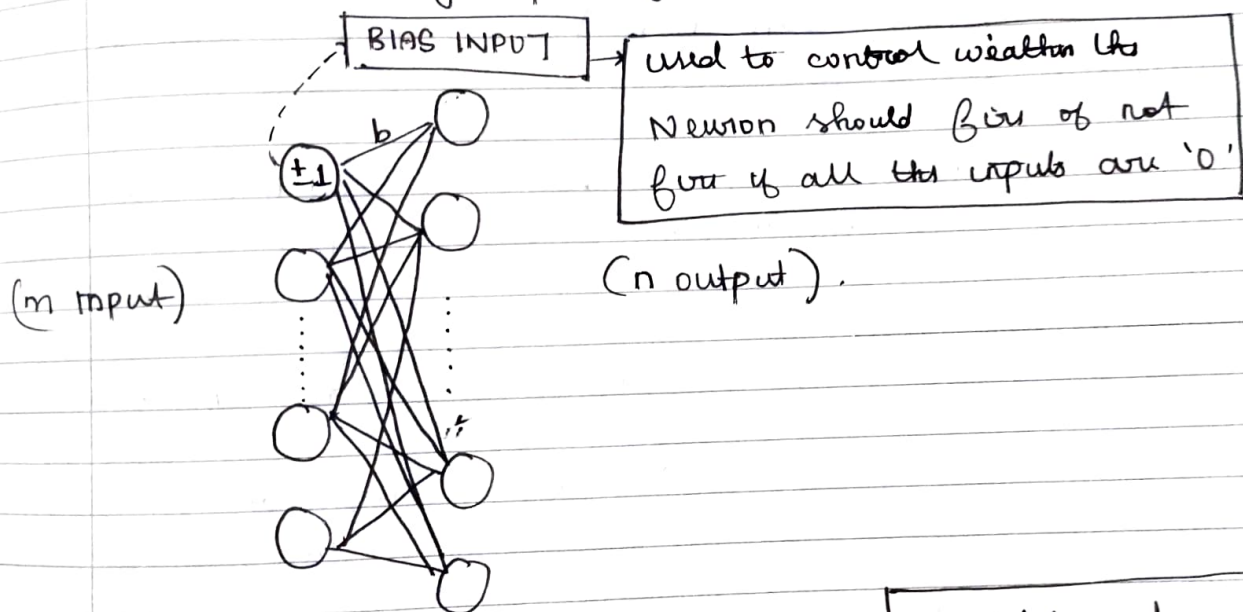# * PERCEPTRON : (SINGLE LAYER PERCEPTRON)

→ A linear model used for binary Classification

→ Step Function for activation i.e $f(x) = \begin{cases} 0, & x \le \alpha \\ 1, & x > \alpha \end{cases}$

$\qquad\qquad\qquad\qquad \alpha \Rightarrow$ THRESHOLD VALUE.

→ Contains only input layer and output layer.

| BIAS INPUT | used to control weather the Neuron should fire or not fire if all the inputs are '0' |

(m input)

(n output).

At each output neuron K we have,

| $x_D = \pm 1$ and $W_{0K} = b$ which is weight of bias |

$$\boxed{\hat{y}_K = f\left(\sum_{i=D}^{m} x_i \, w_{iK}\right)} \Rightarrow \text{Output of Kth neuron.}$$

⇒ The weights of the connections are update through. Perceptron learning algorithm.

At the $K^{th}$ output neuron the error $E_K$ is Given by

$$\boxed{E_K = \hat{y}_K - y_K} \quad \text{where } \hat{y}_K = \text{predicted value}$$

$\qquad\qquad\qquad\qquad y_K = $ actual value.

if $E_k > 0$ ⟹ weights should be reduced

if $E_k < 0$ ⟹ weights should be increased

⟹ $$\Delta w_{ik} = -x_i (\hat{y}_k - y_k)$$

Therefore the weight updation is given by.

$$w_{ik} = w_{ik} + \eta \Delta w_{ik} = w_{ik} - \eta \cdot x_i (\hat{y}_k - y_k)$$

The parameter $\eta$ is called as learning Rate. It determines show fast should the perceptron learn by controlling the amount of change is weights

if $\eta$ is too high ⟹ network becomes unstable

if $\eta$ is too low ⟹ learning will take large amount of time

$$0.1 < \eta < 0.4 \Rightarrow \text{optimal Range}$$

* MULTILAYER PERCEPTRON (MULTILAYER FEED FORWARD NETWORK):

→ Neuros are arranged into groups called as layer

→ Contains 3 main layers: ① Input layer
② Hidden layers.
③ Output layer

⇒ One a given layer all the Neurons will have Same activation function

⇒ For input layer, input will be the actual data neeto For other layers, input will be activation of previous layer

⇒ Number of neurons in the input layer will be equal to number of features in the input.

↝ we have two main Steps:

    ① Forward Propogation (Predicting Output)

    ② Backward Propogation (Reducing Error by (Back propogatt of error) tuning weight and bias value).

⇒ each neuron can be given a Seperate bias value (or) each layer can have a common bias value, This is Based on the implementation

⇒ BACKPROPOGATION: it is the method of calculating and propagation error from output layer to all Hidden. layer.

    → To reduce the error we using "optimisation algorithm" ⇓ Ⓧ

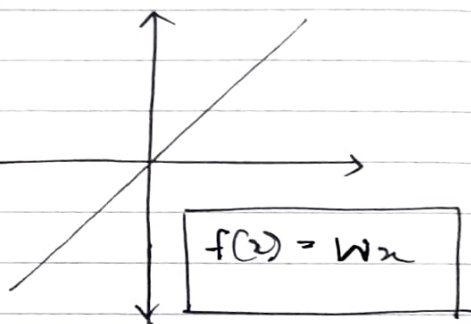| |
|---|
| error is calculated using "loss function" ⇒ The goal of optimiser to to is minimise the loss function |

linear    Non-linear

⇒ ACTIVATION FUNCTION : used to introduce non-linearity in the decision boundary (Because not all. problems are linearly Seperable)

↳ without ~~activation~~ activation function the neural network will just performs linear weighted Sum (Similar to linear Regression).

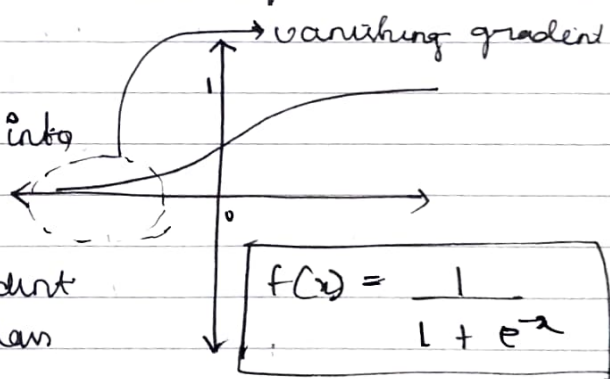① LINEAR :  $f(x) = Wx$.

⇒ identity function
⇒ used in input layer ?

$$f(x) = Wx$$

② SIGMOID :

→ vanishing gradient

⇒ Reduce the input into the Range  $0 - 1$.

⇒ Outputs independent probability for each class
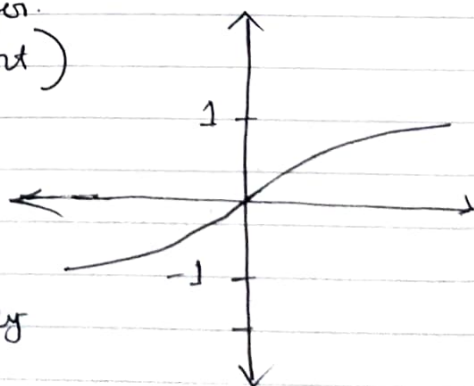
$$f(x) = \frac{1}{1 + e^{-x}}$$

~~....~~ :  ↳ und in binary classification output layer.

③ TANH : (Hyperbolic Tangent)

⇒ Reduces input to the Range −1 to 1

⇒ it can deal more easily with negative values.

$$f(x) = \frac{e^{x} - e^{-x}}{e^{x} + e^{-x}}$$

Both Sigmoid and tanh can cause vanishing Gradient problem (For extreme high / low value) ⟹ So we use ReLU.
But it can cause exploding Gradient
↳ need to Scale i/p

④ SOFTMAX :

→ provides the probability distribution of each class.
→ used in Output layer for multi-class classification

$f(x) = \dfrac{e^?}{\sum_i e^{x_i}}$  ⟹ if output layer contains M neurons and X be the output value/activation

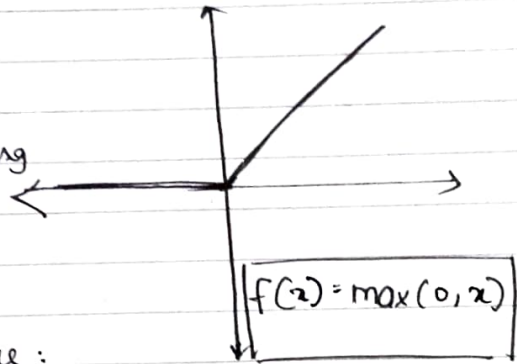⟹ $f(x_j) = \dfrac{e^{x_j}}{\sum_i^M e^{x_i}}$    (Sum of all probability will be 1)
(class with maximum probability will be the Result).

⑤ RECTIFIED LINEAR UNIT (ReLU) :

⟹ overcomes vanishing / exploding gradient problem.
⟹ ReLU is highly used in hidden layer
⟹ Some varient of ReLU are :

$f(x) = max(0, x)$

(i) Leaky ReLU = $f(x) = max(0.01x, x)$.

(ii) Parametric ReLU = $f(x) = max(\alpha x, x)$.

(iii) Exponential ReLU : $f(x) = \begin{cases} x, & \text{if } x \geq 0 \\ \alpha(e^x - 1), & \text{if } x < 0 \end{cases}$

⟶ $max(\alpha(e^x - 1), x)$.

LOSS FUNCTION: It quantifies how close the predicted value is to the Real value.

For Given weights W and Bias B we have L (W,B) as loss function

$$\hat{Y_i} = h_{W,B}(X) \Rightarrow \text{predicted value for } i\text{th Sample}$$
$$Y_i \Rightarrow \text{Real value for } i\text{th Sample}$$

| LOSS FUNCTION FOR REGRESSION | LOSS FUNCTION FOR CLASSIFICATION |
|---|---|
| ① Mean Squared Error | Cross entropy loss. (OR) |
| $$L(W,B) = \frac{1}{N} \sum_{i=1}^{N} (\hat{Y_i} - Y_i)^2$$ | Negative Log Likelihood. |
| if more than one feature is der as output | ↳ Binary ↳ Categorical. |
| $$\Rightarrow \frac{1}{N} \sum_{i=1}^{N} \left( \sqrt{\sum_{j=1}^{M}(\hat{y}_{ij} - y_{ij})^2} \right)^2$$ | $$BCE = -\frac{1}{N} \sum_{i=1}^{N} Y_i \log \hat{Y_i} + (1-Y_i) \log(1-\hat{Y_i})$$ |
| $$\Rightarrow \frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{M} (\hat{y}_{ij} - y_{ij})^2$$ | $$CCE = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{M} Y_{ij} \log(\hat{Y_{ij}})$$ |
| ② Mean Absolute Error. | Hing loss: |
| $$L(W,B) = \frac{1}{N} \sum_{i=1}^{N} |\hat{Y_i} - Y_i|$$ | ↳ designed for SVM |
| $$\Rightarrow \frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{M} |\hat{y}_{ij} - y_{ij}|$$ | Result should be -1 / 1. |
| ③ Root Mean Squared | $$L(W,B) \frac{\sum\limits_{i=}^{N} max(0, 1-Y_i \times \hat{Y})}{N}$$ |
| $$L(WB) = \sqrt{MSE}$$ | $$\Rightarrow \frac{1}{N} \sum_{i=1}^{N}$$ |

$$\sum_{i=1}^{N} \sum_{j=1}^{M} \max\left(0, y_{ij} \times \hat{y}_{ij}\right) \Rightarrow \text{only used for binary}$$
classification

MSE and MAE are prone to outliers due to which
we can use $\boxed{\text{HUBER LOSS}} \rightarrow$ For Regression

$$\hookrightarrow L(W, B) = \begin{cases} \frac{1}{2}(Y_i - \hat{Y}_i)^2, & |Y_i - \hat{Y}_i| \leq \delta \\ \delta|Y_i - \hat{Y}_i| - \frac{1}{2}\delta^2, & \text{else.} \end{cases}$$

Here $\delta$ is a hyperparameter which depicts the
tolerance to outliers.

> Huber loss is like a
> Combination MSE and MAE

\* **GRADIENT DESCENT :**

→ An optimisation algorithm for minimising the
cost / loss function

    (i) Batch Gradient Descent
    (ii) Stochastic Gradient Descent
    (iii) Mini Batch Gradient Descent

⇒ In (i) the complete Training Set is used for every iteration/
epoch.
⇒ In (ii) Only Single Sample is used for one down (i.e.
Number of iterations/epoch = Number of Samples in training
Set → minimum).
⇒ In (iii) A Batch of Samples used for each iteration
from the training dataset.

⇒ Some other optimisers are : RMS Props, ADADelta, ADAGrad,
Adam.

# * HYPERPARAMETERS:

⇒ Configuration options that can be varied in order to improve the performance of model

① No of Hidden Layers
② No of Neurons in each layer
③ Activation Functions
④ Loss Functions
⑤ Optimiser Function
⑥ No of ephochs / Iterations
⑦ Type of Regularisation
⑧ Weight initialisation
⑨ Learning Rate & Momentum.


# * BIAS AND VARIANCE: (Trade off).

⇒ when a model is trained it captures the patterns in the training data and try to generalise it.

⇒ BIAS is the measure of our models inability to learn patterns from data set.

⇒ when model has HIGH BIAS ⇒ did not train well in training Set

when model has LOW BIAS ⇒ it has Identified new patterns from Traing Set properly

⇒ VARIANCE is the measure of models fluctuation in Performance when there is a change in data set.

⇒ UNDERFITTING: when a model does not perform well in training data.
(HIGH BIAS).

⇒ OVERFITTING: when a model performs well in train data but it lacks performance in testing data
(LOW. BIAS, HIGH VARIANCE)

⇒ A Good model should have LOW BIAS and LOW VARIANCE.

MACHINE LEARNING: A Algorithm is said to be learning from experience E with Respect to some task T and performance measure P if the performance. P. increases with increase in experience E.

TYPES OF MACHINE LEARNING:

① SUPERVISED → CLASSIFICATION
→ REGRESSION

② UNSUPERVISED (→ CLUSTERING.).

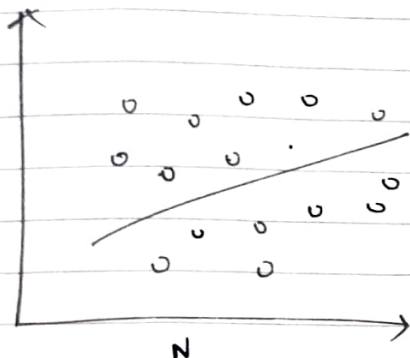③ SEMI-SUPERVISED

④ REINFORCEMENT.

LINEAR SEPERABILITY:

A data is said to be linearly Seperable if only a single line/Plane/Hyperplane can partition the data into 2 classes.

Eg: Perceptron can only work on a data Set that is linearly Seperable.

→ One of the main reason for MLP improvement

## LINEAR REGRESSION:

consider the following equation $y = \theta_0 x + \theta_1$



> Linear regression tries to fit a line of the form on the Given dataset.

$$E = \frac{1}{N} \sum_{i=1}^{N} (\hat{Y}_i - Y_i)^2 \Rightarrow \text{where } \hat{Y}_i \text{ is the predicted}$$

value and Y is the actual value

The error E is called as MEAN SQUARED ERROE. we have to Reduce this here this is called LEAST SQUARE REGRESSION.

$$E = \frac{1}{N} \sum_{i=1}^{N} (\hat{Y}_i - Y_i)^2 = \frac{1}{N} \sum_{i=1}^{N} (\theta_0 x_i + \theta_1 - Y_i)^2.$$

$$\frac{\partial E}{\partial \theta_0} = 2 \times \frac{1}{N} \times \sum_{i=1}^{N} (\theta_0 x_i + \theta_1 - Y_i)(x_i + 0 - 0) = 0$$

$$= \theta_0 \overline{x^2} + \theta_1 \overline{x} - \overline{x \cdot y} = 0 \quad —①$$

$$\frac{\partial E}{\partial \theta_1} = 2 \times \frac{1}{N} \times \sum_{i=1}^{N} (\theta_0 x_i + \theta_1 - Y_i)(0 + 1 - 0) = 0$$

$$= \theta_0 \bar{X} + \theta_1 - \bar{Y} = 0$$

$$\Rightarrow \theta_1 = \bar{Y} - \theta_0 \bar{X} \quad -\textcircled{2}$$
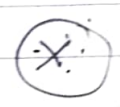
From ① and ② we have.

$$\theta_0 \bar{X^2} + (\bar{Y} - \theta_0 \bar{X}) \cdot \bar{X} - \overline{X \cdot Y} = 0$$

$$\theta_0 \bar{X^2} + \bar{X} \cdot \bar{Y} - \theta_0 \bar{X}^2 - \overline{X \cdot Y} = 0$$

$$\theta_0 = \frac{\overline{X \cdot Y} - \bar{X} \cdot \bar{Y}}{\bar{X^2} - \bar{X}^2}$$

$$\Rightarrow$$

$$\boxed{\theta_0 = \frac{\overline{X \cdot Y} - \bar{X} \cdot \bar{Y}}{\bar{X^2} - \bar{X}^2}}$$

$$\boxed{\theta_1 = \bar{Y} - \theta_0 \cdot \bar{X}}$$

$$\Rightarrow \text{why Huber loss for Regression ?}$$

MSE is Highly Sensitive to outliers because we Square the over all error values. MAE just take the absolute value of error it give same weightage to all errors due to which outlier can be completely Neglected. Huber Loss Provide Balance (S).

# ✳ BAYESIAN STATISTICS :

↳ (methods that used Bayes Theorm)

→ **BAYES THEORM** : Finds the probability of an event
occuring, given the probability of another event
that has already occured
↳ (called as conditional probability

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

⇒ Here, ·) P(B) is called as evidence.

·) P(A) is called as prior probability
(i·e, probability occurance of event A
without considering event B).

·) P(A|B) is called as postirior Probability
(i·e, probability of occurance of event A
with considering event B)

⇒ we can P(Y|X) using above method.

Y ⇒ Target Class and X = Feature Vector.

⇒ $P(Y_j | X_1, X_2, ..., X_M) = $ ~~P(X₁)P(X₂)...P(Xₘ)~~

$$= \frac{P(Y_i|X_1) \cdot P(Y_i|X_2) \cdots P(Y_i|X_M) \cdot P(Y_i)}{P(X_1) \cdot P(X_2) \cdots P(Y_M)}$$

↳ since evidence is same.

for all target we can ignore
this

⇒ Result will be class j having max $P(Y_j | X)$.

→ There are various variation of Naive Bayes algorithms like

      (i) Gaussian / Normal NB
      (ii) Multinomial NB
      (iii) Bernoulli NB
      ⋮

⇒ Important assumption for Naive Bayes algorithm is all the features are independent.

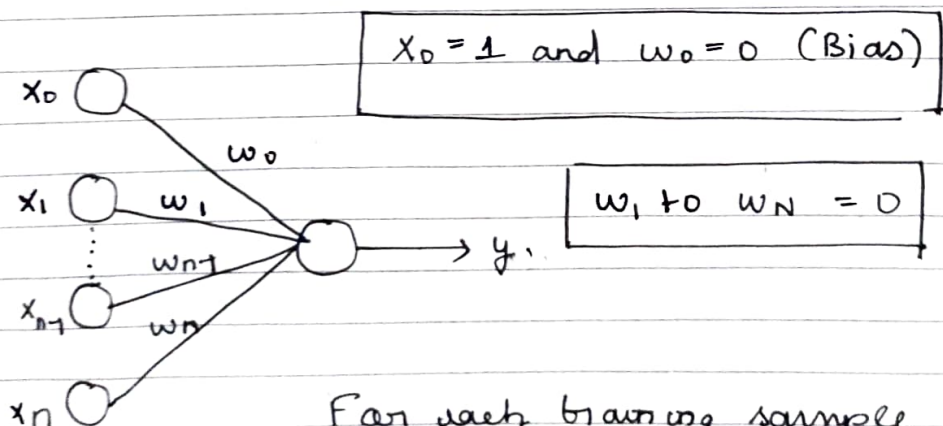$$\boxed{\text{i.e} \quad P(x_i | x_j) = P(x_i)} \quad \overset{.}{\underset{.}{\bigcirc}}X\overset{.}{\phantom{.}}$$

———— × ————

**\* HEBBIAN LEARNING RULE (HEBB NETWORK)**

→ inputs can only be 1 and −1

→ initialise all weights and bias to 0. in Starting

$$\boxed{X_0 = 1 \text{ and } W_0 = 0 \text{ (Bias)}}$$



$X_0$

$W_0$

$X_1$   $W_1$

$$\boxed{W_1 \text{ to } W_N = 0}$$

$W_{n-1}$

$X_{n-1}$   $W_n$   → $y$.

$X_n$

For each training sample we do the following

$\Rightarrow \quad \Delta w_i = x_i y \quad$ and $\quad w_i = w_i + \Delta w_i$

$$\underline{\qquad} \; \times \; \underline{\qquad}$$

# ✳ FORMULAS

① PERCEPTRON WEIGHT UPDATION:

$$W_{ij} = W_{ij} + \eta\,(\hat{y}_j - y_j)(x_i)$$

② LOSS FUNCTIONS:

(i) REGRESSION:

⇒ MEAN SQUARED ERROR:

$$E = \frac{\sum_{i=1}^{N}\sum_{j=1}^{M}(\hat{Y}_{ij} - Y_{ij})^2}{N} \quad (OR) \quad E = \frac{\sum_{i=1}^{N}(\hat{Y}_i - Y_i)^2}{N}$$

⇒ MEAN ABSOLUTE ERROR:

$$E = \frac{\sum_{i=1}^{N}\sum_{j=1}^{M}|\hat{Y}_{ij} - Y_{ij}|}{N} \quad (OR) \quad E = \frac{\sum_{i=1}^{N}(\hat{Y}_i - Y_i)}{N}$$

⇒ ROOT MEAN SQUARED ERROR

$$E = \sqrt{\frac{\sum_{i=1}^{N}\sum_{j=1}^{M}(\hat{Y}_{ij} - Y_{ij})^2}{N}} \quad (OR) \quad E = \sqrt{\frac{\sum_{i=1}^{N}(\hat{Y}_i - Y_i)^2}{N}}$$

⇒ $$\boxed{E_{RMSE} = \sqrt{E_{MSE}}}$$

⇒ HUBER LOSS (A Balance between MSE and MAE).

$$E = \begin{cases} \frac{1}{2}(y_i - \hat{y}_i)^2, & \text{if } |y_i - \hat{y}_i| \le \delta \\ \delta|y - \hat{y}_i| - \frac{1}{2}\delta^2, & \text{otherwise} \end{cases}$$

δ value controls the behaviour of
huber loss (Hyperparameter).

(ii) CLASSIFICATION :

⇒ HINGE LOSS :

⇒ CROSSENTROPY LOSS :

$$E = -\sum_{i=1}^{N} \sum_{j=1}^{M} Y_{ij} \log \hat{Y}_{ij}$$ ( Multiple Class ⇒ Categorical CE ).

$$E = -\sum_{i=1}^{N} Y_i \log \hat{Y}_i + (1-Y_i) \log (1-\hat{Y}_i).$$

( Binary Class ⇒ Binary CE ).

③ SOME ACTIVATION FUNCTION :

(i) Sigmoid :

$$g(x) = \frac{1}{1 + e^{-x}}$$

(ii) TanH :

$$g(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

(iii) ReLU :

$$g(x) = \max(0, x)$$

→ Parametric : $g(x) = \max(ax,$

→ Leaky : $g(x) = \max(0.01x, x)$

→ Exponential : $g(x) = \max(\alpha(e^x-1), x$

(iv) Softmax:

$$g(x_i) = \frac{e^{x_i}}{\sum\limits_{j=1}^{M} e^{x_j}}$$

(v) LINEAR:

$$g(x) = \alpha x.$$

④ NAIVE BAYES:

For M target classes ~~formark~~ we have to find max of

~~PCCooDoooe-PQx~~

$$P(Y_j \mid X_1, X_2, X_3, \ldots, X_K) = \frac{P(X_1 \mid Y_j) \cdot P(X_2 \mid Y_j) \cdots P(X_K \mid Y_j) \cdot P(}{\boxed{P(X_1) \cdot P(X_2) \cdots P(X_K)}}$$

↪ Same for all So can be ignored.

⑤ MULTILAYER PERCEPTRON :

~~FOOl PORROTO OR OOIOOORO~~

$$E_K = O_K (1 - O_K)(O_K - T_K). \quad \text{(Output layer)}$$

$$E_j = O_j (1 - O_j)\left(\sum\limits_{K=1}^{N} W_{jK} E_K\right). \quad \text{(Hidden layer)}.$$

$$W_{ij} = W_{ij} + \eta \, O_i E_j$$