# Study Material - Youtube

#### **Document Information**

• Generated: 2025-08-01 22:49:37

• Source: https://youtu.be/2Sp0BqAWWXY

• Platform: Youtube

• Word Count: 1,556 words

• Estimated Reading Time: ~7 minutes

• Number of Chapters: 3

• Transcript Available: Yes (analyzed from video content)

## **Table of Contents**

1. DDPMs as Score Predictors: A Deep Dive

2. Key Takeaways from This Video

3. Self-Assessment for This Video

# Video Overview

This lecture, titled "DDPMs as score-predictors," presents a powerful alternative interpretation of Denoising Diffusion Probabilistic Models (DDPMs). The instructor demonstrates that a DDPM, which is typically trained to predict the noise added during the forward process, can be equivalently viewed as a model that predicts the **score function** of the noisy data distribution. This connection is formally established using **Tweedie's formula**, a classical result from statistics. Understanding DDPMs as score predictors provides a deeper insight into their mechanism and bridges the gap between DDPMs and another important class of generative models known as score-based models. This perspective is particularly valuable for advanced applications like conditional image generation.

#### Learning Objectives

Upon completing this lecture, students will be able to: - **Define and understand the score function**  $(\nabla \log p(x))$  of a probability distribution and its intuitive meaning. - **Comprehend Tweedie's formula** and how it relates the mean of a Gaussian distribution to an observed sample and the score function. - **Apply Tweedie's formula to the DDPM forward process** to establish a direct mathematical link between the added noise  $(\epsilon_t)$  and the score function. - **Recognize that a DDPM is implicitly learning the score function**, even when its objective is to predict noise. - **Appreciate the equivalence** of various DDPM training objectives, such as predicting the original data, the added noise, or the score function.

# Prerequisites

To fully grasp the concepts in this lecture, students should have a solid understanding of: - **Denoising Diffusion Probabilistic Models (DDPMs):** Familiarity with the forward (diffusion) and reverse (denoising) processes. - **Probability and Statistics:** Concepts of Gaussian distributions, conditional probability, expectation, and log-likelihood. - **Multivariable Calculus:** A strong grasp of gradients ( $\nabla$ ) and their interpretation. - **Basic Linear Algebra:** Understanding of vectors and matrices, particularly the identity matrix.

#### **Key Concepts Covered**

- Score Function
- Tweedie's Formula
- DDPM as a Score Predictor

• Equivalence of DDPM Training Objectives

# DDPMs as Score Predictors: A Deep Dive

This section explores an alternative and highly insightful interpretation of what a Denoising Diffusion Probabilistic Model (DDPM) learns. While we have previously seen DDPMs as models that predict the original data  $(x_0)$  or the added noise  $(\epsilon_t)$ , we will now demonstrate their equivalence to **score-based models**.

#### The Score Function and Tweedie's Formula

## **Intuitive Foundation**

At the heart of this new perspective is the **score function**. For any given probability distribution, the score function at a particular point tells us the direction in which the probability density increases most steeply. Imagine a landscape where the height represents probability. The score function is a vector that always points "uphill" towards the nearest peak (mode) of the distribution.

A fundamental result from statistics, **Tweedie's formula**, provides a remarkable connection between this score function and the mean of a Gaussian distribution. It essentially states that if we have a data point sampled from a Gaussian, we can make a better guess about the distribution's mean by starting at our data point and taking a small step in the direction indicated by the score function.

## Mathematical Analysis of Tweedie's Formula (00:51)

The instructor introduces Tweedie's formula as a cornerstone for understanding DDPMs as score predictors.

1. Statistical Setup: Let's consider a random variable t drawn from a multivariate Gaussian distribution with mean  $\mu_t$  and covariance  $\Sigma_t$ . We denote this as:

$$t \sim \mathcal{N}(t; \mu_t, \Sigma_t)$$

The probability density function (PDF) for this distribution is denoted by p(t).

2. The Score Function Definition (02:51): The score function is defined as the gradient of the log-probability density with respect to the random variable t. > Definition: Score Function > The score function of a distribution p(t) is given by: >

$$score(t) = \nabla_t \log p(t)$$

> It's crucial to note that the gradient is taken with respect to the data variable t, not the model parameters.

3. Tweedie's Formula (01:21): Tweedie's formula establishes a relationship between the conditional expectation of the mean  $\mu_t$  (given an observation t) and the score function at that observation.

Tweedie's Formula: For a random variable  $t \sim \mathcal{N}(t; \mu_t, \Sigma_t)$ , the conditional expectation of the mean is:

$$\mathbb{E}[\mu_t|t] = t + \Sigma_t \cdot \nabla_t \log p(t)$$

Intuitive Breakdown of the Formula: -  $\mathbb{E}[\mu_t|t]$ : This is our best estimate of the true mean  $\mu_t$  after observing a single sample t. - t: Our starting point is the observed sample. -  $\nabla_t \log p(t)$ : This is the score function, which provides the direction to move from t to get closer to the mean (the region of highest probability). -  $\Sigma_t$ : The covariance matrix acts as a scaling factor. A larger variance means we should take a larger step, as the distribution is more spread out.

## Connecting DDPMs and Score Prediction

We can now apply Tweedie's formula to the forward process of a DDPM to reveal its connection to score prediction.

#### Step 1: Applying Tweedie's Formula to the DDPM Forward Process (05:02)

Recall the distribution of the noisy image  $x_t$  at timestep t, given the original image  $x_0$ :

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)$$

Let's map this to the terms in Tweedie's formula: - The random variable is  $x_t$ . - The mean is  $\mu = \sqrt{\bar{\alpha}_t} x_0$ . - The covariance is  $\Sigma_t = (1 - \bar{\alpha}_t)I$ .

Applying Tweedie's formula, we get an expression for the conditional expectation of the mean:

$$\mathbb{E}[\mu|x_t] = x_t + (1 - \bar{\alpha}_t) \nabla_{x_t} \log q(x_t)$$

The best estimate for this conditional expectation is the true mean itself. Therefore:

$$\sqrt{\bar{\alpha}_t}x_0 = x_t + (1 - \bar{\alpha}_t)\nabla_{x_t}\log q(x_t)$$

## Step 2: Deriving the True Score Function (16:47)

We now have two different ways to express  $x_0$ : 1. From Tweedie's Formula (rearranged):

$$x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( x_t + (1 - \bar{\alpha}_t) \nabla_{x_t} \log q(x_t) \right)$$

#### 2. From the DDPM Forward Process Definition:

$$x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_t \right)$$

By equating these two expressions for  $x_0$  and simplifying, we arrive at a profound result for the true score function:

$$\begin{split} (1-\bar{\alpha}_t)\nabla_{x_t}\log q(x_t) &= -\sqrt{1-\bar{\alpha}_t}\epsilon_t \\ \nabla_{x_t}\log q(x_t) &= -\frac{\epsilon_t}{\sqrt{1-\bar{\alpha}_t}} \end{split}$$

**Key Insight:** The true score of the noisy data distribution  $q(x_t)$  is simply the **negatively scaled** noise  $\epsilon_t$  that was added to create  $x_t$ .

This means that when we train a neural network  $\epsilon_{\theta}(x_t, t)$  to predict the noise  $\epsilon_t$ , we are implicitly training it to predict a scaled version of the score function.

## Step 3: The DDPM Objective as Score Matching (13:45)

The standard DDPM loss function aims to match the predicted noise with the true noise:

$$L_{noise} \propto ||\epsilon_t - \epsilon_\theta(x_t, t)||_2^2$$

Given our new insight, we can re-interpret this. If we define a score-predicting network  $S_{\theta}(x_t, t)$  such that:

$$S_{\theta}(x_t,t) \approx \nabla_{x_t} \log q(x_t) = -\frac{\epsilon_t}{\sqrt{1-\bar{\alpha}_t}}$$

This implies that our noise predictor is related to the score predictor by:

$$\epsilon_{\theta}(x_t,t) \approx -\sqrt{1-\bar{\alpha}_t}S_{\theta}(x_t,t)$$

Substituting this into the loss function reveals that minimizing the noise prediction error is equivalent to minimizing the score matching error:

$$L_{score} \propto ||\nabla_{x_t} \log q(x_t) - S_{\theta}(x_t, t)||_2^2$$

This confirms that a DDPM is fundamentally a score-based model.

The following diagram illustrates this alternative interpretation:

```
flowchart TD
A["Input<br>x_t, t"] --> B{U-Net<br>S<sub>&theta;</sub>(x<sub>t</sub>, t)};
B --> C["Predicted Score<br>S<sub>&theta;</sub>(x<sub>t</sub>)"];
D["True Score<br>><sub>x<sub>t</sub></sub> log p(x<sub>t</sub>)"] --> E{Loss Calculation};
C --> E;
E --> F["Minimize<br>|| <sub>x<sub>t</sub></sub> log p(x<sub>t</sub>) - S<sub>&theta;</sub>(x<sub>t</sub>)
```

**Figure 1:** A DDPM viewed as a regressor on the score function. The U-Net takes the noisy data  $x_t$  and timestep t to predict the score, which is then compared to the true score to compute the loss.

# Key Takeaways from This Video

- **DDPMs** are Implicit Score Predictors: The central message is that training a DDPM to denoise an image (by predicting the noise  $\epsilon_t$ ) is mathematically equivalent to training it to predict the score function  $(\nabla_{x_*} \log q(x_t))$  of the noisy data distribution.
- True Score is Scaled Noise: For the DDPM forward process, the true score is elegantly shown to be the negatively scaled version of the noise that was added, i.e.,  $\nabla_{x_t} \log q(x_t) \propto -\epsilon_t$ .
- Equivalence of Objectives: The lecture highlights that the objectives of predicting noise, predicting the original data  $x_0$ , or predicting the score are all mathematically equivalent, differing only by scaling and shifting. This provides flexibility in how DDPMs are formulated and understood.
- Connection to Score-Based Models: This interpretation formally connects DDPMs to the broader family of score-based generative models, unifying different approaches to generative modeling.

# Self-Assessment for This Video

Test your understanding of the concepts covered in this lecture.

**Question 1:** In your own words, what is the "score function" of a probability distribution, and what is its geometric interpretation?

Question 2: You are given a data point t=5 sampled from a 1D Gaussian distribution  $p(t)=\mathcal{N}(t;\mu,\sigma^2=4)$ . If the score at this point is  $\nabla_t \log p(t) = -0.5$ , what is the best estimate for the mean  $\mu$  according to Tweedie's formula?

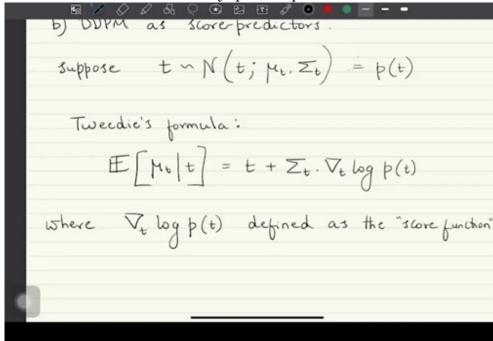
Question 3: Explain the logical steps that connect the DDPM noise prediction objective,  $||\epsilon_t - \epsilon_{\theta}(x_t, t)||_2^2$ , to the score matching objective,  $||\nabla_{x_t} \log q(x_t) - S_{\theta}(x_t, t)||_2^2$ .

**Question 4:** Why is the insight that "true score = negatively scaled noise" so important for understanding DDPMs?

Question 5: If a DDPM is trained to predict the score function  $S_{\theta}(x_t, t)$ , how can you recover the predicted noise  $\epsilon_{\theta}(x_t, t)$  from it?

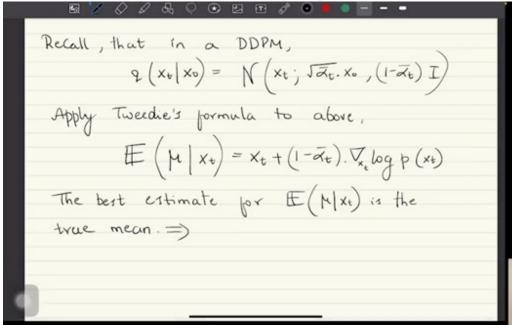
## Visual References

A diagram illustrating the score function ( $\log p(x)$ ) as vectors on a probability density land-scape. This visual explains the intuition that the score function always points 'uphill' towards ar-



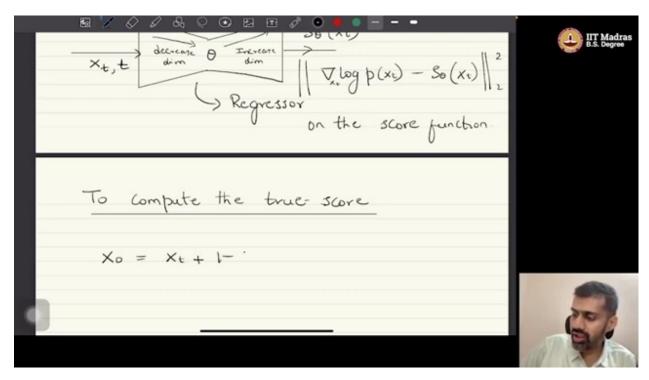
eas of higher probability. (at 03:15):

The formal mathematical statement of Tweedie's formula. This slide presents the key equation that connects the mean of a Gaussian distribution to an observed sample and the score function of



**the prior.** (at 07:30):

\*\*The final step of the derivation applying Tweedie's formula to the DDPM forward process. This screenshot shows the crucial equation establishing the direct relationship between the added noise  $(\underline{\phantom{x}}t)$  and the score function of the noisy data distribution  $(\underline{\phantom{x}}t)$  tog  $\underline{\phantom{x}}t$  (at 15:45):



A summary slide or concept map that visually demonstrates the equivalence of different DDPM training objectives. It shows how predicting the noise, predicting the original data (x\_0), and predicting the score function are all mathematically linked and valid approaches. (at 22:10):

