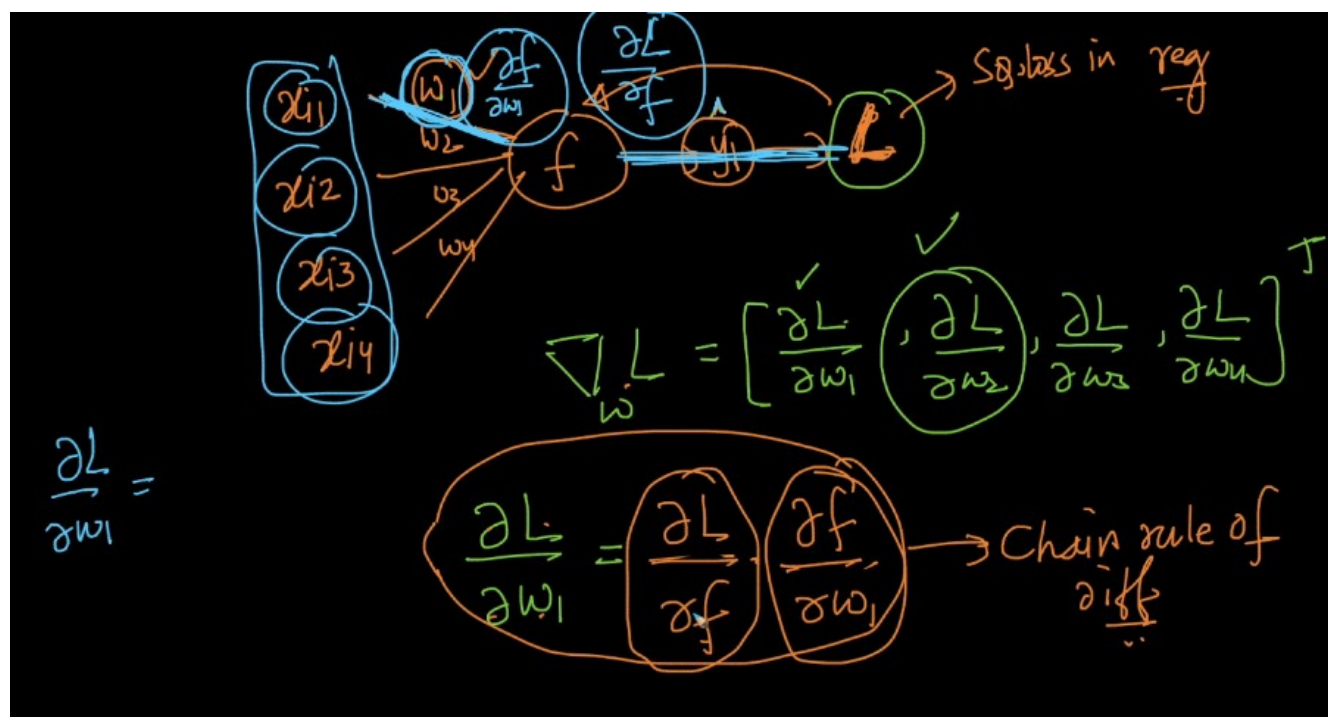
Small batch of pts → batch SGD

In regression we have squared loss for the network.



These are not ratios, they are derivatives :P.

Pictorial representation of derivatives:

$$f(w^T x_i) = w^T x_i$$

$$L = \sum_{i=1}^{n} (y_i - \hat{y_i})^2 + \cancel{\lambda g}$$

$$= \sum_{i=1}^{n} (y_i - f \quad)^2$$

$$\frac{\partial L}{\partial f} = -\sum_{i=1}^{n} 2(y_i - f(w^T x_i))$$

$$\frac{\partial f}{\partial w_1} = \boxed{x_i}$$

Training MLP : Chain Rule