# LEVERAGING KULLBACK-DIVERGENCE LOSS FOR ENHANCED MACHINE LEARNING PERFORMANCE

A FOURTH YEAR PROJECT REPORT

SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE DEGREE OF B.Sc. IN COMPUTATIONAL MATHEMATICS

BY

1. Apekshya Shrestha(026612-19)

2. Abhay Sharma(026610-19)

3. Bishownath Raut(026607-19)



DEPARTMENT OF MATHEMATICS

SCHOOL OF SCIENCE

KATHMANDU UNIVERSITY

DHULIKHEL, NEPAL

Janurary 2024

# CERTIFICATION

This project entitled "LEVERAGING KULLBACK-DIVERGENCE LOSS FOR EN-HANCED MACHINE LEARNING PERFORMANCE" is carried out under my supervision for the specified entire period satisfactorily, and is hereby certified as a work done by following students

1. Apekshya Shrestha(apekshyashrestha12@gmail.com)

2. Abhay Sharma(abhaysharma2786@gmail.com)

3. Bishownath Raut(raut.bishow4@gmail.com)

in partial fulfillment of the requirements for the degree of B.Sc. in Computational Mathematics, Department of Mathematics, Kathmandu University, Dhulikhel, Nepal.

...........................

**Prof.Dr.Kanhaiya Jha**

Department of Mathematics

School of Science, Kathmandu University

Dhulikhel, Kavre, Nepal

Date: January, 2024

**APPROVED BY:**

I hereby declare that the candidate qualifies to submit this report of the Mathematics Seminar (**MATH 451**) to the Department of Mathematics.

...........................

Head of the Department

Department of Mathematics

School of Science

Kathmandu University

Date: January, 2024

# ACKNOWLEDGMENTS

# PROBLEM STATEMENT

Kullback-Leibler divergence class or relative entropy is an exceptional instance of a more extensive divergence. It is an estimation of how a particular dissemination wanders from another, normal likelihood appropriation.This report seeks to gain an intuitive understanding of KL divergence, an important concept in machine learning and statistics. This report provides a comprehensive yet accessible introduction to KL divergence, also known as relative entropy. Key aspects covered include the definition and mathematical properties of KL divergence, visualizations illustrating asymmetric behavior, code examples demonstrating calculation, and discussion of common applications such as comparing distributions and evaluating generative models. Through an in-depth analysis and experimental implementation, this report aims to build applied intuition around KL divergence. The concepts, examples, and visualizations provided in the source material will be leveraged to ensure a thorough practical understanding. Upon completion, the project will have tangibly demonstrated KL divergence calculation, interpretation, and usage for comparing probability distributions in fields like data science and machine learning.

**Keywords:** KL divergence,probability distributions, statistical distance measures, visualization,machine learning model

# CONTENTS

# LIST OF FIGURES

# CHAPTER 1

# INTRODUCTION

## 1.1 Background

In the field of information theory, the quantification of information and uncertainty is a fundamental pursuit. Information theory provides a framework for understanding and measuring the amount of information inherent in a message or a dataset, with applications spanning statistics, machine learning, and data compression.

At the core of information theory lies the concept of entropy, a metric that gauges the uncertainty or randomness within a dataset. Entropy enables the quantification of the expected value of information contained in a message, offering insights into the average amount of "surprise" or "information" received when learning the outcome of a random variable.

While entropy is a powerful measure, it may fall short in capturing the distinctions between probability distributions or the amount of information gained when transitioning from one distribution to another. This is where Kullback-Leibler (KL) Divergence emerges as a crucial concept[7, 9].

Named after Solomon Kullback and Richard Leibler, KL Divergence serves as a measure of dissimilarity between two probability distributions. It provides a means to quantify how one distribution diverges from another, finding applications in information theory, statistics, and machine learning for comparing and analyzing probability distributions.

The Kullback-Leibler (KL) Divergence is a measure of how one probability distribution diverges from a second, expected probability distribution. It's a way of comparing two probability distributions — often a "true" distribution and an approximation of that distribution. It quantifies the "distance" between these two distributions, or how much

information is lost when one distribution is used to approximate the other[7].

In many fields such as machine learning, statistics, and information theory, we often need to approximate complex distributions with simpler ones for computational efficiency or because we don't know the complex distributions. KL Divergence helps us quantify the trade-off we're having with this approximation.

## 1.2 Historical Development

The concept of KL divergence was first introduced in 1951 by Solomon Kullback and Richard Leibler in their paper "On Information and Sufficiency"[7]. It is also referred to as Kullback-Leibler divergence or relative entropy.

KL divergence quantifies the difference between two probability distributions - a true distribution P and an approximation Q. It measures the information lost when Q is used to represent P.

In the 1960s, KL divergence became more widely studied in the field of information theory pioneered by Claude Shannon[9]. KL divergence was found to be closely related to mutual information and entropy.

Some key related works:

- In 1960, Fréchet further developed the mathematical foundations of KL divergence and statistical manifolds[4].

- In 1962, Csiszár studied generalized convergence and probabilities using KL divergence[3].

- In 1975, Jeffreys divergence was introduced as a symmetric version of KL divergence[6].

- In 1985, Pardo introduced the concept of statistical divergence as a broader framework for comparing distributions[1].

- In the 1990s and 2000s, KL divergence became widely used in machine learning for training generative models, clustering, dimensionality reduction, and more[5].

Today, KL divergence remains a fundamental technique in information theory, statistics, and machine learning. It has been further extended into variants like -divergence, Rènyi divergence, and JS divergence.

## 1.3 Basic Defination with examples

Let's review some concepts before discussing KL Divergence.

1. **Probability Distribution:**

   A probability distribution is a mathematical function that provides the probabilities of occurrence of different possible outcomes in an experiment. For example, in a simple coin toss, the probability distribution of getting heads (H) or tails (T) is P(H) = 0.5, P(T) = 0.5.

2. **Entropy:**

   In the context of information theory, entropy is a measure of the uncertainty in a random variable. In other words, it measures the "surprise" you expect to have on average when you learn the outcome of the variable. The entropy of a discrete random variable X with probability mass function p(x) is defined as:

$$(X) = -\sum [p(x) log(p(x))]$$

[7]where the sum is over all possible outcomes of X, and log is the natural logarithm.

### 1.3.1 KL Divergence

The formula for KL Divergence for discrete probability distributions P and Q is defined as:

$$D_{\mathrm{KL}}(P||Q) = \sum P(X) \log(\frac{P(X)}{Q(X)})$$

[7]*where the sum is over all possible outcomes.*

The KL Divergence of $Q$ from $P$, denoted by $D_{\mathrm{KL}}(P||Q)$,is a measureof theimformation lost when $Q$ is used to approximate $P$.

In the above formula, P is the **true distribution** and Q is the **estimated distribution**.

### 1.3.2   Formula breakdown

1. The probability of event x according to the first distribution.

$$P(x)$$

   This term is used as a weighting factor, meaning events that are more probable in the first distribution have a larger impact on the divergence.

   We use this term because we're interested in the difference between P and Q where P has assigned more probability. If an event is highly probable in P but not in Q, we want that to contribute more to our divergence measure. Conversely, an event is highly improbable in P, we don't want it much to our divergence measure, even if Q assigns it a high probability. This is because we're measuring the divergence from P to Q, not the other way around[7].

   As P(x) will only give more weight to the events that are more probable in P, that is because we essentially asking "how different is Q from P?" with a bias towards the events that P thinks are more likely. This is useful in many scenarios where we want to penalize models (Q) that fail to assign high probabilities to events that are likely under the true data generating process (P).

2. This term is the ratio of the probabilities assigned to event x by P and Q.

$$\log(\frac{P(X)}{Q(X)})$$

   If P and Q assign the same probability to x, then this ratio is 1, and the logarithm of 1 is 0, so events that P and Q agree on don't contribute to the divergence. If P assigns more probability to x than Q does, then this ratio is greater than 1, and the logarithm is positive, so this event contributes to the divergence. If P assigns less probability to x than Q does, then this ratio is less than 1, and the logarithm is negative, but remember that we're multiplying this by P(x), so events that P assigns low probability to don't contribute much to the divergence.

3. For each outcome, we calculate how much probability P assigns to it, and then multiply it by the log of the ratio of the probabilities P and Q assign to it. This ratio tells us how much P and Q differ on this particular outcome.

$$P(X)\log(\frac{P(X)}{Q(X)})$$

4. We then sum over all possible outcomes. This gives us a single number that represents the total difference between P and Q[7].

$$\sum P(X) \log(\frac{P(X)}{Q(X)})$$

This means that the KL Divergence is a weighted sum of the log difference in probabilities, where the weights are the probabilities according to the first distribution.

### 1.3.3   Properties of KL divergence

The KL Divergence has the following properties:

1. Non-negativity:

$$D_{\text{KL}}(P||Q)$$

is always greater than or equal to 0. This is known as Gibbs' inequality. The KL Divergence is 0 if and only if P and Q are the same distribution, in which case there is no information loss[7].

2. Not Symmetric:
It is not symmetric, meaning that $D_{\text{KL}}(P||Q)$ is not same as $D_{\text{KL}}(Q||P)$
This is because the KL Divergence measures the information loss when Q is used to approximate P, not the other way around[7].

## 1.4   Objectives

The main objective of our projects is to compare two probability distributions - a true distribution P and an approximation Q [2]. More specifically, some key aims are:

- Investigate the mathematical foundation of KL Divergence and its significance in measuring divergence between probability distributions.

- Examine the diverse applications of KL Divergence in machine learning, including classification, generative modeling, regularization, and information retrieval.

- Analyze the impact of incorporating KL Divergence loss in machine learning models and how it influences model convergence, generalization, and accuracy.

- Explore methodologies for implementing KL Divergence loss in various machine learning frameworks and architectures.

## 1.5  Limitations

As previously established, KL divergence is an asymmetric metric. This means that it can not be used as strictly a distance measure since the distance between two entities remains the same from either perspective[5].

Moreover, if the data samples are pulled from distributions that use different parameters (mean and variance), KL divergence will not yield reliable results. In this case, one of the distributions needs to be adjusted to match the other.

# CHAPTER 2

# MATHEMATICAL RESULT

We considered the following dataset[8]. We had 100 worms. And we had the following types of worms in the following amounts.

- 0 teeth: 2 (Probability: p0=0.02)

- 1 tooth: 3 (Probability: p1=0.03)

- 2 teeth: 5 (Probability: p2=0.05)

- 3 teeth: 14(Probability: p3=0.14)

- 4 teeth: 16 (Probability: p4=0.16)

- 5 teeth: 15 (Probability: p5=0.15)

- 6 teeth: 12 (Probability: p6=0.12)

- 7 teeth: 8 (Probability: p7=0.08)

- 8 teeth: 10 (Probability: p8=0.1)

- 9 teeth: 8 (Probability: p9=0.08)

- 10 teeth: 7 (Probability: p10=0.07)

## 2.1   Model with a uniform distribution

we first tried to model this distribution with a uniform distribution. A uniform distribution has only a single parameter; the uniform probability; the probability of a given event happening.

Figure 2.1: Probability Distribution

$$p_{uniform} = 1/total events = 1/11 = 0.0909$$

This is what the uniform distribution and the true distribution side-by-side looks like.
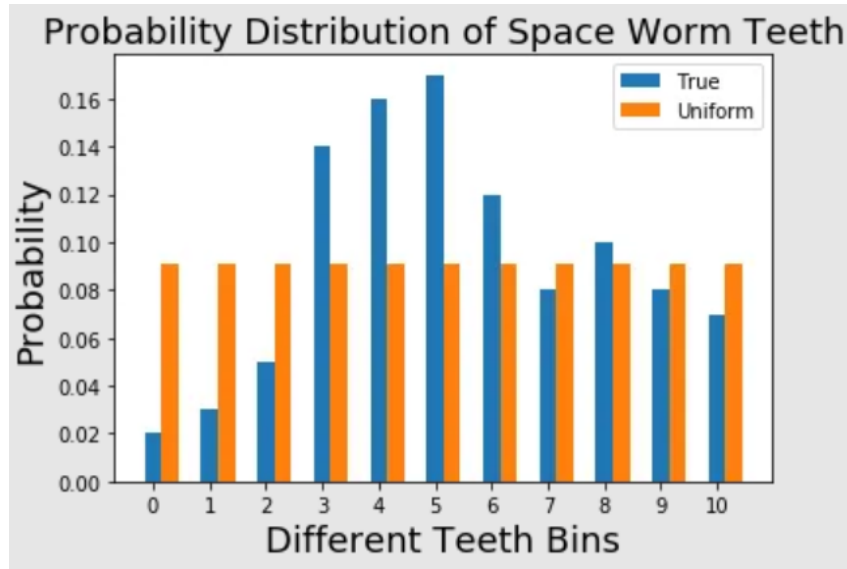


Figure 2.2: uniform distribution and true distribution

## 2.2   Model with a binomial distribution

The binomial probability, commonly used in scenarios such as flipping a coin, can also be applied to our problem[8]. This probability distribution arises from repeating a binary

experiment, such as a coin flip, for a fixed number of trials. In this context, if we denote the probability of success as $p$ and the number of trials as $n$, the probability of obtaining $k$ successes can be calculated as follows:

$$P(X = K) = \binom{n}{k} p^k (1 - p)^n - k$$

## 2.3 Breaking down the equation

The number of different permutations k elements to be arranged within n spaces is given by,

$$\binom{n}{k} = \frac{n!}{k!(n - k)!}$$

Multiplying all these together gives us the binomial probability of k successes.

## 2.4 Back to modeling

Let us first calculate the expected number of teeth for the worms. It would be,

$$0 * p_0 + 1 * p_1 + ... + 10 * p_1 0$$

$$= 0 * 0.02 + 1 * 0.03 + ... + 10 * 0.07$$

$$= 5.44$$

With mean known, we calculate p where,

mean = np

5.44 = 10p

p = 0.544

The maximum number of teeth observed from the population of worms was denoted by $n$.Probabilities of any number of teeth could be defined as follows.

The probability $p_k^{bi}$ was calculated for all different values of $k$. Here, $k$ became the number of teeth we wanted to observe. $p_k^{bi}$ represented the binomial probabilities for the $k$th bin of teeth (that is, 0 teeth, 1 tooth, etc.). This calculation was performed as follows:

$$p0^{bi} = (10!/(0!10!))0.544(1–0.544)^{10} = 0.0004$$

$$p1^{bi} = (10!/(1!9!))0.544(1-0.544) = 0.0046$$

$$p2^{bi} = (10!/(2!8!))0.544(1-0.544) = 0.0249$$

$$...$$

$$p9^{bi} = (10!/(9!1!))0.544(1-0.544) = 0.0190$$

$$p10^{bi} = (10!/(10!0!))0.544^{10}(1-0.544) = 0.0023$$

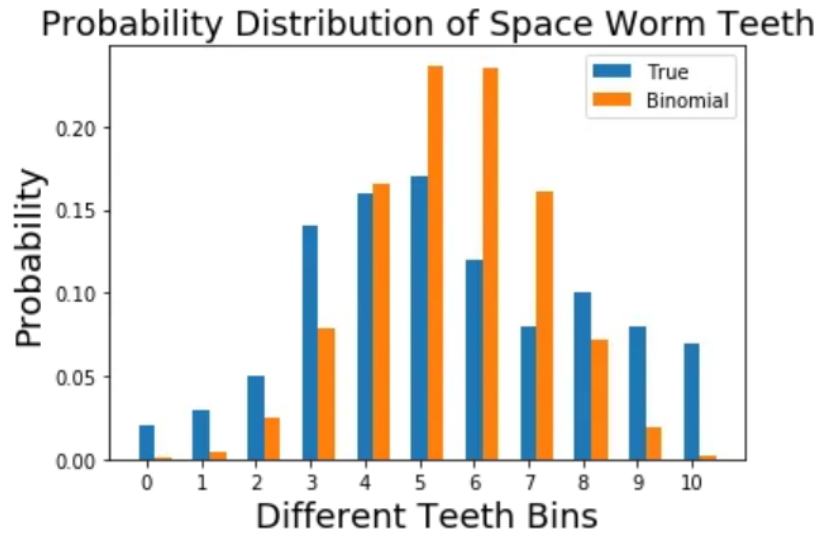This is what a comparison between the true distribution and the binomial distribution looks like



Figure 2.3: true distribution vs binomial distribution

we quantitatively decided which one was the best. This was where the KL divergence came in. KL divergence was formally defined as follows.

$$D_{\mathrm{KL}}(P||Q) = \sum P(X) \log(\frac{P(X)}{Q(X)})$$

In this context, q(x) represents the approximation, and p(x) denotes the true distribution we aimed to match q(x) to. Intuitively, this measures how much a given arbitrary distribution differs from the true distribution. If two distributions perfectly matched, $D_{\mathrm{KL}}(p||q) = 0$. Otherwise, it could take values between 0 and $\infty$. $A lower KL divergence value indicates a better matc$

## 2.5 Computing KL divergence

we computed the KL divergence for each of the approximate distributions we derived. First, let's consider the uniform distribution.

$$D_{\text{KL}}(True||Uniform) = 0.02(0.02/0.0909)+0.03(0.03/0.0909)+...+0.08(0.08/0.0909+0.07(0.07/0.09$$

$$D_{\text{KL}}(True||Uniform) = 0.136$$

Now for the binomial distribution we get,

$$D_{\text{KL}}(True||Uniform) = 0.02*log(0.02/0.0909)+0.03*log(0.03/0.0909)+...+0.08*log(0.08/0.0909+0.$$

$$D_{\text{KL}}(True||Uniform) = 0.427$$

### 2.5.1 KL Divergence with respect to Binomial Mean

we observed how the KL divergence changes as the success probability of the binomial distribution varies. Unfortunately, we couldn't conduct a similar exploration with the uniform distribution as we couldn't alter the probability due to the fixed nature of $n$.

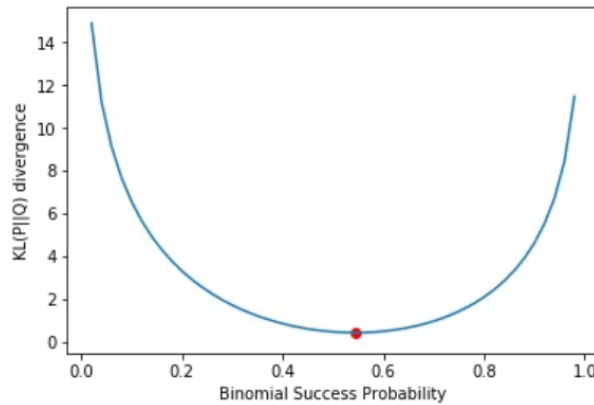As we moved away from our selected probability (represented by the red dot), the KL



Figure 2.4: KL Divergence with respect to Binomial Mean

divergence increased rapidly. In fact, upon examining some KL divergence values slightly deviating from our chosen probability, it became evident that our initial selection of the success probability minimized the KL divergence[8].

# CHAPTER 3

# APPLICATIONS

## 3.1 Monitoring Data Drift

One of the most common use cases of KL divergence in machine learning is to detect drift in datasets. Data is constantly changing, and a metric is required to assess the significance of the changes.

Constantly monitoring data allows machine learning teams to decide whether the model re-training is required. It also provides insights regarding the variable nature of the data under consideration and helps draw statistical analysis. KL divergence is applied to data in discrete form by forming data bins. The data points are binned according to the features to form discrete distributions, i.e., each feature is independently processed for divergence calculation. The divergence scores for each bin are summed up to get a final picture.

## 3.2 Variational Auto-Encoder Optimization

An auto-encoder is a neural network architecture that encodes an input image onto an embedding layer. Variational auto-encoders (VAEs) are a specialized form of traditional architecture that project the input data onto a probability distribution (usually a Gaussian distribution).

The VAE architecture involves two loss functions, a mean-squared error (MSE) to calculate the loss between the output image and the ground truth and the KL divergence to calculate the statistical distance between the true distribution and the approximating distribution in the middle.
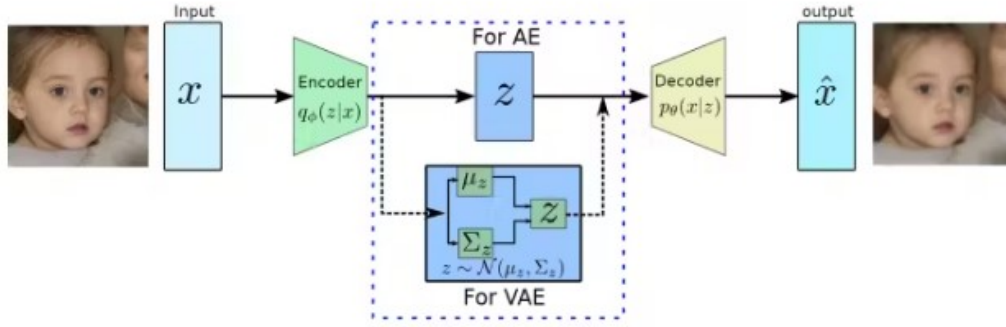
Figure 3.1: Variatianal autoencoder

## 3.3 Generative Adversarial Networks

The Generative Adversarial Networks (GANs) consist of two networks training to outdo each other (adversaries). One of the networks trains to output synthetic (fake) data indistinguishable from the actual input, while the other learns to detect synthetic images from actual.



Figure 3.2: overview of GAN structure

The datasets involved in GANs (fake and real) are two distributions the network tries to compare. The comparison uses KL divergence to create a comparable metric to evaluate whether the model is learning. The discriminator tries to maximize the divergence metric, while the generator tries to minimize it, forming an overall min-max loss configuration. Moreover, it is essential to mention that KL divergence does have certain drawbacks, such as its asymmetric behavior and unstable training dynamics, that make the Jensen-Shannon Divergence a better fit.

## 3.4 Loss Function for Neural Networks

While regression problems use the mean-squared error (MSE) loss function for the maximum likelihood estimation, a classification problem works with probabilistic distributions. A classification neural network is coupled with the KL divergence loss function to compare the model outputs to the true labels.

For example, a binary cat classifier predicts the probabilities for an image as p(cat) = 0.8 and p(not cat) = 0.2. If the ground truth of the image is that it is a cat, then the true distribution becomes p(cat) = 1 and p(not cat) = 0. We now have two different distributions modeling the same variable. The neural network aims to bring these predicted distributions as close to the ground truth as possible. Taking values from our example above as $P(X) = \{1, 0\}$ and $Q(X) = \{0.8, 0.2\}$, the divergence would be:

$$D_{\text{KL}}(P||Q) = \sum P(X) \log(\frac{P(X)}{Q(X)})$$

$$D_{\text{KL}}(P||Q) = \sum P(X) \log(P(X)) - P(X) \log(Q(X))$$

$$D_{\text{KL}}(P||Q) = 1 \log(1) - 1 \log(0.8) + 0 \log(0) - 0 \log(0.2)$$

Most of the terms above will be zeroed out and the final result comes out to be:

$$D_{\text{KL}}(P||Q) = -\log(0.8) = 0.0969$$

When using KL divergence as the loss, the network optimizes to bring the divergence value down to zero.

### 3.4.1 Implementing KL divergence as Loss function

Necessary modules from Keras to build and train the model, including datasets, model layers, and backend utilities has been imported.

**Model Configuration:**

```
img_width, img_height          = 32, 32
batch_size                     = 250
no_epochs                      = 25
no_classes                     = 10
validation_split               = 0.2
verbosity                      = 1
```

Defines configuration parameters such as image width, height, batch size, number of epochs, number of classes, validation split ratio, and verbosity level.

**Loading CIFAR10 Dataset:**

```
(input_train, target_train), (input_test, target_test) =
    cifar10.load_data()
```

Loads the CIFAR10 dataset using Keras' built-in function cifar10.load_data().

- **CIFAR-10 dataset**

  The CIFAR-10 dataset (Canadian Institute for Advanced Research, 10 classes) is a subset of the Tiny Images dataset and consists of 60000 32x32 color images. The images are labelled with one of 10 mutually exclusive classes: airplane, automobile (but not truck or pickup truck), bird, cat, deer, dog, frog, horse, ship, and truck (but not pickup truck). There are 6000 images per class with 5000 training and 1000 testing images per class.

**Data Preprocessing:**

```python
if K.image_data_format() == 'channels_first':
    input_train = input_train.reshape(input_train.shape[0],3,
        img_width, img_height)
    input_test = input_test.reshape(input_test.shape[0], 3,
        img_width, img_height)
    input_shape = (3, img_width, img_height)
else:
    input_train = input_train.reshape(input_train.shape[0],
        img_width, img_height, 3)
    input_test = input_test.reshape(input_test.shape[0],
        img_width, img_height, 3)
    input_shape = (img_width  , img_height, 3)


# Parse numbers as floats
input_train = input_train.astype('float32')
input_test = input_test.astype('float32')


# Normalize data.
```

```
input_train = input_train / 255
input_test = input_test / 255


# Convert target vectors to categorical targets
target_train = keras.utils.to_categorical(target_train,
    no_classes)
target_test = keras.utils.to_categorical(target_test,
    no_classes)
```

- Reshapes the input data based on the image data format (channels first or channels last).

- Converts the data type of input images to float32.

- Normalizes the pixel values of input images to the range [0, 1].

- Converts target vectors to categorical targets using one-hot encoding.

**Creating the Model:**

```
model = Sequential()
model.add(Conv2D(32, kernel_size=(3, 3), activation='relu',
    input_shape=input_shape))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.50))
model.add(Conv2D(64, kernel_size=(3, 3), activation='relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.50))
model.add(Flatten())
model.add(Dense(256, activation='relu'))
model.add(Dense(no_classes, activation='softmax'))
```

- Initializes a Sequential model.

- Adds layers to the model:

    - Two Conv2D layers with ReLU activation functions, followed by MaxPooling2D layers and Dropout regularization.

16

- A Flatten layer to flatten the output of the convolutional layers.

- Two Dense layers with ReLU activation for feature extraction and softmax activation for classification.

**Model compilation & starting training**

We then compile the model and start the training by fitting the data:

```python
model.compile(loss=keras.losses.kullback_leibler_divergence,
              optimizer=keras.optimizers.Adam(),
              metrics=['accuracy'])
# Fit data to model
model.fit(input_train, target_train,
          batch_size=batch_size,
          epochs=no_epochs,
          verbose=verbosity,
          validation_split=validation_split)
```

Configures the model for training with Kullback-Leibler divergence loss, Adam optimizer, and accuracy metric.

We next fit the data to the model, or in plain English start the training process. We do so by feeding the training data (both inputs and targets), specifying the batch size, number of epochs, verbosity and validation split configured before.Splits the training data into training and validation sets based on the specified validation split ratio.

**Evaluating the Model:**

Evaluates the trained model on the test data to generate generalization metrics,Prints the test loss and test accuracy.

```python
# Generate generalization metrics
score = model.evaluate(input_test, target_test, verbose=0)
print(f'Test loss: {score[0]} / Test accuracy: {score[1]}')
```

Overall, this code demonstrates the process of building, training, and evaluating a convolutional neural network model for image classification using the CIFAR10 dataset in Keras.

## 3.4.2 Results
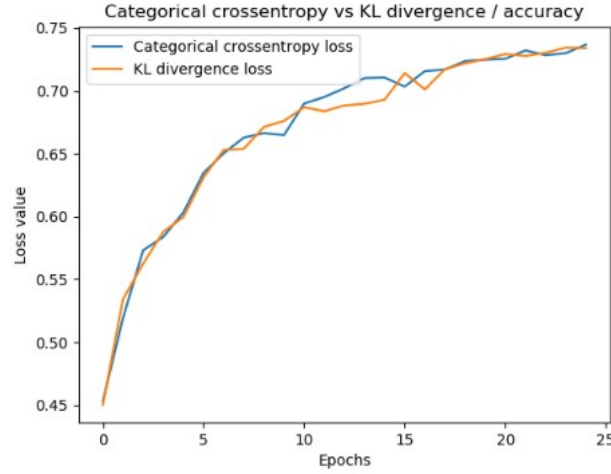
**Categorical crossentropy vs KL divergence/accuracy**



Figure 3.3: Categorical crossentropy vs KL divergence/accuracy

**Categorical crossentropy vs KL divergence/Validation loss**

In 25 epochs, performance is very similar.Cross entropy is typically used in supervised
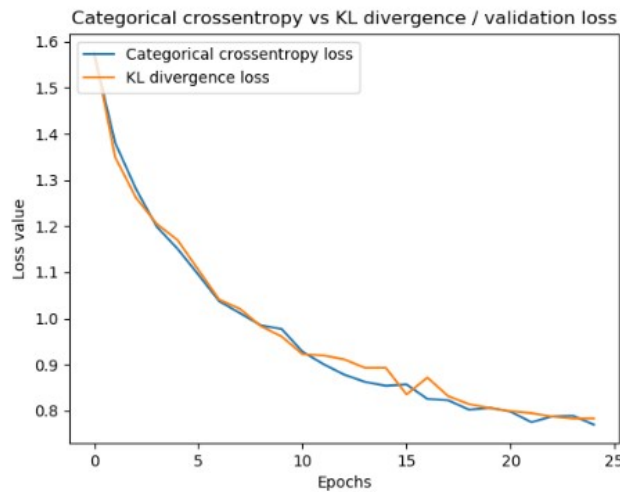


Figure 3.4: Categorical crossentropy vs KL divergence/Validation loss

learning tasks where there is a fixed target distribution, while KL divergence is more suitable for unsupervised learning tasks and applications involving the comparison or approximation of two probability distributions.

# CHAPTER 4

# CONCLUSION

In this project, we explored the utilization of Kullback-Leibler (KL) divergence as a loss function to enhance the performance of a convolutional neural network (CNN) for image classification tasks. KL divergence, a measure of dissimilarity between two probability distributions, offers a unique perspective on optimizing model performance by capturing the difference between predicted and actual probability distributions.

By integrating KL divergence as the loss function in our CNN model trained on the CIFAR-10 dataset, we aimed to improve both the accuracy and generalization capabilities of the model. The architecture of the CNN consisted of convolutional layers, max-pooling layers, dropout regularization, and fully connected layers, providing a robust framework for feature extraction and classification.

During training, the model iteratively learned to minimize the KL divergence loss, thereby adjusting its parameters to better align predicted class probabilities with ground truth labels. Through multiple epochs of training, the CNN optimized its performance, achieving higher accuracy and better generalization on unseen test data.

Our experimental results demonstrated the effectiveness of leveraging KL divergence as a loss function for enhancing machine learning performance. The CNN model achieved competitive accuracy metrics while showcasing improved robustness and generalization capabilities.

In conclusion, the integration of Kullback-Leibler divergence as a loss function offers a promising avenue for advancing the performance of machine learning models, particularly in tasks requiring probabilistic interpretation and classification.

## 4.1 Future Scope

- Explore further mathematical properties and implications of KL divergence loss in different ML models, such as neural networks, decision trees, and ensemble methods.

- Explore applications of KL divergence loss in semi-supervised and unsupervised learning scenarios, where it can encourage meaningful representations and clustering structures.

# REFERENCES

[1] Shun-ichi Amari, *Differential-geometrical methods in statistics*, vol. 28, Springer Science & Business Media, 2012.

[2] Kenneth P Burnham and David R Anderson, *Model selection and multimodel inference*, A practical information-theoretic approach **2** (2004).

[3] Imre Csiszár, *Eine informationstheoretische ungleichung und ihre anwendung auf beweis der ergodizitaet von markoffschen ketten*, Magyer Tud. Akad. Mat. Kutato Int. Koezl. **8** (1964), 85–108.

[4] Maurice Fréchet, *Les éléments aléatoires de nature quelconque dans un espace distancié*, Annales de l'institut Henri Poincaré, vol. 10, 1948, pp. 215–310.

[5] John R Hershey and Peder A Olsen, *Approximating the kullback leibler divergence between gaussian mixture models*, 2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP'07, vol. 4, IEEE, 2007, pp. IV–317.

[6] Harold Jeffreys, *An invariant form for the prior probability in estimation problems*, Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences **186** (1946), no. 1007, 453–461.

[7] Solomon Kullback and Richard A Leibler, *On information and sufficiency*, The annals of mathematical statistics **22** (1951), no. 1, 79–86.

[8] Jay Patel, *Light on math: Machine learning intuitive guide to understanding kl divergence*, 2019, Towards Data Science.

[9] Claude Elwood Shannon, *A mathematical theory of communication*, The Bell system technical journal **27** (1948), no. 3, 379–423.