**Name:** Janmejay Mohanty      **Course:** CS 583 A      **Quiz1A**
**CWID:** 20009315      **Course Name:** Deep Learning      **Date:** 11th October 2022
**Email:** jmohanty@stevens.edu

**Solutions**

**Ans1:**



$$w_1 = \begin{bmatrix} 1 & 1 \end{bmatrix} \qquad\qquad w_2 = \begin{bmatrix} 1 & 1 \end{bmatrix}$$

$$b_1 = \begin{bmatrix} -\frac{1}{2} & -\frac{1}{2} \end{bmatrix} \qquad\qquad b_2 = \frac{1}{2}$$

**Activation Function:** $f(h) = \begin{cases} 1, & if \ w_x + b > 0 \\ 0, & else \end{cases}$

| $X_1$ | $X_2$ | $h_1$ | $h_2$ | $y$ |
|-------|-------|-------|-------|-----|
| 1 | 1 | 1 | 1 | 1 |
| 1 | 0 | 1 | 1 | 1 |
| 0 | 1 | 1 | 1 | 1 |
| 0 | 0 | 0 | 0 | 0 |

**Name:** Janmejay Mohanty      **Course:** CS 583 A      **Quiz1A**
**CWID:** 20009315      **Course Name:** Deep Learning      **Date:** 11th October 2022
**Email:** jmohanty@stevens.edu

**Ans2:**

   **(1)**

$Pr(C) \rightarrow \boldsymbol{Prior\ Probability}$: The prior probability of an event is the probability of the event computed before the collection of new data. One begins with a prior probability of an event and revises it in the light of new data.

$Pr(X|C) \rightarrow \boldsymbol{Class\ Conditional\ Probability}$: Conditional probability is a measure of probability to an event that is occurring, given other event has already occurred.

$Pr(C|X) \rightarrow \boldsymbol{Discrimitive\ Model}$: Also referred as conditional models. This is a class of logistical models used for classification or regression. It should distinguish decision boundary through observed data.

$Pr(C).Pr(X|C) \rightarrow \boldsymbol{Generative\ Model}$: A generative model describes how a dataset is generated, in terms of a probabilistic model. By sampling this model, we generate data.

   **(2)**

Beam Size $k = 1$: $\langle BOS \rangle$ Montreal a great playground.

Beam Size $k = 2$: $\langle BOS \rangle$ Montreal, giant playground.

$\boldsymbol{Pros}$: More words would be traversed and chosen. Therefore, a higher probability to get better results.

$\boldsymbol{Cons}$: More computational resources and more memory are required.


**Ans3:**

   **(1) Trigram model of language modelling:** A trigram model restricts the conditional information to the previous two words. Using this method, the conditional distribution can calculate a certain word combination frequency based on the previous two words.
   **(2) Procedure of 5-fold cross-validation:** Dataset is split into 5 sets. Each one set is taken as test set by turn while other four are training sets. Eventually, we'll have 5 accuracies and the average is the accuracy of 5-fold cross-validation.
      **Pros:** Avoiding the randomness and bias by training and testing all the data.
   **(3) Bagging:** Bootstraps the training set, estimates many copies of a model on the resulting samples and then averages their predictions.
      **Boosting:** Sequentially reweights the training samples forcing the model to attend the training examples with higher loss.
      **Stacking:** Used a separate validation set to train a meta-model that combines predictions of multiple models.