# Multimodal Deep Learning

Purvanshi Mehta
purvanshi.mehta11@gmail.com
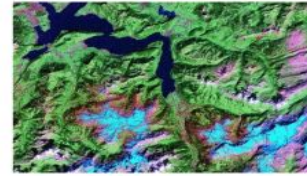
# Table of contents

- Introduction
- Challenges
  - Representation
  - Translation
  - Alignment
  - Fusion
  - Co-Learning
- Interpretability in Multimodal Deep Learning

# Aim of the presentation

- Identify challenges particular to Multimodal Learning
- Popular research topics in the field
- Brief of the problem I have been working on - Interpretability in Multimodal Deep Learning

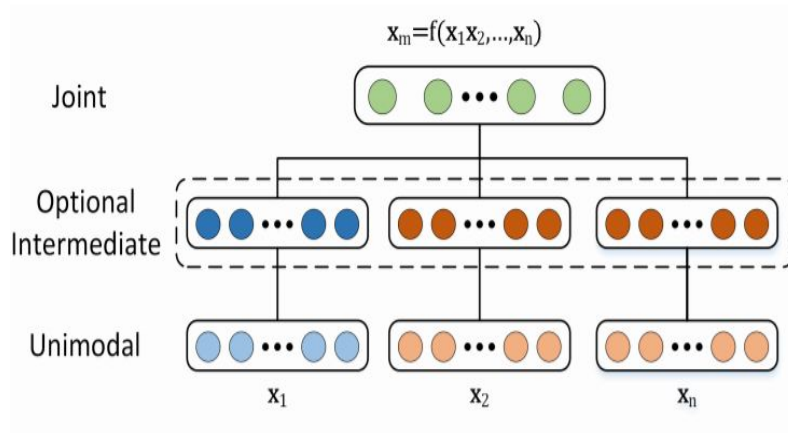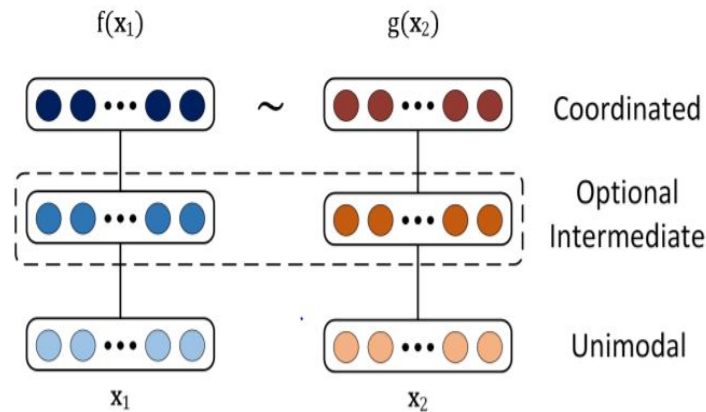# Multimodal Learning - Heterogeneous Information Sources

# Challenge - 1) Representation

- How to combine the information from multiple sources?
- How to deal with different levels of noise?
- How to deal with missing data?

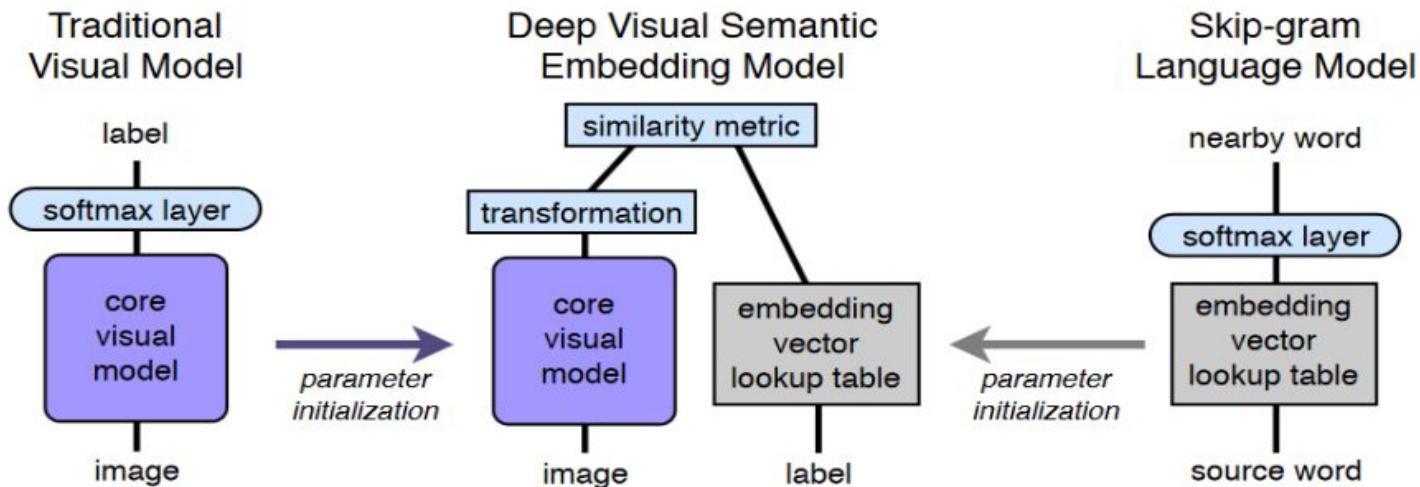# 1) Representation - Ways of learning



(a) Joint representation

(b) Coordinated representations

# Coordinated Representation

**DeViSE — a deep visual-semantic embedding - Similarity model**

# Challenge - 2) Translation

{...a **multicolored table** in the middle of the room...,
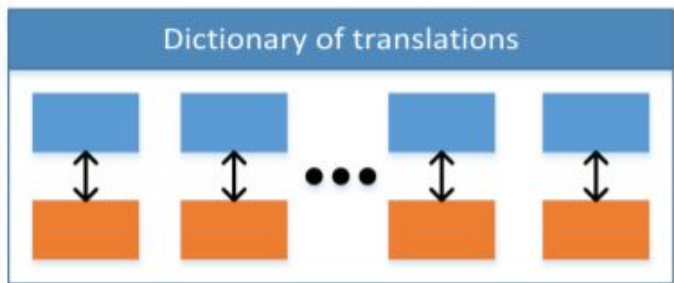...four red and white chairs and a **colorful table**, ...}

{...**L-shaped room** with walls that have 2 tones of gray...,
A **dark room** with a pool table...}

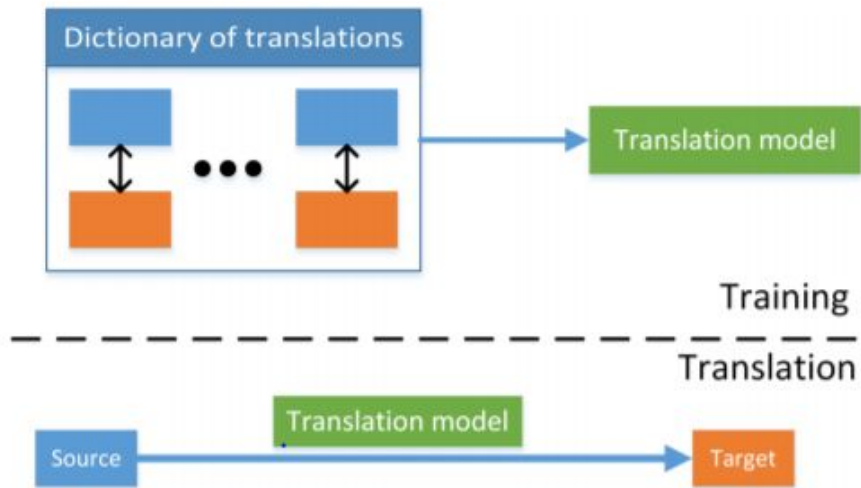# Translation



(a) Example-based

(b) Generative

# How to evaluate translations



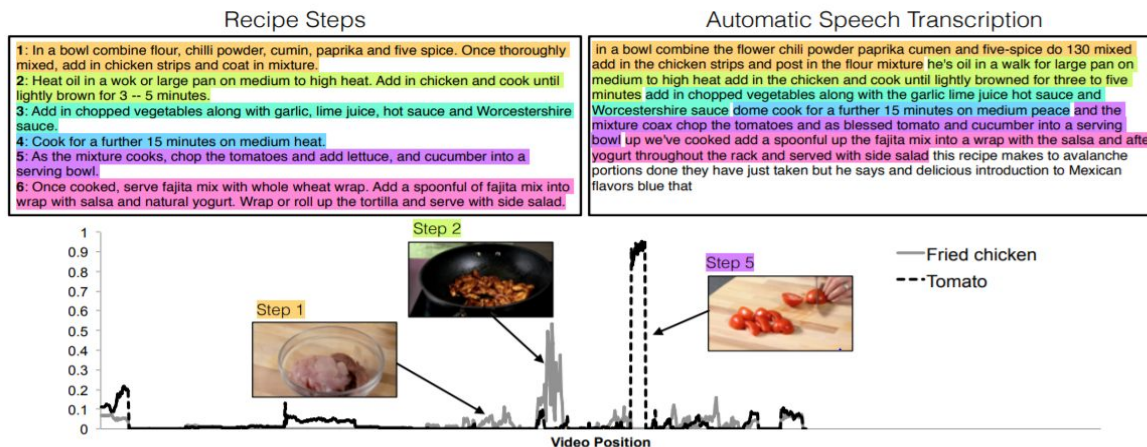**Candidate**: Football players gathering to contest something to collaborating officials.
**Reference**: A football player in red and white is holding both hands up.

# Challenge - 3) Alignment

1) Identify the direct relations between (sub)elements from two or more different modalities.

# Alignment

Given an image and a caption we want to find the areas of the image corresponding to the caption's words or phrases

# Alignement problems

- Few datasets with explicitly annotated alignments
- It is difficult to design similarity metrics between modalities
- Existence of multiple possible alignments and not all elements in one modality have correspondences in another

# Challenge - 4) Fusion



Early Multimodal Fusion     Intermediate Multimodal Fusion     Dense Multimodal Fusion

x modality    y modality     x modality    y modality     x modality    y modality

# Fusion

- Signals not temporarily aligned
- Lack of interpretability of where is the prediction coming from
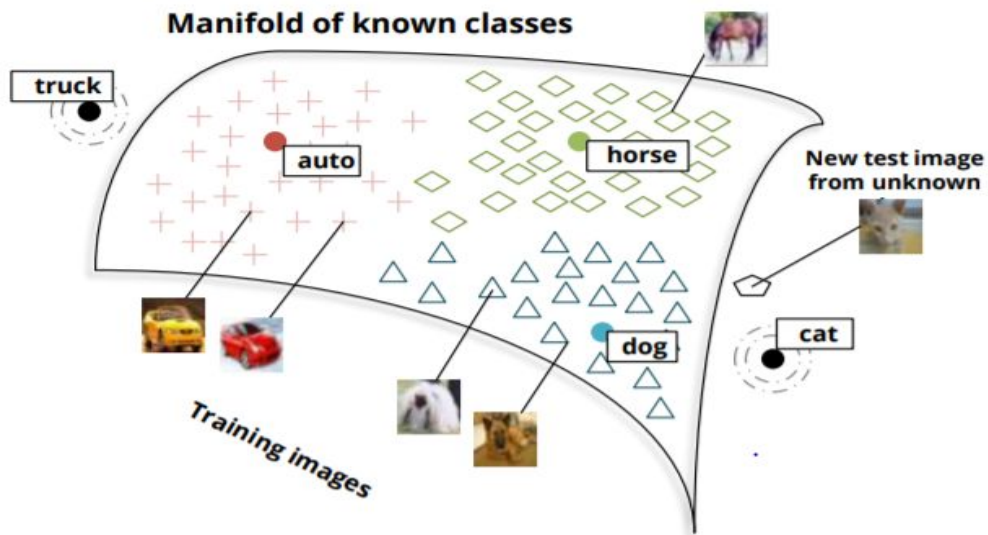  - (This is what I am working on)

# Challenge - 5) Co-Learning

- Aiding the modeling of a (resource poor) modality by exploiting knowledge from another (resource rich) modality.
- When one modality has lack of annotated data, noisy inputs and unreliable labels.

# Colearning - Zero Shot learning

Using text embeddings to classify unseen classes of images

# Interpretability in Multimodal Deep Learning

Problem statement -

Not every modality has <span style="color:red">equal contribution</span> to the prediction

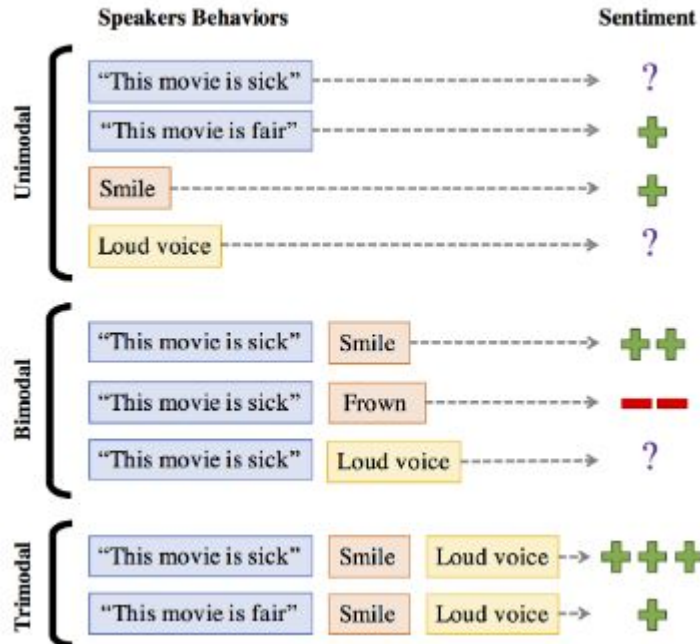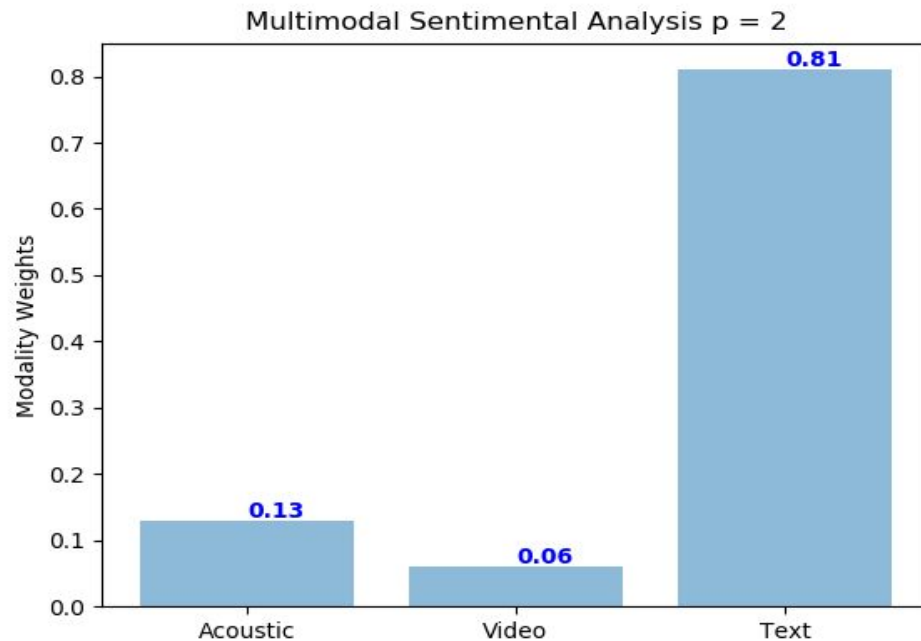# Interpretability in Multimodal Deep Learning

Solution -

We give different weights to different modalities

# Real data experiments - Multimodal Sentiment Analysis(MOSI Dataset)

# MUltimodal Sentiment Analysis



Multimodal Sentimental Analysis p = 2

# P1 and P2 contribution of each modality



$$N_p = \sum_{m=1}^{M} p_m < f^m_{w^m_1, w^m_2 .. w^m_{L-1}}(x), W^m_L > + b$$

Θ

Θ

# Modified loss function with $\beta$ weight given to each modality

$$\min_{\substack{w_1, w_2, \ldots, w_L, \beta \\ \beta \in \mathbb{R}^M, \beta \geq 0, \|\beta\|_p \leq 1}} \left( \sum_{i=1}^{n} \ell\left( \sum_{m=1}^{M} \sqrt{\beta_m} \langle w_L^m, f^m(x_i^m) \rangle + b, y_i \right) + \Lambda \sum_{l=1}^{L} \sum_{m=1}^{M} \|w_l^m\|_2^2 \right)$$

# Modified loss function with $\beta$ weight given to each modality

$$\min_{\substack{w_1,w_2,\ldots,w_L,\beta \\ \beta \in \mathbf{R}^M, \beta \geq 0, \|\beta\|_p \leq 1}} \left( \sum_{i=1}^{n} \ell \left( \sum_{m=1}^{M} \sqrt{\beta_m} \langle w_L^m, f^m(x_i^m) \rangle + b, y_i \right) + \Lambda \sum_{l=1}^{L} \sum_{m=1}^{M} \|w_l^m\|_2^2 \right)$$

$$w_L^m \leftarrow \sqrt{\beta_m} w_L^m, \qquad \text{Let's make } \boldsymbol{\beta} \text{ trainable}$$

$$\min_{\substack{w_1,w_2,\ldots,w_L,\beta \\ \beta \in \mathbf{R}^M, \beta \geq 0, \|\beta\|_p \leq 1}} \sum_{i=1}^{n} \ell \left( \sum_{m=1}^{M} \langle w_L^m, f_{w_1^m,w_2^m,\ldots,w_{L-1}^m}^m (x_i^m) \rangle + b, y_i \right)$$

$$+ \Lambda \sum_{l=1}^{L-1} \sum_{m=1}^{M} \|w_l^m\|_2^2 + \Lambda \sum_{m=1}^{M} \frac{\|w_L^m\|_2^2}{\beta_m}.$$

# Modified loss function with $\beta$ weight given to each modality

$$\min_{\substack{w_1, w_2, \ldots, w_L, \beta \\ \in \mathbf{R}^M, \beta \geq 0, \|\beta\|_p \leq 1}} \sum_{i=1}^{n} \ell \left( \sum_{m=1}^{M} \langle w_L^m, f_{w_1^m, w_2^m, \ldots, w_{L-1}^m}^m (x_i^m) \rangle + b, y_i \right)$$

$$+ \Lambda \sum_{l=1}^{L-1} \sum_{m=1}^{M} \|w_l^m\|_2^2 + \Lambda \sum_{m=1}^{M} \frac{\|w_L^m\|_2^2}{\beta_m}.$$

Using the Lemma

$$\min_{\beta > 0, \|\beta\|_p^p \leq 1} \sum_{m=1}^{M} \frac{A_m}{\beta_m} = \left( \sum_{m=1}^{M} A_m^{\frac{p}{p+1}} \right)^{\frac{p+1}{p}}$$

$$\min_{w_1, w_2, \ldots, w_L} \sum_{i=1}^{n} \ell \left( \sum_{m=1}^{M} \langle w_L^m, f^m(x_i^m) \rangle + b, y_i \right)$$

$$+ \Lambda \sum_{l=1}^{L-1} \sum_{m=1}^{M} \|w_l^m\|_2^2 + \Lambda \left( \sum_{m=1}^{M} \|w_L^m\|_2^q \right)^{\frac{2}{q}},$$

# Final Optimization problem

$$\min_{\substack{w_1,w_2,\ldots,w_L,\beta \\ \beta\in\mathbb{R}^M,\beta\geq 0,\|\beta\|_p\leq 1}} \left( \sum_{i=1}^{n} \ell\left( \sum_{m=1}^{M} \sqrt{\beta_m}\langle w_L^m, f^m(x_i^m)\rangle + b, y_i \right) + \Lambda \sum_{l=1}^{L}\sum_{m=1}^{M} \|w_l^m\|_2^2 \right)$$

$$\min_{w_1,w_2,\ldots,w_L} \sum_{i=1}^{n} \ell\left( \sum_{m=1}^{M} \langle w_L^m, f^m(x_i^m)\rangle + b, y_i \right)$$
$$+\Lambda \sum_{l=1}^{L-1}\sum_{m=1}^{M} \|w_l^m\|_2^2 + \Lambda \left( \sum_{m=1}^{M} \|w_L^m\|_2^q \right)^{\frac{2}{q}},$$

$$\beta_m = \frac{\|w_l^m\|_2^{\frac{2}{p+1}}}{\left( \sum_{\tilde{m}=1}^{M} \|w_l^{\tilde{m}}\|_2^{\frac{2p}{p+1}} \right)^{\frac{1}{p}}}$$
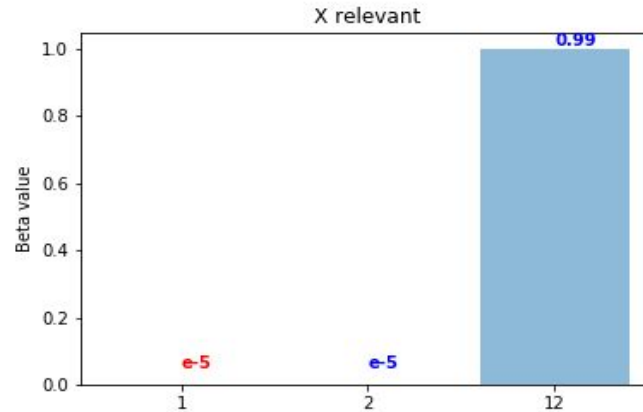
# Tensor fusion

To allow interaction between modalities - take tensor product between modalities.

$$
\begin{bmatrix} a & b \\ c & d \end{bmatrix} \otimes \begin{bmatrix} a^* & c^* \\ b^* & d^* \end{bmatrix} = \begin{bmatrix} a\begin{bmatrix} a^* & c^* \\ b^* & d^* \end{bmatrix} & b\begin{bmatrix} a^* & c^* \\ b^* & d^* \end{bmatrix} \\ c\begin{bmatrix} a^* & c^* \\ b^* & d^* \end{bmatrix} & d\begin{bmatrix} a^* & c^* \\ b^* & d^* \end{bmatrix} \end{bmatrix}
$$

$$
= \begin{bmatrix} aa^* & ac^* & ba^* & bc^* \\ ab^* & ad^* & bb^* & bd^* \\ ca^* & cc^* & da^* & dc^* \\ cb^* & cd^* & db^* & dd^* \end{bmatrix}
$$

# Problem with interpretability and tensorfusion - Degree Inflation

Beta weights came to be higher for higher dimension weights - Always high for tensor fusion.

# Degree Inflation - reason

When one modality learns constant features.

The information from the first modality can be

expressed in the tensorfusion.

Relevant source-with information

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \otimes \begin{bmatrix} a^* & c^* \\ b^* & d^* \end{bmatrix} = \begin{bmatrix} a\begin{bmatrix} a^* & c^* \\ b^* & d^* \end{bmatrix} & b\begin{bmatrix} a^* & c^* \\ b^* & d^* \end{bmatrix} \\ c\begin{bmatrix} a^* & c^* \\ b^* & d^* \end{bmatrix} & d\begin{bmatrix} a^* & c^* \\ b^* & d^* \end{bmatrix} \end{bmatrix}$$

$$= \begin{bmatrix} aa^* & ac^* & ba^* & bc^* \\ ab^* & ad^* & bb^* & bd^* \\ ca^* & cc^* & da^* & dc^* \\ cb^* & cd^* & db^* & dd^* \end{bmatrix}$$

Irrelevant source-constant

# Iterative batch normalisation

We derived a batch norm which does not allow the lower dimension information to be represented in higher dimension(TensorFusion).

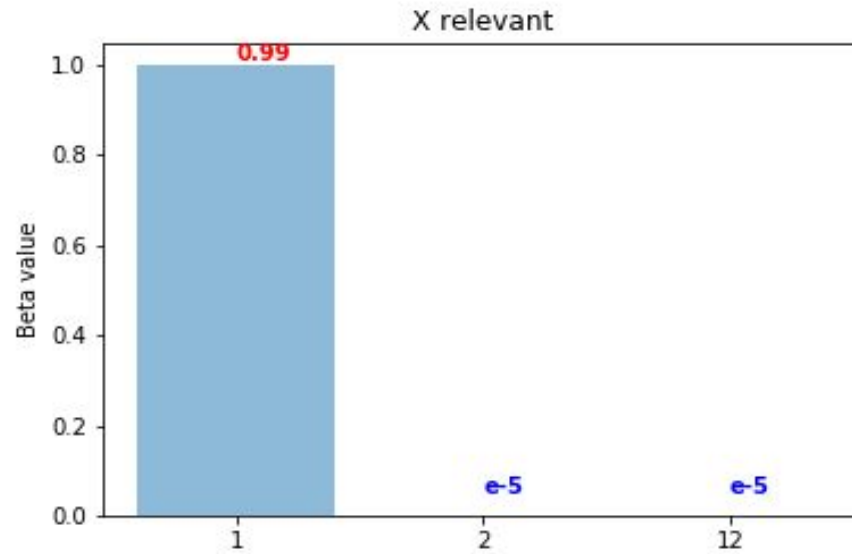This method helps to remove noise from the data and give better accuracies.

# **Synthetic data experiment**

- Relevant data - Sequences of length 100 (composed of 4 letters - ATGC) with a signature which defines the label

- Irrelevant data - Random 100 length sequences

# Example - 1 as relevant source

# Thank You