MALLA REDDY COLLEGE OF ENGINEERING & TECHNOLOGY

Autonomous Institution – UGC, Govt. of India



Department of CSE Artificial Intelligence and Machine Learning

B. TECH (R-22 Regulation)

(IV YEAR - I SEM)

2024-25

COMPUTER VISION (R22A6606)



LECTURE NOTES

MALLA REDDY COLLEGE OF ENGINEERING & TECHNOLOGY

(Autonomous Institution – UGC, Govt. of India)

Recognized under 2(f) and 12(B) of UGC ACT 1956 (Affiliated to JNTUH, Hyderabad, Approved by AICTE-Accredited by NBA & NAAC – 'A' Grade - ISO 9001:2015 Certified) Maisammaguda, Dhulapally (Post Via. Hakimpet), Secunderabad–500100, Telangana State, India

Department of Computer Science and Engineering

(Artificial Intelligence and Machine Learning)

Vision

To be a premier center for academic excellence and research through innovative interdisciplinary collaborations and making significant contributions to the community, organizations, and society as a whole.

Mission

- To impart cutting-edge Artificial Intelligence technology in accordance with industry norms.
- To instill in students a desire to conduct research in order to tackle challenging technical problems for industry.
- To develop effective graduates who are responsible for their professional growth, leadership qualities and are committed to lifelong learning.

QUALITY POLICY

- To provide sophisticated technical infrastructure and to inspire students to reach their full potential.
- To provide students with a solid academic and research environment for a comprehensive learning experience.
- To provide research development, consulting, testing, and customized training to satisfy specific industrial demands, thereby encouraging self-employment and entrepreneurship among students.

For more information: www.mrcet.ac.in

IV Year B. Tech CSE (AI & ML) -I Sem

3/-/-/3

(R22A6606) COMPUTER VISION

COURSE OBJECTIVES

- 1. To introduce various components of image processing techniques for computer vision.
- 2. To understand filters and computing Image Gradient.
- 3. To understand segmentation, model fitting and tracking
- 4. To impart knowledge about object registration and object matching
- 5. To implement various techniques available for object recognition.

UNIT-I

IMAGE FORMATION: Geometric Camera Models, Intrinsic and Extrinsic Parameters, Geometric Camera Calibration – Linear and Non – linear approach, Light and Shading - Inference from, Modeling Inter reflection, Human Color Perception.

UNIT-II

EARLY VISION: Linear Filters - Convolution, Fourier Transforms, Sampling and Aliasing, Filters as Templates, Correlation, Local Image Features - Computing the Image Gradient, Gradient Based Edge Detectors, Orientations, Texture - Local Texture Representations Using Filters, Shape from Texture.

UNIT-III

MID-LEVEL VISION: Segmentation by Clustering - Basic Clustering Methods, The Watershed Algorithm, Segmentation Using K-means, Grouping and Model Fitting - Fitting Lines with the Hough Transform, Fitting Curved Structures, Tracking - Tracking by Detection, Tracking Translations by Matching, Tracking Linear Dynamical Models with Kalman Filters.

UNIT-IV

HIGH-LEVEL VISION: Registration, Registering Rigid and Deformable Objects, Smooth Surfaces and Their Outlines - Contour Geometry, Koenderink's Theorem, The Bitangent Ray Manifold, Object Matching using Interpretation Trees and Spin Images, Classification, Error, and Loss.

UNIT-V

OBJECT DETECTION AND RECOGNITION: Detecting Objects in Images - The Sliding Window Method, Face Detection, Detecting Humans, Boundaries and Deformable Objects, Object Recognition–Categorization, Selection, Applications – Tracking People, Activity Recognition.

TEXT BOOKS:

- 1. Forsyth, Jean Ponce David A. "Computer Vision: A Modern Approach", Second Edition, Pearson Education Limited 2015.
- 2. Szeliski, Richard, "Computer vision: algorithms and applications", Springer Science & Business Media, 2010.

REFERENCE BOOKS:

- 1. Hau, Chen Chi, "Handbook of pattern recognition and computer vision", World Scientific, Fifth Edition, 2015.
- 2. Muhammad Sarfraz, "Computer Vision and Image Processing in Intelligent Systems and Multimedia Technologies", IGI Global, 2014.
- 3. Theo Gevers, Arjan Gijsenij, Joost van de Weijer, Jan-Mark Geusebroek "Color in Computer Vision: Fundamentals and Applications", Wiley, 2012.

4. Kale, K. V, Mehrotra S.C, Manza. R.R., "Advances in Computer Vision and Information Technology", IK International Pvt Ltd, 2013.

COURSE OUTCOMES:

- 1. Understand various image formation models.
- 2. Extract shape, texture and edge-based features.
- 3. Detect region of interest using image segmentation and object localization techniques.
- 4. Identify and recognize objects using image registration and classification.
- 5. Explore various case studies on vision-based applications.

MALLA REDDY COLLEGE OF ENGINEERING AND TECHNOLOGY CSE (Artificial Intelligence and Machine Learning)

&

B. TECH (Artificial Intelligence Machine Learning)

INDEX

<u>S. No</u>	<u>Unit</u>	<u>Topic</u>	Page No
1	Ι	Introduction	1-5
2	Ι	Geometric Camera Models	5-6
3	Ι	Intrinsic & Extrinsic Parameters	6-9
4	Ι	Geometric Camera Calibration	9-10
5	T	a) Linear & Non – Linear Approach	11.12
5	1 	b) Light and Shading – Interence from Modelling	11-15
0	1	Human Color Perception	13-15
7	II	Linear Filters	16-29
8	II	Local Image Features	23-26
		a) Computing the Image Gradient	
9	II	b) Gradient Based Edge Detectors	26-30
10	II	Texture	30-32
11	III	Mid – Level Vision, Segmentation by Clustering	33-38
12	III	Grouping and Model Fitting a) Fitting Lines with Hough Transform	38-42
13	III	b) Fitting Curved Structures	42-43
14	III	Tracking	43-46
		a) Detection and Matching	
15	III	b) Linear Dynamic Models with Kalman Filters	46-48
16	IV	Registration – Rigid and Deformable Objects	49-50
17	IV	Smooth Surfaces & their Outlines	50-54
18	IV	Object Matching using Interpretation Trees & Spin Images	54-56
19	IV	Classification	56
20	IV	Error and Loss	56-58
	± 1		59-64
21	V	Detecting Objects in Images	-
22	V	Detecting Humans	64-68
23	V	Object Recognition	68-72

UNIT-I

IMAGE FORMATION: Geometric Camera Models, Intrinsic and Extrinsic Parameters, Geometric Camera Calibration – Linear and Non – linear approach, Light and Shading - Inference from, Modeling Inter reflection, Human Color Perception.

<u>Computer Vision</u>

Computer vision is a field of artificial intelligence (AI) that uses machine learning and neural networks to teach computers and systems to derive meaningful information from digital images, videos and other visual inputs—and to make recommendations or take actions when they see defects or issues. If AI enables computers to think, computer vision enables them to see, observe and understand.

Working of Computer Vision:

- Computer vision needs lots of data. It runs analyses of data over and over until it discerns distinctions and ultimately recognize images. For example, to train a computer to recognize automobile tires, it needs to be fed vast quantities of tire images and tire-related items to learn the differences and recognize a tire, especially one with no defects.
- Two essential technologies are used to accomplish this: a type of machine learning called deep learning and a convolutional neural network (CNN).
- Machine learning uses algorithmic models that enable a computer to teach itself about the context of visual data. If enough data is fed through the model, the computer will "look" at the data and teach itself to tell one image from another. Algorithms enable the machine to learn by itself, rather than someone programming it to recognize an image.
- A CNN helps a machine learning or deep learning model "look" by breaking images down into pixels that are given tags or labels.
- It uses the labels to perform convolutions (a mathematical operation on two functions to produce a third function) and makes predictions about what it is "seeing."
- The neural network runs convolutions and checks the accuracy of its predictions in a series of iterations until the predictions start to come true

Examples:

Here are some examples of computer vision:

- **Facial recognition**: Identifying individuals through visual analysis.
- Self-driving cars: Using computer vision to navigate and avoid obstacles.
- **Robotic automation**: Enabling robots to perform tasks and make decisions based on visual input.
- **Medical anomaly detection**: Detecting abnormalities in medical images for improved diagnosis.
- **Sports performance analysis**: Tracking athlete movements to analyze and enhance performance.
- **Manufacturing fault detection**: Identifying defects in products during the manufacturing process.
- **Agricultural monitoring**: Monitoring crop growth, livestock health, and weather conditions through visual data.



OpenCV (Open-Source Computer Vision)

- It is a cross-platform and free to use library of functions is based on real-time Computer Vision which supports Deep Learning frameworks that aids in image and video processing.
- In Computer Vision, the principal element is to extract the pixels from the image to study the objects and thus understand what it contains. Below are a few key aspects that Computer Vision seeks to recognize in the photographs:
 - > **Object Detection:** The location of the object.
 - > **Object Recognition:** The objects in the image, and their positions.
 - > **Object Classification:** The broad category that the object lies in.
 - > **Object Segmentation:** The pixels belonging to that object.

IMAGE FORMATION

- Computer vision is a fascinating field that seeks to develop mathematical techniques capable of reproducing the three-dimensional perception of the world around us.
- Vision is an inverse problem, where we seek to recover unknown information from insufficient data to fully specify the solution.
- To solve this problem, it is necessary to resort to models based on physics and probability, or machine learning with large sets of examples.



active environ and sensor

How an Image is Formed?

• Before analyzing and manipulating images, it's essential to understand the image formation process. As examples of components in the process of producing a given image:

1. **Perspective projection:** The way three-dimensional objects are projected onto a twodimensional image, taking into account the position and orientation of the objects relative to the camera.

2. **Light scattering after hitting the surface:** The way light scatters after interacting with the surface of objects, influencing the appearance of colors and shadows in the image.

3. **Lens optics:** The process by which light passes through a lens, affecting image formation due to refraction and other optical phenomena.

4. **Bayer color filter array:** A color filter pattern used in most digital cameras to capture colors at each pixel, allowing for the reconstruction of the original colors of the image.

Focus and Focal Length

Focus is one of the main aspects of image formation with lenses. The focal length, represent f by is the distance between the center of the lens and the focal point, where light rays parallel to the optical axis converge after passing through the lens.



The focal length is directly related to the lens's ability to concentrate light and, consequently, influences the sharpness of the image. The focus equation is given by:

$$\frac{1}{f} = \frac{1}{z} + \frac{1}{e}$$

Areas where mathematical concepts play vital role in image formation:

Here is a high-level overview of the main mathematical components:

- **Coordinate Systems:** Images are represented in a discrete coordinate system. In a 2D image, each point is identified by its (x, y) coordinates. The origin (0, 0) is typically located at the top-left corner of the image.
- **Camera Models:** Cameras capture images by projecting 3D points in the world onto a 2D image plane. The pinhole camera model is commonly used in computer vision. It assumes that light travels through a small aperture (pinhole) and creates an inverted image on the image plane.
- **Intrinsic Parameters:** Intrinsic parameters describe the internal characteristics of the camera. These parameters include the focal length (f), principal point (c_x, c_y), and lens distortion coefficients (k1, k2, etc.). These parameters affect the transformation from 3D world coordinates to 2D image coordinates.

- **Projection Matrix:** The projection matrix combines intrinsic and extrinsic parameters to perform the projection from 3D world coordinates to 2D image coordinates. It is typically represented by a 3x4 matrix.
- **Homogeneous Coordinates:** Homogeneous coordinates are used to represent both 2D and 3D points in computer vision. Homogeneous coordinates use an extra dimension, typically denoted as w, to represent points. This allows for efficient matrix transformations.
- **Perspective Projection:** Perspective projection maps 3D points onto a 2D plane, simulating how objects appear smaller as they move farther away from the camera. It involves dividing the 3D coordinates by the depth (Z) of the point to obtain normalized device coordinates (NDC).
- **Distortion Correction:** Lens distortion occurs due to imperfections in the camera lens, resulting in image distortion. Distortion correction is applied to remove these distortions using distortion coefficients and geometric transformations.
- **Image Rectification:** Image rectification is a transformation applied to images to make them appear as if they were taken from a standard viewpoint, usually by aligning epipolar lines. This is often used in stereo vision for depth estimation.

• Mathematical Formulation:

- 1. *Ray Formation:* To determine the ray of light that intersects the object and passes through the pinhole, we can subtract the camera position from the object position. This gives us a direction vector for the ray: (X C_x, Y C_y, Z C_z).
- Ray Projection: The next step is to project the ray onto the image plane. We can achieve this by scaling the direction vector by the distance f and dividing it by the magnitude of the vector. This normalization step ensures that the vector represents a unit direction: (f * (X − C_x) / ||P||, f * (Y − C_y) / ||P||, f * (Z − C_z) / ||P||).
- 3. *Image Coordinates:* Now we have a ray in 3D space that passes through the pinhole and intersects the object. To obtain the corresponding image coordinates, we need to find the intersection point of the ray with the image plane. Let's denote the image coordinates as (u, v). We can compute them using similar triangles:

$$\begin{split} &u = (f * (X - C_x) / ||P||) / (f * (Z - C_z) / ||P||) \\ &v = (f * (Y - C_y) / ||P||) / (f * (Z - C_z) / ||P||) \end{split}$$

Simplifying the equations, we get:

 $u = (X - C_x) / (Z - C_z)$

 $v = (Y - C_y) / (Z - C_z)$

These equations give us the image coordinates (u, v) for a given object point (X, Y, Z) in the 3D world. By repeating this process for each object point, we can generate the image formed by the pinhole camera.

Challenges

- When it comes to forming images for computer vision, there are several challenges that researchers and developers often encounter. Here are some of the common challenges:
- 1. **Variability in lighting conditions:** Lighting conditions can greatly affect the appearance of an image, making it challenging to extract meaningful information. Shadows, reflections, and uneven illumination can distort or obscure the objects of interest.

- 2. **Variability in scale and viewpoint:** Objects can appear at different scales and viewpoints in images. This variation makes it difficult to develop algorithms that can recognize objects reliably under different perspectives or sizes.
- 3. **Occlusions:** Objects in real-world scenes are often partially or completely occluded by other objects or by the scene itself. Occlusions can make it challenging to accurately detect and recognize objects in an image.
- 4. **Background clutter:** Images can contain complex and cluttered backgrounds that can distract or confuse computer vision algorithms. It becomes difficult to separate the objects of interest from the surrounding clutter.
- 5. **Intra-class variability:** Objects belonging to the same class can exhibit significant variations in appearance, shape, texture, and color. For example, different breeds of dogs or variations in handwritten characters can pose challenges in accurately classifying or recognizing them.
- 6. **Limited training data:** Collecting and annotating large-scale datasets for training computer vision models can be time-consuming and expensive. Limited training data can lead to overfitting or poor generalization performance of the models.
- 7. **Computational complexity:** Many computer vision tasks, such as object detection or semantic segmentation, require analyzing and processing large amounts of data. These tasks can be computationally demanding and may require specialized hardware or efficient algorithms to achieve real-time performance.
- 8. **Robustness to noise:** Images can be corrupted by various types of noise, including sensor noise, compression artifacts, or environmental factors. Ensuring that computer vision algorithms are robust to noise and can provide accurate results is a significant challenge.
- 9. **Ethical and privacy concerns**: Computer vision systems have the potential to invade privacy or be used for unethical purposes. Addressing concerns related to data privacy, bias, fairness, and accountability is crucial for the responsible development and deployment of computer vision technologies.

Geometric Camera Models

- Computer Vision is a sub-area of AI which enables computers to analyze images or videos and extract meaningful information from them, mimicking human vision.
- Some of the applications of CV are face recognition, self-driving cars, object detection, etc.
- When we click a picture, a real scenario which is in 3D is captured by a real camera in 2D. So here, 3D to 2D conversion takes place which means we have lost information about a dimension.
- This is where computer vision enters the scene.
- Using CV, we can try to recover the missing information, or you can say, a missing dimension and gain high level understanding of the image.
- Real Scene (3D) \rightarrow Real Cameras (2D) \rightarrow CV Output (3D)

Pinhole Camera Model

- Before lenses came into existence, pinholes were used to capture images.
- World's first camera model was invented using this model.
- For a pinhole camera, a hole of the size of a pin is created on one side of a box and a thin paper on the other side of the box.
- The light entering this hole will then project the image of the world on the paper.
- The image captured will be upside down i.e. an inverted image.

Below is the example of an image captured by a pinhole camera.



Image formation on the backplate of a photographic camera. Figure from US NAVY MANUAL OF BASIC OPTICS AND OPTICAL INSTRUMENTS, prepared by the Bureau of Naval Personnel, reprinted by Dover Publications, Inc. (1969). Taken from the book 'Computer Vision — A Modern Approach — D. Forsyth, J. Ponce'.

Intrinsic and Extrinsic Parameters:

- Images are one of the most commonly used data in recent deep learning models.
- Cameras are the sensors used to capture images. They take the points in the world and project them onto a 2D plane which we see as images.
- This transformation is usually divided into two parts: *Extrinsic and Intrinsic*.
- The extrinsic parameters of a camera depend on its location and orientation and have nothing to do with its internal parameters such as focal length, the field of view, etc.
- On the other hand, the intrinsic parameters of a camera depend on how it captures the images. Parameters such as focal length, aperture, field-of-view, resolution, etc govern the intrinsic matrix of a camera model.
- These intrinsic and extrinsic parameters are transformation matrices that convert points from one coordinate system to the other. In order to understand these transformations, we first need to understand what are the different coordinate systems used in imaging.

Commonly used coordinate systems in CV



The commonly used coordinate systems in Computer Vision are

- 1. World coordinate system (3D)
- 2. Camera coordinate system (3D)
- 3. Image coordinate system (2D)
- 4. Pixel coordinate system (2D)

The extrinsic matrix is a transformation matrix from the world coordinate system to the camera coordinate system, while the intrinsic matrix is a transformation matrix that converts points from the camera coordinate system to the pixel coordinate system.

World coordinate system (3D):

[Xw, Yw, Zw]: It is a 3D basic cartesian coordinate system with arbitrary origin. For example a specific corner of the room. A point in this coordinate system can be denoted as Pw = (Xw, Yw, Zw).



Object/Camera coordinate system (3D):

[Xc, Yc, Zc]: It's the coordinate system that measures relative to the object/camera's origin/orientation. The z-axis of the camera coordinate system usually faces outward or inward to the camera lens (camera principal axis) as shown in the image above (z-axis facing inward to the camera lens). One can go from the world coordinate system to object coordinate system (and vice-versa) by Rotation and Translation operations.



Image coordinate system (2D) [Pinhole Model]:

[Xi, Yi]: A 2D coordinate system that has the 3D points in the camera coordinate system projected onto a 2D plane (usually normal to the z-axis of the camera coordinate system — shown as a yellow plane in the figures below) of a camera with a Pinhole Model.

- The rays pass the center of the camera opening and are projected on the 2D plane on the other end.
- The 2D plane is what is captured as images by the camera.
- It is a lossy transformation, which means projecting the points from the camera coordinate system to the 2D plane cannot be reversed (the depth information is lost Hence by looking at an image captured by a camera, we can't tell the actual depth of the points).
- The X and Y coordinates of the points are projected onto the 2D plane. The 2D plane is at f (focal-length) distance away from the camera. The projection Xi, Yi can be found by the law of similar triangles (the ray entering and leaving the camera center has the same angle with the x and y-axis, alpha and beta respectively).

Pixel coordinate system (2D):

[**u**, **v**]: This represents the integer values by discretizing the points in the image coordinate system. Pixel coordinates of an image are discrete values within a range that can be achieved by dividing the image coordinates by pixel width and height (parameters of the camera — units: meter/pixel).



Transformations:

- 1. World-to-Camera: 3D-3D projection. Rotation, Scaling, Translation
- **2.** Camera-to-Image: 3D-2D projection. Loss of information. Depends on the camera model and its parameters (pinhole, f-theta, etc)
- 3. Image-to-Pixel: 2D-2D projection. Continuous to discrete. Quantization and origin shift.

$$\begin{pmatrix} u \\ v \\ w \end{pmatrix} = \begin{pmatrix} f/\rho_u & 0 & c_x & 0 \\ 0 & f/\rho_v & c_y & 0 \\ 0 & 0 & 1 & 0 \\ \text{Camera Intrinsic Matrix} \end{pmatrix} \begin{pmatrix} R_{3\times3} & t_{3\times1} \\ 0_{1\times3} & 1_{1\times1} \end{pmatrix}_{(4\times4)} \begin{pmatrix} X_W \\ Y_W \\ Z_W \\ 1 \end{pmatrix}$$

Camera matrices - Image by Author

Camera Extrinsic Matrix (World-to-Camera):

Converts points from world coordinate system to camera coordinate system. Depends on the position and orientation of the camera.

Camera Intrinsic Matrix (Camera-to-Image, Image-to-Pixel):

Converts points from the camera coordinate system to the pixel coordinate system. Depends on camera properties (such as focal length, pixel dimensions, resolution, etc.)

Geometric Camera Calibration

Camera calibration is a fundamental task in computer vision_crucial in various applications such as 3D reconstruction, object tracking, augmented reality, and image analysis. Accurate calibration ensures precise measurements and reliable analysis by correcting distortions and estimating intrinsic and extrinsic camera parameters.

Central Projection

If world and image points are represented by homogeneous vectors, central projection is a linear transformation:



Internal Camera Parameters



Light and Shading - Inference from

Light designates the interaction between materials and light sources. Shading is the process of determining the color of a pixel.

Light and shading play crucial roles in computer vision, influencing how images are captured, processed, and interpreted. Here's a breakdown of their importance and how they affect computer vision:

1. Light and Its Properties:

- **Illumination Variability:** The amount, direction, and color of light can vary greatly, affecting how objects appear in images. For example, a scene illuminated from the side can cast shadows and create highlights that alter the perceived shape and texture of objects.
- **Light Sources:** Different types of light sources (e.g., natural sunlight, incandescent bulbs, fluorescent lights) have distinct spectral properties and intensity levels, which can affect the color and brightness of the captured image.

2. Shading and Its Effects:

- **Shadows:** Shadows provide depth information and can be used to infer the 3D structure of objects. However, shadows can also complicate object detection and recognition if they obscure parts of an object or create misleading visual cues.
- **Highlights:** Bright spots where light directly reflects off surfaces can give clues about material properties and surface orientation. They can also affect how textures are perceived.

3. Challenges in Computer Vision:

- **Illumination Invariance:** For many computer vision applications, such as object recognition, it's important to design algorithms that are invariant to changes in lighting. Techniques like histogram equalization, adaptive thresholding, and illumination normalization can help mitigate these effects.
- **Specular Highlights and Shadows:** These can interfere with object detection and segmentation. Advanced methods may include modeling these effects to distinguish between true object features and lighting artifacts.
- **Color Constancy:** Ensuring that colors are perceived consistently across different lighting conditions is a significant challenge. Algorithms for color constancy attempt to correct for varying illumination to maintain consistent color representation.

4. Applications:

- **3D Reconstruction:** Understanding light and shading is essential for reconstructing 3D shapes from 2D images. Techniques like photometric stereo use variations in shading to infer surface details.
- **Object Detection and Recognition:** Algorithms must account for varying lighting conditions to accurately detect and recognize objects. This can involve using robust features or training models on diverse lighting scenarios.
- **Augmented Reality (AR):** Accurate rendering in AR systems requires real-time adjustment of virtual objects' shading to match the lighting of the real world, enhancing the realism of the experience.

5. Techniques and Tools:

- **Machine Learning:** Modern computer vision systems often use machine learning and deep learning to learn how to handle variations in lighting and shading automatically. Convolutional Neural Networks (CNNs) and other architectures can be trained on large datasets with varied lighting conditions.
- **Physical Models:** Some approaches involve modeling the physical properties of light and shading, such as using the Lambertian reflectance model or more complex Bidirectional Reflectance Distribution Functions (BRDFs).

Modelling Inter reflection

In computer vision, inter-reflection refers to the phenomenon where light rays bounce off surfaces and affect the appearance of neighboring surfaces. Modeling inter-reflection is crucial for accurate scene understanding and rendering. Here are some common approaches and techniques used to model inter-reflection in computer vision:

1. Physically Based Rendering (PBR):

• PBR models simulate the behavior of light in a physically accurate way, including interreflections. They use techniques like ray tracing or path tracing to calculate how light interacts with surfaces and how those interactions affect neighboring surfaces.

2. Radiosity:

• Radiosity is a method for computing global illumination in scenes with diffuse surfaces. It accounts for inter-reflections by solving the rendering equation, considering how light is emitted, reflected, and absorbed by surfaces, and how this affects the light arriving at other surfaces.

3. Monte Carlo Methods:

• Monte Carlo methods, such as Monte Carlo ray tracing, are used to simulate the behavior of light rays in a scene. They can simulate complex light interactions, including inter-reflections, by tracing paths of light rays and computing their interactions with surfaces and other light sources.

4. Light Transport Models:

• These models simulate the transport of light through a scene, considering how light from different sources interacts with surfaces and how that light is redistributed through reflections and refractions. Techniques like Bidirectional Reflectance Distribution Function (BRDF) and Bidirectional Surface Scattering Distribution Function (BSSRDF) are used to model these interactions.

5. Image-based Techniques:

• In some cases, inter-reflections can be estimated directly from images using computer vision techniques. For example, methods based on photometric stereo or shape-from-shading can infer surface properties and inter-reflection effects from multiple images of the same scene under different lighting conditions.

6. Deep Learning Approaches:

• Recent advances in deep learning have also been applied to modeling inter-reflections. Neural networks can learn to infer inter-reflection effects from images or to enhance the realism of rendered scenes by considering light transport properties.

Color perception



Three types of cones

- Each is sensitive to a different region of the spectrum but regions overlap
 - Short (S) corresponds to blue
 - Medium (M) corresponds to green
 - Long (L) corresponds to red
- Different sensitivities: we are more sensitive to green than red
 varies from person to person (and with age)
- · Colorblindness-deficiency in at least one type of cone

RGB color cube (additive color model)



• R, G, B values normalized to (0, 1) interval

• Human perceives gray for triples along the diagonal; origin=black

• Additive: Mix RGB to get colors



Color triangle (CIE system) and normalized RGB



Note: To represent the full gamut of colors (e.g., black), you need to include brightness and therefore you are back in a 3D space (like the RGB cube)

HSI (or HSV) Model (Color hexagon)

Hue: Distinguishes between colors (angle between 0 and 2π).

Saturation: Purity of color (distance on vertical axis (0 to 1)).

Intensity: Light versus dark shades of a color (height along the vertical axis (0 to 1)



CIELAB

Designed to approximate human vision

One luminance channel (L) and two color-opponent channels (a and b).

In this model, the color differences which you perceive correspond to Euclidian distances in CIELab.

The a axis extends from green (-a) to re (+a) and the b axis from blue (-b) to yellow (+b). The brightness (L) increases from the bottom to the top of the three-dimensional model.



UNIT-II

EARLY VISION: Linear Filters - Convolution, Fourier Transforms, Sampling and Aliasing, Filters as Templates, Correlation, Local Image Features - Computing the Image Gradient, Gradient Based Edge Detectors, Orientations, Texture - Local Texture Representations Using Filters, Shape from Texture.

Early Vision

- Early vision in computer vision refers to the initial stages of processing visual information, where the focus is on low-level features and basic image characteristics.
- These stages are crucial for extracting fundamental elements from visual data, which serve as building blocks for higher-level interpretation and understanding.

Linear Filters

- In the context of early vision in computer vision and image processing, linear filters play a crucial role in basic operations that mimic initial stages of human visual perception.
- These filters are typically applied to raw pixel data to extract fundamental features or enhance certain aspects of images. Here are some key applications of linear filters in early vision.

1. Edge Detection:

- Sobel Filter: Detects edges by computing the gradient magnitude in both horizontal and vertical directions.
- Prewitt Filter: Similar to Sobel but uses slightly different coefficients.
- Roberts Cross Operator: Uses a 2x2 kernel to detect edges at 45-degree angles.

2. Smoothing and Noise Reduction:

- Gaussian Filter: Smooths images by applying a Gaussian kernel, which reduces high-frequency noise and blurs the image slightly.
- Mean Filter: Replaces each pixel value with the average of its neighboring pixels, effectively reducing noise and smoothing the image.

3. Feature Extraction:

- a) Laplacian of Gaussian (LoG) Filter: Enhances edges and detects features by first applying a Gaussian smoothing filter and then computing the Laplacian of the resulting image.
- b) Difference of Gaussians (DoG) Filter: Similar to LoG but approximated by the difference between two Gaussian-blurred versions of the image.

4. Texture Analysis:

Gabor Filter: Used for texture analysis and feature extraction. It is a complex wavelet filter that resembles the response of simple cells in the human visual cortex.

5. Corner Detection: Harris Corner Detector: Utilizes a local autocorrelation function to identify corners or interest points in an image.

6. Frequency Analysis:

- a) Fourier Transform Filters: While not strictly linear in the spatial domain, filters based on Fourier transforms are used to analyze frequency components of images.
- b) Linear filters are characterized by their convolution operation, where a small matrix (kernel) is applied to each pixel of the image. The output value for each pixel is typically a linear combination of the neighboring pixel values weighted by the kernel coefficients. These filters are computationally efficient and widely used in early vision tasks such as preprocessing for higher-level vision algorithms, feature extraction, and image enhancement.

Convolution

- **a.** Convolution in image processing involves applying a small matrix called a **kernel** or **filter** to an image.
- b. The kernel is typically a square matrix (e.g., 3x3 or 5x5) with numerical values that define how the operation modifies the image.
- c. Convolution is carried out by sliding this kernel over the image and computing a weighted sum of the pixels that overlap with the kernel at each position.

Steps Involved in Convolution:

1.Kernel Definition: Define a kernel matrix that represents the operation you want to perform on the image. For example, a Gaussian blur can be achieved using a Gaussian kernel matrix.

2.Positioning the Kernel: Place the kernel matrix over the top-left corner of the image.

3.Element-wise Multiplication: Multiply each element of the kernel matrix with the corresponding pixel value in the image.

4. Summation: Sum up all the multiplied values to get the new pixel value for the center position of the kernel.

5. Slide Over: Slide the kernel to the right (and possibly down) to repeat the process for all positions in the image, computing a new value for each pixel based on its neighborhood.

Key Uses of Convolution in Computer Vision:

- 1. **Smoothing and Blurring**: Gaussian and mean filters use convolution to blur images, reducing noise and detail.
- 2. **Edge Detection**: Filters like Sobel and Prewitt use convolution to highlight edges by detecting changes in intensity.
- 3. **Feature Extraction**: Convolution with specialized kernels like Gabor filters can extract features such as textures, orientations, and frequencies from images.
- 4. **Image Sharpening**: Convolution can enhance edges and details in images by emphasizing high-frequency components.

Properties of Convolution:

• **Linear Operation**: Convolution is a linear operation, meaning it satisfies the principles of superposition and homogeneity.

- **Translation Invariant**: Convolution is translation invariant, which means the same kernel can be applied at different locations in the image to achieve the same effect.
- **Commutativity**: The order of convolution with multiple kernels can be interchanged without affecting the final result.

In practical terms, convolution in computer vision is often implemented efficiently using algorithms such as Fast Fourier Transform (FFT) for large kernels to speed up the computation process. This allows real-time or near-real-time applications such as video processing and augmented reality. In summary, convolution in computer vision is a powerful technique used for filtering, feature extraction, and other image processing tasks by applying a kernel matrix over an image to modify its pixels based on their neighborhoods. It forms the basis for many algorithms and techniques that enhance the understanding and analysis of visual data.

Fourier Transform

Frequency in images is the rate of change of intensity values. Thus, a high-frequency image is the one where the intensity values change quickly from one pixel to the next. On the other hand, a low- frequency image may be one that is relatively uniform in brightness or were intensity changes very slowly. Frequency" means the rate of change of intensity per pixel. Let's say you have some region in your image that changes from white to black. If it takes many pixels to undergo that change, it's low frequency. The fewer the pixels it takes to represent that intensity variation, the higher the frequency.

- A Fourier Transform maps a signal into its component frequencies.
- It does not change the orginal signal, only its representation. It is an extremely useful operator used in many fields.
- In the context of images, the 2D Fourier Transform converts an image from its spatial domain (pixel intensity values) to its frequency domain (amplitudes and phases of different spatial frequencies).
- Fourier Transform is a mathematical model which helps to transform the signals between two different domains, such as transforming signal from frequency domain to time domain or vice versa.



$$f(x) = \int_{-\infty}^{\infty} F(u) e^{i2\pi u x} du$$

 Since of a real function is generally complex, we use magnitude and phase



Important Properties:

• Fourier transform is linear

$$F(g(x) + h(x)) = F(g(x)) + F(h(x))$$
$$F(kg(x)) = kF(g(x))$$

• FT and Convolution

$$f(x) * g(x) \Leftrightarrow F(u)G(u)$$
$$f(x)g(x) \Leftrightarrow F(u) * G(u)$$

• FT of a Gaussian is a Gaussian

Sampling and Aliasing

- Sampling and aliasing are fundamental concepts in computer vision and signal processing that are crucial for understanding how images and videos are captured, processed, and represented.
- Here's a breakdown of these concepts and their relevance to computer vision:

Sampling

- Sampling is the process of converting a continuous signal (such as a visual scene) into a discrete set of values.
- In the context of images, this involves measuring the intensity of light at specific locations on a grid.

Sampling in Images:

- **Pixels:** An image is sampled into a grid of pixels, where each pixel represents the color or intensity at that point.
- **Resolution:** The resolution of an image is determined by the number of pixels in the grid. Higher resolution means more pixels and more detailed information.

Sampling Rate:

• **Spatial Sampling Rate:** In images, this is the density of the pixel grid. Higher spatial sampling rates result in higher resolution images.

• **Temporal Sampling Rate:** In video, this refers to the number of frames per second (fps). Higher frame rates lead to smoother motion.

Aliasing: Aliasing occurs when a signal is sampled at a rate insufficient to capture its variations accurately, leading to distortions or artifacts.

Aliasing in Images:

• **Visual Artifacts:** When an image is under sampled, patterns that are finer than the pixel grid may not be represented accurately. This can lead to artifacts such as moiré patterns or jagged edges (also known as "jaggies").

• **Nyquist Theorem:** To avoid aliasing, the sampling rate should be at least twice the highest frequency present in the image (known as the Nyquist rate). For images, this means the pixel density should be high enough to capture the details.

To overcome Aliasing

1. Anti-Aliasing:

- **Filtering:** To reduce aliasing, anti-aliasing filters can be applied before sampling. For images, this often involves applying a low-pass filter to smooth out high-frequency details that could cause aliasing.
- **Subsampling:** Techniques such as super sampling or multisampling can help by sampling at a higher resolution and then averaging the results to reduce artifacts.

2. Image and Video Processing:

- **Resampling:** When resizing images or videos, interpolation methods (such as bilinear or bicubic interpolation) are used to estimate pixel values and reduce aliasing.
- **Image Enhancement:** Techniques like sharpening can be applied carefully to enhance details without introducing significant aliasing artifacts.
- sampling and aliasing are crucial for designing effective computer vision systems, as it affects image quality, processing algorithms, and overall system performance.
- Proper sampling ensures that images and videos are captured and processed with minimal distortion, leading to better visual information and more accurate analysis.

Filters as Templates

- In computer vision, filters are essential tools used to manipulate and analyze images.
- Filters as templates, are often referred to the concept of *template matching* or using filters in the context of convolutional operations.

1. Template Matching

- Template matching is a technique used to find a sub-region within an image that matches a given template. Here's how it works:
- Template: A small image or pattern that you want to locate in a larger image.
- Matching Process: The template is slid over the larger image (or vice versa) to check how well it matches different parts of the image. This is typically done using similarity measures like correlation or distance metrics.
- Result: The location where the template best matches the region in the image is identified.

• Template matching is commonly used in tasks such as object detection and facial recognition.

2. Convolutional Filters

- In the context of image processing, convolutional filters (or kernels) are used to apply various transformations to an image.
- These filters can be thought of as templates that are convolved (slid over) the image to produce different effects. Here's how convolution works:
- **Kernel/Filter**: A small matrix (template) used to perform convolution with the image. This matrix defines the type of transformation to apply, such as edge detection, blurring, or sharpening.
- **Convolution Operation**: The kernel is applied to every pixel in the image by computing a weighted sum of the pixel values within the kernel's area. This operation generates a new image based on the filter's effect.
- **Result**: The output is a transformed image where features such as edges, textures, or patterns are enhanced or detected.

Filters in Deep Learning

- In deep learning, particularly in Convolutional Neural Networks (CNNs), filters (or kernels) are used in convolutional layers to automatically learn features from images:
- **Training**: During training, CNNs learn the values of these filters by backpropagation, optimizing them to detect specific features like edges, textures, or more complex patterns.
- **Feature Extraction**: As the network layers process the image, different filters capture various hierarchical features, leading to robust and complex representations used for tasks like object recognition and classification.
- Filters as templates in computer vision are versatile tools used for image analysis and transformation. Whether through template matching or convolutional operations, they play a crucial role in understanding and manipulating visual data.

Correlation

• Correlation having similar operating steps as convolution operation which also takes an input image and another kernel and traverses the kernel window through the input by computing a weighted combination of pixel neighborhood values with the kernel values and producing the output image.

Correlation

$$\boldsymbol{g}(x, y) = \boldsymbol{f} \star \boldsymbol{K} = \sum_{u=-h}^{h} \sum_{v=-h}^{h} \boldsymbol{f}(x+u, y+v) \boldsymbol{K}(u, v)$$
Convolution
$$h \ge h \ker k \operatorname{kernel} \boldsymbol{K}$$

$$\boldsymbol{g}(x, y) = \boldsymbol{f} \ast \boldsymbol{K} = \sum_{u=-h}^{h} \sum_{v=-h}^{h} \boldsymbol{f}(x-u, y-v) \boldsymbol{K}(u, v)$$

• From above it is clear that only difference between correlation and convolution is that convolution flips the kernel twice (with regards to the horizontal and vertical axis) before computing the weighted combination.

• Whereas correlation maintains the original orientation of the kernel. It's often used in tasks like template matching, where the orientation of the kernel matters for finding occurrences of the template

Why correlation?

- **Commutative in Nature:** It means it does not distinguish between the input signal and the filter kernel, whereas convolution is not commutative. This means correlation is more suitable for tasks where directionality is not a concern, such as template matching or feature extraction.
- **Template matching capabilities:** Correlation is performer's better template matching as compared to convolution as template matching tasks, where a smaller template image is compared with regions of a larger image to find occurrences of the template. In template matching, the orientation of the template in the larger image may not matter, so correlation is preferred over convolution.
- **Cross-Correlation in Registration:** In image registration tasks, cross-correlation is often used to align two images. Cross-correlation measures the similarity between two images as one is shifted relative to the other. Convolution wouldn't be suitable for this task as it requires flipping one of the images, which isn't desirable for registration purposes.
- **Simplicity:** Sometimes, using correlation may lead to simpler algorithms or interpretations compared to convolution. For certain tasks, especially in introductory contexts, correlation may be easier to understand and implement.

Applying 1D correlation for simpler understanding

$$J(x) =_{(i = -N)^{N} H(i).I(x+i)}$$

Let image be I , I= [10, 20, 10, 50, 60]

Indexes of the image are 0, 1, 2, 3 and 4.

Filter be H, H = [1/3, 1/3, 1/3]

Indexes of the filter are -1, 0 and 1.

Apply correlation between image and mask at index=2 in the image.

 $J(2) = I(1) \cdot H(-1) + I(2) \cdot H(0) + I(3) \cdot H(1)$ Indexes are represented in the parentheses.

J(2) = 20 x 1/3 + 10 x 1/3 + 50 x 1/3

J(2) = 1/3(20+10+50)

J(2) = 80/3

Local Image Features

- Local image features are a crucial component in computer vision, especially for tasks like object recognition, image matching, and scene understanding.
- These features capture information about small regions of an image, which can then be used to identify and analyze patterns and structures within the image.

Key Concepts

- 1. **Feature Points**: These are distinctive points in an image that can be reliably detected and matched. They often correspond to corners, edges, or blobs.
- 2. **Descriptors**: Once feature points are identified, descriptors are used to describe the appearance of the local region around each feature point. These descriptors help in comparing and matching feature points across different images.

Detection and Description:

- **Detection**: Algorithms like Harris Corner Detector or the Scale-Invariant Feature Transform (SIFT) are used to identify key points in an image.
- **Description**: Once key points are detected, descriptors like SIFT, Speeded-Up Robust Features (SURF), or Binary Robust Invariant Scalable Key points (BRISK) provide a robust representation of the local area around these points.

Popular Algorithms

1. SIFT (Scale-Invariant Feature Transform):

- 1. **Detection**: Finds key points in the image and assigns them a scale and orientation.
- 2. **Description**: Generates a descriptor based on the local image gradients around each key point, which is invariant to scale and rotation.

2. SURF (Speeded-Up Robust Features):

- **Detection**: Similar to SIFT but optimized for speed using integral images and a different key point detector.
- **Description**: Uses a Hessian matrix-based approach for the descriptor, which is also scale and rotation invariant.

3.ORB (Oriented FAST and Rotated BRIEF):

- **Detection**: Uses the FAST key point detector to find features and then computes orientation.
- **Description**: Employs the BRIEF descriptor, modified to be rotation invariant, making it efficient and suitable for real-time applications.

4. AKAZE (Accelerated KAZE):

- **Detection**: Aimed at detecting features in nonlinear scale spaces, improving the detection of features in different scales.
- **Description**: Uses binary descriptors, which are more computationally efficient.

5. BRIEF (Binary Robust Independent Elementary Features):

• **Description**: A binary descriptor that uses a set of random intensity comparisons. It's fast but less robust compared to SIFT or SURF.

Computing the Image Gradient

• Computing image gradients provides essential **information about the structure and boundaries within an image**, making it a foundational element in many computer vision tasks.

Key Concepts

- 1. **Image Gradient**: The gradient of **an image at a point** represents the **direction and magnitude of the greatest rate of intensity change at that point**. It's essentially a vector that points in the direction of the most significant change in intensity.
- 2. Gradient Components:
 - Magnitude: Measures the strength of the intensity change. It is computed as $\sqrt{G_x^2 + G_{y'}^2}$ where G_x and G_y are the gradients in the x and y directions, respectively.
 - Direction: Indicates the direction of the intensity change, often given by atan2(G_y, G_x), which provides the angle of the gradient vector.

Computation Methods

- 1. **Finite Differences**: This method approximates the gradient by computing differences between neighboring pixel values.
- Forward Difference:
 - For G_x (horizontal gradient): $G_x = I(x+1,y) I(x,y)$
 - For G_y (vertical gradient): $G_y = I(x,y+1) I(x,y)$
- Central Difference: Provides a more accurate approximation by averaging the forward and backward differences:
 - For G_x : $G_x = \frac{I(x+1,y) I(x-1,y)}{2}$
 - For G_y : $G_y = \frac{I(x,y+1) I(x,y-1)}{2}$

- Convolution with Sobel Kernels: The Sobel operator is a common method for computing image gradients. It uses convolution with specific kernels to estimate the gradient.
 - Sobel Kernels:
 - Horizontal Gradient G_x :

$$\begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

• Vertical Gradient G_y :

$\left[-1\right]$	-2	-1
0	0	0
1	2	1

Convolution with these kernels provides approximations for G_x and G_y .

- Prewitt Operator: Another method similar to Sobel, but with different kernels. It also estimates gradients in the horizontal and vertical directions.
 - Prewitt Kernels:
 - Horizontal Gradient G_x:

$$\begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$$

• Vertical Gradient G_y :

[−1	$^{-1}$	-1
0	0	0
1	1	1

- 4. Roberts Cross Operator: A simpler method for edge detection, using a pair of 2x2 kernels.
 - Roberts Kernels:
 - Diagonal Gradient G_x:

 $\begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$

• Diagonal Gradient G_y :

 $\begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$

5. Gaussian Derivatives: For more robust gradient estimation, especially in noisy images, Gaussian smoothing can be applied before computing the gradient. This involves convolving the image with a Gaussian filter followed by applying gradient operators.

Applications

- **Edge Detection**: The gradient magnitude helps identify edges in an image, which are areas of high intensity change. Edge detection algorithms like Canny use gradients to find and refine edges.
- **Feature Detection**: Gradients are used in algorithms for detecting and describing local features, such as in the SIFT or SURF algorithms.
- **Image Segmentation**: Gradients can help segment an image into different regions based on intensity changes.
- **Image Enhancement**: Techniques like sharpening or contrast adjustment often rely on gradient information.

Gradient Based Edge Detectors

- Gradient-based edge detectors are essential tools in computer vision for identifying and locating edges within an image.
- These detectors rely on calculating the gradient of the image intensity to find areas of significant change, which typically correspond to edges.

1. Sobel Operator

• **Description**: The Sobel operator calculates the gradient of the image intensity at each pixel using convolution with Sobel kernels. It provides estimates of the gradient in the horizontal and vertical directions.

Kernels:

- Horizontal Gradient (G_x):
- Vertical Gradient (G_u):

-1	0	1
-2	0	2
-1	0	1

$\lceil -1 \rceil$	-2	-1
0	0	0
1	2	1

Process:

- 1. Convolve the image with the Sobel kernels to get G_x and G_y .
- 2. Compute the gradient magnitude: $\sqrt{G_x^2 + G_y^2}$.
- 3. Optionally, compute the gradient direction: $\operatorname{atan2}(G_y, G_x)$.
- Applications: Edge detection, image segmentation.

2. Prewitt Operator

- Description: Similar to the Sobel operator, the Prewitt operator also estimates image gradients using convolution with Prewitt kernels. It's simpler and less sensitive to noise compared to Sobel.
- Kernels:
 - Horizontal Gradient (G_x):

-1	0	1
-1	0	1
$\left\lfloor -1 \right\rfloor$	0	1

• Vertical Gradient (Gy):

$\left[-1\right]$	-1	-1
0	0	0
1	1	1

Process:

- 1. Convolve the image with Prewitt kernels to obtain gradients GxG_xGx and GyG_yGy.
- 2. Compute the gradient magnitude and direction.
- **Applications**: Edge detection, basic image processing tasks.
- 3. Roberts Cross Operator
- **Description**: The Roberts Cross operator uses diagonal kernels to detect edges. It's particularly good at detecting edges in images with sharp intensity changes.
 - Kernels:
 - Diagonal Gradient (G_x):
 - Diagonal Gradient (G_y):

 $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

 $\begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$

 $\begin{bmatrix} 0\\ -1 \end{bmatrix}$

4. Canny Edge Detector

- Description: The Canny edge detector is a multi-stage algorithm that is more sophisticated and provides more accurate edge detection by using gradient information in combination with additional steps like Gaussian smoothing and non-maximum suppression.
- Steps:
 - 1. Smoothing: Apply a Gaussian filter to reduce noise.
 - 2. Gradient Computation: Use Sobel operators to find gradient magnitude and direction.
 - Non-Maximum Suppression: Thin the edges by suppressing non-maximum gradient magnitudes.
 - Edge Tracking by Hysteresis: Apply double thresholding to identify strong and weak edges and then track edges by connecting weak edges that are connected to strong edges.

5. Laplacian of Gaussian (LoG) Detector

- **Description**: The LoG operator detects edges by first smoothing the image with a Gaussian filter and then applying the Laplacian operator to detect zero-crossings where the intensity changes sign.
- Process:
 - **Smoothing**: Convolve the image with a Gaussian filter.
 - Laplacian: Apply the Laplacian operator to the smoothed image to detect edges.
 - Zero-Crossing: Identify edges as zero-crossings in the Laplacian response.
- **Applications**: Edge detection, particularly in noisy images.

6. Gaussian Derivative Filters

- **Description**: These filters compute image gradients by convolving the image with Gaussian derivative filters.
- They are useful for capturing edges at different scales and are less sensitive to noise compared to the Sobel operator.

Filters:

- Horizontal Derivative (G_x): Gaussian filter derivative in the x direction.
- Vertical Derivative (G_y): Gaussian filter derivative in the y direction.
- Process:
 - 1. Convolve the image with the Gaussian derivative filters.
 - 2. Compute gradient magnitude and direction.

Orientations

Orientation refers to **the direction of an edge or feature within an image**. Understanding orientation is crucial for tasks such as edge detection, object recognition, and image alignment.

1. Orientation of Edges:

- Definition: The orientation of an edge is the angle of the edge relative to a reference direction, typically the horizontal axis. It indicates the direction in which the intensity gradient changes the most.
- Calculation: For a given edge, the orientation can be calculated using the gradient components G_x (horizontal gradient) and G_y (vertical gradient). The orientation angle θ is given by:

$$\theta = \operatorname{atan2}(G_y, G_x)$$

Here, atan2 is a function that computes the angle from the arctangent, taking into account the sign of both components to determine the correct quadrant.

- 2. Gradient Magnitude:
 - Definition: The magnitude of the gradient represents the strength of the edge. It is calculated as:

$$ext{Magnitude} = \sqrt{G_x^2 + G_y^2}$$

Computing Orientation with Different Methods

- 1. Sobel Operator:
 - Process: Convolve the image with Sobel kernels to compute G_x and G_y. Calculate the orientation using the formula:

$$\theta = \operatorname{atan2}(G_y, G_x)$$

- 2. Canny Edge Detector:
 - Process: Compute the gradient magnitude and direction at each pixel using the Sobel operators or other gradient operators. Apply non-maximum suppression to thin out the edges, preserving the direction information.

- 3. Harris Corner Detector:
 - Orientation Calculation: The Harris detector is used primarily for corner detection but can also provide orientation information at corners based on the eigenvalues of the structure tensor.
- 4. Gabor Filters:
 - Process: Gabor filters can be used to detect edges and their orientations by convolving the image with Gabor kernels oriented at different angles. This helps in capturing texture and orientation information in different frequency bands.

Texture

• In computer vision, "texture" refers to the surface quality or pattern of an object that can be described by its spatial arrangement of colors, intensities, or patterns.

• Texture is crucial in various computer vision tasks because it provides information that helps in object recognition, classification, and scene understanding.

Some key concepts related to texture in computer vision are:

1. Texture Features

• **Statistical Features:** These include properties like mean, variance, skewness, and kurtosis of pixel intensity values. The Gray-Level Co-occurrence Matrix (GLCM) is often used to derive statistical features such as contrast, correlation, energy, and homogeneity.

• **Structural Features:** These describe the arrangement of texture elements or patterns. Examples include the size, shape, and orientation of texture elements.

• **Spectral Features:** Texture can be analyzed in the frequency domain using techniques like Fourier Transform, where different textures are represented by their frequency components.

Texture Analysis Techniques

- **Gray-Level Co-occurrence Matrix (GLCM):** Computes how often pairs of pixels with specific values occur in a specified spatial relationship. This matrix can be used to extract texture features.
- Local Binary Patterns (LBP): A simple yet effective texture descriptor that compares each pixel with its neighbors to form a binary pattern. These patterns are then used to characterize texture.
- **Gabor Filters:** These are used to analyze textures by convolving the image with a set of filters at different orientations and frequencies.
- **Wavelet Transform:** This method breaks down an image into components at different scales, capturing texture information at multiple levels.

Applications of Texture Analysis

- **Object Recognition:** Identifying and classifying objects based on their surface patterns.
- **Image Segmentation:** Dividing an image into regions with similar texture characteristics.
- **Medical Imaging:** Analyzing textures in medical scans (e.g., MRI, CT) to detect abnormalities or diagnose diseases.
- **Quality Control:** Inspecting surfaces in manufacturing to detect defects or inconsistencies.

Local Texture Representations Using Filters

- Local texture representations using filters are a fundamental concept in image processing and computer vision.
- They involve analyzing and extracting features from images by applying different filters.
- These filters help to capture various aspects of the local texture, such as edges, patterns, and variations in intensity. Here's a breakdown of the key ideas and methods involved:

1. Purpose of Local Texture Representations

- Local texture representations are used to understand and describe the texture patterns within a small region of an image. This can be crucial for various applications such as:
- Image Segmentation: Differentiating between different regions based on texture.
- **Object Recognition:** Identifying objects based on their texture patterns.
- **Image Classification:** Categorizing images based on texture characteristics.

2. Types of Filters Used

1. **Convolutional Filters:** These are small matrices (kernels) that are convolved with the image to detect specific features. Common types include:

- **Sobel Filters:** Used for edge detection by emphasizing gradients in the horizontal and vertical directions.
- **Gaussian Filters:** Smooth the image by blurring, which helps to reduce noise and detail.
- Laplacian Filters: Detects regions of rapid intensity change and is used for edge detection.

2. **Gabor Filters:** These are used for texture representation and analysis. They are particularly effective in capturing frequency and orientation information from the image. Gabor filters can be seen as a bank of filters with different frequencies and orientations.

3. Local Binary Patterns (LBP): LBP is a simple yet effective method for texture classification. It compares each pixel with its neighbors and encodes the result as a binary number. This captures local texture information by focusing on the contrast between neighboring pixels.

4. **Gray Level Co-occurrence Matrix (GLCM):** This method analyzes the spatial relationship between pixels. It computes various texture features like contrast, correlation, and homogeneity by looking at how often pixel pairs with specific values occur in the image.

5. **Wavelet Transforms:** These are used to analyze texture at multiple scales. Wavelet transforms decompose an image into different frequency components, allowing for multi-resolution analysis of textures.

<u>Shape from Texture</u>: "Shape from texture" is a technique in computer vision and image processing used to infer the 3D shape of an object based on its texture patterns.
- This method relies on the understanding that textures can provide clues about the underlying surface structure. Here's a comprehensive look at the concept and how it works:
- Textures in images are not uniformly distributed; they often follow certain patterns that change with the surface's orientation and curvature.
- By analyzing these texture patterns, we can infer information about the 3D shape of the surface on which the texture is applied.

Key Principles

- 1. **Texture Gradient**: As an object's surface curves or tilts, the texture pattern changes in a predictable way. A texture gradient refers to the change in texture pattern over distance. By analyzing these gradients, you can estimate surface orientation.
- 2. **Texture Compression and Expansion**: When a surface is sloped, textures appear compressed or expanded. For example, texture elements might appear larger or smaller, or more densely packed, as the surface tilts away from the viewer. By measuring these changes, one can infer surface curvature.
- 3. **Surface Orientation**: The orientation of the surface can be inferred by observing how texture elements align. If a texture is applied to a flat surface, its orientation remains consistent, but if the surface is curved or inclined, the texture will deform accordingly.

Techniques for Shape from Texture

1. Pattern Analysis:

- **Parallel Texture Lines**: If textures are parallel lines, their apparent convergence or divergence can indicate surface tilt or curvature.
- **Repeated Patterns**: If the texture consists of repeating elements, changes in their apparent size or spacing can indicate changes in surface depth or orientation.

2. **Texture Gradient Method**:

- **Texture Gradient Extraction**: Compute gradients in the texture pattern to understand how the texture changes across the surface. This involves detecting variations in texture elements and their relative distances.
- **Depth Estimation**: Use these gradients to estimate the depth of the surface at different points. This can be done by comparing the observed gradient with the expected gradient from a known texture pattern.

3. Photometric Stereo:

• **Multiple Light Sources**: Capture multiple images of the same scene under different lighting conditions. Changes in texture shading due to lighting variations can help infer surface normal and, subsequently, the shape of the object.

4. Structure from Motion (SfM) with Texture:

• **3D Reconstruction**: Combine texture information with camera motion data to reconstruct the 3D shape of an object. This technique involves capturing multiple views of the object and using the texture patterns to assist in depth estimation.

5. Machine Learning Approaches:

• **Deep Learning Models**: Train NN's to recognize texture patterns & infer shape information. CNNs are used to learn texture features & their relationship to 3D shapes from large datasets.

UNIT-III

MID-LEVEL VISION: Segmentation by Clustering - Basic Clustering Methods, The Watershed Algorithm, Segmentation Using K-means, Grouping and Model Fitting - Fitting Lines with the Hough Transform, Fitting Curved Structures, Tracking - Tracking by Detection, Tracking Translations by Matching, Tracking Linear Dynamical Models with Kalman Filters.

Mid-Level Vision

- In computer vision, "mid-level vision" refers to the processing and analysis of visual information that falls between low-level features and high-level object recognition.
- It involves intermediate stages of visual understanding that help bridge the gap between raw pixel data and high-level semantic interpretation.
- Mid-level vision typically includes the following:
- 1. Feature Extraction: This involves identifying and describing significant features in an image, such as edges, textures, and corners. Techniques like edge detection (e.g., Canny edge detector) and feature descriptors (e.g., SIFT, SURF) are often used.
- 2. Segmentation: This process divides an image into meaningful regions or segments, often based on characteristics like color, texture, or intensity. Common methods include thresholding, clustering (e.g., k-means), and more advanced techniques like semantic segmentation using neural networks.
- 3. Object Detection and Localization: This step involves identifying objects within an image and determining their locations. It includes techniques such as bounding box detection and region proposals.
- 4. Pattern Recognition: At this level, patterns within the segmented regions or objects are analyzed. This might involve recognizing shapes, textures, or repetitive patterns that are relevant for further analysis.
- 5. Scene Understanding: This involves analyzing the spatial relationships and contextual information between different objects in an image. Techniques might include understanding object interactions or scene layout.
- 6. Feature Matching: This is the process of matching features across different images or within the same image to establish correspondences. This can be crucial for tasks like image stitching, 3D reconstruction, or tracking.

Segmentation by Clustering

- Segmentation by clustering is a popular mid-level vision technique used to partition an image into regions or segments based on similarity.
- The core idea is to group pixels that are similar to each other into clusters, with the goal of identifying meaningful structures or regions within the image. Here's how it typically works:

1. Feature Extraction:

- **Color Features:** Extract color information from each pixel, which could be in RGB, HSV, or another color space.
- **Texture Features:** Analyze texture patterns, which can be captured using methods like Local Binary Patterns (LBP) or Gray-Level Co-occurrence Matrix (GLCM).
- **Other Features:** Depending on the application, you might use other features like intensity, gradients, or more complex descriptors.

2. Choosing a Clustering Algorithm:

- Several clustering algorithms can be used for image segmentation including
- **K-Means Clustering:**
- **Algorithm:** K-Means partitions the pixels into K clusters by minimizing the variance within each cluster.
- **Process:** Initialize K cluster centroids, assign each pixel to the nearest centroid, update centroids based on the mean of assigned pixels, and repeat until convergence.
- **Advantages:** Simple and effective for images where regions have distinct color or texture differences.
- **Disadvantages:** Requires the number of clusters K to be specified and can be sensitive to initialization.

Mean Shift Clustering:

- **Algorithm:** Mean Shift finds clusters by iteratively shifting data points towards the mean of their local neighborhood.
- **Process:** Start with a set of points, shift each point towards the mean of its neighborhood, and update clusters based on the convergence of points.
- **Advantages:** Does not require specifying the number of clusters beforehand and can handle varying cluster shapes.
- **Disadvantages:** Can be computationally expensive and may require careful tuning of parameters like bandwidth.

Gaussian Mixture Models (GMM):

- **Algorithm:** GMM assumes data is generated from a mixture of several Gaussian distributions and uses Expectation-Maximization (EM) to find the best fit.
- **Process:** Estimate the parameters of the Gaussian distributions and assign pixels to clusters based on their probabilities.
- Advantages: Flexible and can model more complex distributions.
- **Disadvantages:** Requires specifying the number of components and can be sensitive to initialization.
- **DBSCAN (Density-Based Spatial Clustering of Applications with Noise):**
- **Algorithm:** DBSCAN groups pixels based on density, identifying clusters of varying shapes and sizes.
- **Process:** Define core points with a minimum number of neighbors and expand clusters from these core points.
- Advantages: Can find clusters of arbitrary shapes and handle noise.
- **Disadvantages:** Sensitive to parameters like the radius and minimum points, and might struggle with varying densities.

3. Post-Processing:

- **Smoothing:** Apply smoothing techniques to refine the boundaries between segments and reduce noise.
- **Merging/Splitting:** Adjust the segments by merging similar regions or splitting large segments if necessary.

• **Edge Detection:** Sometimes combined with edge detection methods to enhance segment boundaries.

4. Evaluation:

- **Quantitative Measures:** Assess segmentation quality using metrics like the Davies-Bouldin index, Silhouette score, or the Adjusted Rand Index.
- **Visual Inspection:** Check the segmented regions visually to ensure they align with the expected structures or regions in the image.

Segmentation by clustering is versatile and can be adapted to various types of images and applications. The choice of clustering algorithm and features will depend on the specific characteristics of the images and the goals of the segmentation task.

The Watershed Algorithm

- The Watershed algorithm is a powerful technique used in mid-level vision for image segmentation.
- It's particularly effective for segmenting images where objects or regions are separated by varying intensity levels.
- The segmentation process will take the similarity with adjacent pixels of the image as an important reference to connect pixels with similar spatial positions and gray values.
- The algorithm is inspired by the concept of watershed in geography, where water flowing from different sources eventually meets at certain points, forming distinct basins.



Concept and Workflow

1. Gradient Computation:

• The Watershed algorithm starts by computing the gradient of the image. The gradient represents the intensity changes and highlights the boundaries between different regions.

• Typically, a gradient image is obtained using edge detection methods like Sobel or using morphological gradient operations. This gradient image helps to find the edges of objects in the image.

2. Markers Initialization:

- Markers are initial seeds or points that indicate the starting regions for segmentation. They are placed in regions that are known to belong to specific segments or objects.
- Markers can be manually defined, automatically generated using techniques like connected component analysis, or determined through other methods like thresholding.

3. Flooding Process:

- Imagine the gradient image as a topographic surface with varying heights. The idea is to simulate "flooding" the surface from the marker points.
- Starting from each marker, the algorithm floods the gradient image, expanding the regions until different markers meet or boundaries are encountered. As the water rises, it fills up the different basins, which correspond to the segmented regions.

4. Region Boundaries:

- The boundaries where different markers meet during the flooding process form the segmented regions.
- The final result is an image where different regions are separated by watershed lines, indicating the boundaries of different objects or segments.

Steps in Detail

1. Preprocessing:

- Convert the image to grayscale if it's in color, as the Watershed algorithm typically operates on single-channel images.
- Apply a smoothing filter (like Gaussian blur) to reduce noise and improve segmentation quality.

2. Gradient Computation:

• Compute the gradient magnitude of the image to highlight edges. This step is crucial as it helps to identify the boundaries between regions.

3. Markers Placement:

• Create markers for the foreground (objects of interest) and background. Foreground markers are placed inside the objects, and background markers are placed in areas that should be segmented as background.

4. Applying the Watershed Algorithm:

• Use the markers and gradient image to run the Watershed algorithm. This process will segment the image into different regions based on the gradient and marker information.

5. Post-processing:

• Clean up the segmented regions by removing small, noisy segments or merging regions if necessary. Morphological operations can be applied to refine the segmentation results.

Segmentation Using K-means

- Segmentation using K-means clustering is a widely used technique in mid-level vision for partitioning images into distinct regions based on color, texture, or other features.
- The K-means algorithm is simple yet effective, particularly when the number of segments is known or can be reasonably estimated.

How K-Means Clustering Works

- 1. Initialization:
 - 1. **Select K:** Choose the number of clusters (K) you want to partition the image into. This number needs to be specified before running the algorithm.
 - 2. **Initialize Centroids:** Randomly select K initial centroids (or cluster centers) from the image data. These centroids represent the initial guesses for the center of each cluster.

2. Assignment Step:

- **Assign Pixels to Clusters:** For each pixel (or feature vector) in the image, calculate its distance to each centroid. The distance metric is usually Euclidean distance.
- **Cluster Assignment:** Assign each pixel to the cluster whose centroid is closest to it. This step effectively partitions the image into K regions based on the initial centroids.

3. Update Step:

• **Recompute Centroids:** After all pixels have been assigned to clusters, recalculate the centroids of each cluster. The new centroid is typically the mean of all pixels assigned to that cluster.

4. Repeat Steps:

• **Iterate:** Repeat the assignment and update steps until the centroids no longer change significantly, or a maximum number of iterations is reached. Convergence is achieved when pixel assignments stabilize and centroids are consistent.

5. Output:

• **Segmented Image:** The result is an image where each pixel is assigned a cluster label, effectively segmenting the image into K regions. Each region corresponds to a cluster of similar pixels based on the chosen features.

Steps for Image Segmentation Using K-Means

1. Preprocessing:

- **Convert to Suitable Feature Space:** If the image is in color, convert it to a suitable feature space like RGB or HSV. You may also choose to work with grayscale images, depending on the application.
- **Flatten Image:** Convert the 2D image into a 1D array of pixel values or feature vectors. Each pixel is treated as a data point with its color (or other features) as coordinates.

2. Apply K-Means Algorithm:

• Select Number of Clusters (K): Choose the number of desired segments or regions.

• **Run K-Means:** Apply the K-means algorithm to the pixel data. This involves initialization, assignment, and update steps as described above.

3. Reshape Output:

- **Reshape Labels:** Convert the cluster labels back into the 2D image shape. Each pixel's label indicates the segment to which it belongs.
- **Visualize Segmentation:** Display or visualize the segmented image, where each segment is represented by a different color or intensity.

4. Post-Processing

- **Refine Segmentation:** Apply additional processing to refine the segmentation results. Techniques like morphological operations (e.g., erosion, dilation) can help clean up the segmented regions and remove small artifacts.
- **Evaluate Results:** Assess the quality of segmentation using metrics such as silhouette score or by visual inspection.

K-means clustering provides a foundational method for image segmentation, offering a balance of simplicity and effectiveness. It is particularly useful when working with images where regions can be approximated as clusters in feature space.

Grouping and Model Fitting

• In computer vision, grouping and model fitting are crucial techniques for understanding and interpreting visual data.

1. Grouping

- Grouping in computer vision refers to the process of clustering or segmenting similar features or objects in an image to simplify analysis.
- This is fundamental for tasks like object detection, recognition, and scene understanding. Key methods include:

a. Segmentation

- **Thresholding:** Simple method where pixels are grouped based on intensity or color values.
- **Clustering-based:** Techniques like K-means or Mean Shift that group pixels with similar attributes.
- **Region-based:** Methods that group neighboring pixels based on similarity criteria (e.g., Region Growing).
- **Graph-based:** Algorithms like Normalized Cuts or Graph Cuts that partition the image into segments based on graph theory.

b. Feature Grouping

- **Feature Matching:** Involves grouping key points or descriptors (e.g., SIFT, SURF) from different images to find correspondences.
- **Object Proposal:** Techniques like Selective Search that generate candidate regions where objects might be located and group these regions for further analysis.

c. Clustering

- K-means Clustering: Groups data into K clusters based on feature similarity.
- **Hierarchical Clustering:** Builds a hierarchy of clusters using either agglomerative (bottom-up) or divisive (top-down) methods.
- **DBSCAN (Density-Based Spatial Clustering of Applications with Noise):** Groups points that are close to each other while marking outliers.

2. Model Fitting

• Model fitting involves finding the best model that describes or fits the observed data. In computer vision, this typically means fitting geometric models to data points or features. Key aspects include:

a. Geometric Transformations

- **Homography:** A transformation used to map points from one plane to another, useful in tasks like image stitching.
- **Affine Transformation:** Used for modeling transformations like rotation, translation, and scaling.
- **Perspective Transformation:** Models how a 3D scene is projected onto a 2D image.

b. Parameter Estimation

- **Least Squares:** Commonly used to minimize the difference between the observed data and the model's predictions.
- **RANSAC (Random Sample Consensus):** A robust method to estimate parameters by iteratively fitting the model to random subsets of data and identifying inliers.
- **Maximum Likelihood Estimation (MLE):** Estimating parameters by maximizing the likelihood of the observed data given the model.

c. Object Detection and Recognition

- **Template Matching:** Fitting a pre-defined template to parts of the image to detect objects.
- **Machine Learning Models:** Using classifiers (e.g., SVM, neural networks) trained on annotated data to recognize objects.

d. Deep Learning Approaches

- **Convolutional Neural Networks (CNNs):** Automatically learn features and fitting functions from raw image data, widely used in modern computer vision tasks.
- **Transformers:** Emerging models for vision tasks that capture long-range dependencies and relationships in images.

Fitting lines with Hough transform

- The Hough Transform is a pivotal algorithm in computer vision and image processing, enabling the detection of geometrical shapes such as lines, circles, & ellipses within images.
- By transforming image space into parameter space, the Hough Transform leverages a voting mechanism to identify shapes through local maxima in an accumulator array.
- Typically, this method detects lines and edges, utilizing parameters like rho and theta to represent straight lines in polar coordinates.

• This algorithm is essential in various applications, from edge detection and feature extraction to more complex tasks like circle detection and generalized shape identification.

Why is it Needed?

- In many circumstances, a pre-processing stage can use an edge detector to obtain picture points or pixels on the required curve in the image space.
- However, there may be missing points or pixels on the required curves due to flaws in either the image data or the edge detector and spatial variations between the ideal line/circle/ellipse and the noisy edge points acquired by the edge detector.
- As a result, grouping the extracted edge characteristics into an appropriate collection of lines, circles, or ellipses is frequently difficult.



1. Original image of Lane



Figure 2: Image after applying edge detection technique. Red circles show that the line is breaking there.



Dept. of CSE AIML/B.Tech AIML

How Does the Hough Transform Work?

- The **accumulator array**, sometimes referred to as the **parameter space or Hough space**, is the first thing that the Hough Transform creates.
- The available parameter values for the shapes that are being detected are represented by this space.
- The **slope (m) and y-intercept (b) of a line**, for instance, could be the parameters in the line detection scenario.
- The **Hough Transform calculates the matching curves in the parameter space** for each edge point in the image.
- This is accomplished by finding the curve that intersects the parameter values at the spot by iterating over all possible values of the parameters.
- **The "votes" or intersections for every combination** of parameters are recorded by the accumulator array.
- In the end, **the program finds peaks in the accumulator array that match the parameters of the shapes it has identified**. These peaks show whether the image contains lines, circles, or other shapes.



• A line can be described analytically in various ways. One of the line equations uses the parametric or normal notion: $x \cos \theta + y \sin \theta = r$. where r is the length of a normal from the origin to this line and θ is the orientation



- The known variables (i.e., x_i, y_i) in the image are constants in the **parametric line** equation, whereas r and are the unknown variables we seek.
- Points in cartesian image space correspond to curves (i.e., sinusoids) in the polar Hough parameter space if we plot the potential (r, θ) values specified by each.
- The Hough Transform algorithm for straight lines is this point-to-curve transformation.
- Collinear spots in the cartesian image space become obvious when examined in the Hough parameter space because they provide curves that overlap at a single (r, θ) point.
- A and b are the circle's center coordinates, and r is the radius.
- The algorithm's computing complexity increases because we now have three coordinates in the parameter space and a 3-D accumulator.
- (In general, the number of parameters increases the calculation and the size of the accumulator array polynomial) As a result, the fundamental Hough approach described here only applies to straight lines.

Fitting Curved Structures

- Fitting curved structures in mid-level vision involves understanding and processing visual information to recognize and model shapes that are not strictly linear.
- Mid-level vision bridges the gap between low-level processes (like edge detection) and high-level object recognition. Here's a general overview of how fitting curved structures is approached in mid-level vision:

1. Curve Detection and Representation

- **Edge Detection:** Initial steps often involve detecting edges using methods like the Canny edge detector or the Sobel operator. These edges provide the starting points for identifying curves.
- **Curve Models:** Common models for representing curves include splines (like B-splines or cubic splines), parametric curves (such as Bézier curves), and mathematical functions (like circle or ellipse equations).

2. Curve Fitting Techniques

- **Least Squares Fitting:** This involves minimizing the difference between the observed data points and the curve model. For example, fitting a quadratic or cubic polynomial to a set of edge points.
- **RANSAC (Random Sample Consensus):** A robust method that can handle outliers. It iteratively fits a model to a subset of the data and evaluates its consistency with the rest of the data.
- **Hough Transform:** A popular technique for detecting curves by transforming the edge points into a parameter space and identifying curves through accumulator cells.

3. Segmentation and Grouping

- **Region Growing:** This technique groups pixels or edges into segments based on similarity criteria and can be used to fit curves within these segments.
- Active Contours (Snakes): A method where a curve evolves to fit the boundary of an object. The curve adjusts iteratively based on energy minimization principles, considering both internal smoothness and external image features.

4. Feature Integration

- **Hierarchical Approaches:** Mid-level vision often involves integrating features from multiple levels. For instance, combining detected curves with texture information or contextual cues to improve the accuracy of curve fitting.
- **Contextual Information:** Leveraging the spatial arrangement and relationships between different curves to enhance the fitting process. This can involve incorporating prior knowledge about object shapes or structures.

Tracking - Tracking by Detection

- In computer vision, "Tracking by Detection" (TbD) is a method used to track objects in video sequences where the object is first detected in each frame, and then tracking algorithms are used to maintain the identity of the object across frames.
- This approach contrasts with "Tracking by Matching" methods, where the focus is on continuously updating object locations without explicit re-detection.
- This is how Tracking by Detection works:

1. Detection

a. Object Detection:

- <u>Initial Detection</u>: Each frame of the video is processed using an object detection algorithm to identify and localize objects. Popular algorithms include YOLO (You Only Look Once), SSD (Single Shot Multi Box Detector), and Faster R-CNN. These algorithms provide bounding boxes and class labels for detected objects.
- <u>Feature Extraction</u>: Key features are extracted from the detected objects to aid in tracking. This might include color histograms, texture patterns, or deep learned embeddings.

2. Tracking

a. Initialization:

• Object Initialization: The detected objects from the first frame are used to initialize the tracking process. Each object's position and appearance are recorded.

b. Tracking Algorithms:

- <u>Kalman Filter</u>: Used to predict the future position of objects based on their current trajectory and movement model. It helps to handle noise and small occlusions.
- <u>Particle Filter</u>: Utilizes a set of samples (particles) to represent the possible states of an object. It's useful for complex motion models and scenarios where objects undergo significant changes.
- <u>Correlation Filters:</u> Track objects by comparing the appearance of the object in the current frame with a reference appearance model. Methods like KLT (Kanade-Lucas-Tomasi) tracker use this approach.
- <u>Optical Flow:</u> Computes the apparent motion of objects by analyzing changes in pixel intensity between frames. It's useful for tracking objects based on how they move across the image.

3. Integration of Detection and Tracking

a. Re-detection:

• <u>Re-detection Mechanism</u>: Periodically, the object detection algorithm is applied again to correct tracking errors, handle drift, or re-acquire objects that have been lost temporarily.

b. Data Association:

• <u>Matching</u>: Detected objects are matched to existing tracks based on features, appearance, and proximity. Techniques such as the Hungarian algorithm for assignment or metrics like Intersection over Union (IoU) are used.

c. Track Management:

- <u>Track Initialization</u>: New tracks are created for newly detected objects that do not match existing tracks.
- <u>Track Termination</u>: Tracks are terminated if the object is no longer detected or has been lost for an extended period.

4. Applications

- <u>Surveillance Systems:</u> Tracking people or vehicles in security footage.
- <u>Autonomous Vehicles:</u> Monitoring pedestrians, other vehicles, and obstacles for navigation.
- <u>Sports Analytics:</u> Tracking players and the ball for performance analysis.
- <u>Augmented Reality:</u> Tracking objects to overlay digital information.

5. Challenges

- <u>Occlusion</u>: Objects may be blocked by other objects or people, making detection and tracking difficult.
- <u>Appearance Changes:</u> Variations in lighting, scale, or object deformation can impact tracking accuracy.
- <u>Computational Complexity:</u> Real-time processing requires efficient algorithms and significant computational resources.
- <u>Drift:</u> Over time, tracking accuracy might degrade, especially if the initial detection was not perfect.

Tracking Translations by Matching

- In computer vision, "Tracking Translations by Matching" is a method used to follow the movement of objects in a video by continuously updating their positions based on matching techniques.
- This approach focuses on tracking objects by comparing their appearance across frames, assuming the object's motion primarily involves translations (i.e., moving in a straight line without significant changes in shape or orientation).
- Here's a detailed overview:

1. Initialization

a. Object Initialization:

• Initial Detection: The object to be tracked is first detected in the initial frame. This can be done using object detection algorithms or manual selection. The object is usually represented by a bounding box or a mask.

b. Feature Extraction:

• Appearance Model: Features are extracted from the object to create an appearance model. This can include color histograms, texture patterns, or features from a pre-trained deep learning model (e.g., using convolutional neural networks).

2. Tracking Process

a. Template Matching:

- Template Creation: A template or reference image of the object is created from the initial frame. This template captures the object's appearance and is used for matching in subsequent frames.
- Template Matching Techniques: Methods like cross-correlation, normalized crosscorrelation, or more advanced techniques like the Sum of Absolute Differences (SAD) or the Sum of Squared Differences (SSD) are used to find the best match for the template in each new frame.

b. Tracking Algorithms:

- Correlation Filters: Algorithms like the Discriminative Correlation Filter (DCF) are used to track objects by learning an appearance model and finding the best match in each frame. The object's position is updated based on the location with the highest correlation score.
- Mean-Shift and CAMShift: Mean-Shift tracking iteratively moves the object's position to the peak of the color histogram or probability distribution in each frame. CAMShift (Continuously Adaptive Mean-Shift) adapts the window size and orientation to handle scale and rotation variations.

c. Motion Models:

- **Constant Velocity Model**: Assumes the object moves with a constant velocity. This can be used to predict the object's position in future frames and improve tracking robustness.
- **Translation Only Model**: Assumes the object primarily undergoes translations with minimal changes in shape or appearance.

3. Integration and Refinement

a. Error Correction:

- **Error Handling**: If tracking errors occur (e.g., due to occlusion or appearance changes), methods like re-initialization or re-detection of the object can be employed to correct errors.
- **Drift Correction**: Techniques like re-aligning the template or adjusting the appearance model can help correct tracking drift over time.

b. Adaptive Models:

• **Model Update**: The appearance model is updated periodically to adapt to changes in the object's appearance due to lighting, scale, or deformation. This helps maintain tracking accuracy over longer sequences.

4. Applications

- **Object Tracking in Videos**: Used in various applications, including surveillance, video analysis, and augmented reality.
- **Human-Computer Interaction**: Tracking hand movements or facial expressions for gesture-based control or emotion recognition.
- **Sports Analytics**: Monitoring player movements and ball trajectory for performance analysis.

Tracking Linear Dynamical Models with Kalman Filters

- Tracking objects using Linear Dynamical Models with Kalman Filters is a powerful technique in computer vision and control systems.
- This method leverages the Kalman Filter algorithm to estimate the state of a dynamic system (such as the position and velocity of an object) over time, based on noisy and partial observations.
- Here's a detailed explanation of how it works:

1. Linear Dynamical Models

State Representation:

State Transition Model:

• Linear Transition: Describes how the state evolves over time. It can be represented as:

$$\mathbf{x}_{k+1} = \mathbf{F}\mathbf{x}_k + \mathbf{w}_k$$

where \mathbf{F} is the state transition matrix that models the dynamics of the system (e.g., constant velocity model), and \mathbf{w}_k is the process noise, assumed to be Gaussian with zero mean and covariance \mathbf{Q} .

Observation Model:

• Measurement Model: Describes how the observations (or measurements) are related to the state. It can be represented as:

$$\mathbf{z}_k = \mathbf{H}\mathbf{x}_k + \mathbf{v}_k$$

where \mathbf{H} is the measurement matrix that maps the state vector to the observation vector, and \mathbf{v}_k is the measurement noise, assumed to be Gaussian with zero mean and covariance \mathbf{R} .

3. Kalman Filter Algorithm

• The Kalman Filter algorithm consists of two main steps: Prediction and Update.

Prediction Step:

• State Prediction: Predict the next state based on the current state and the state transition model:

$$\hat{\mathbf{x}}_{k+1|k} = \mathbf{F}\hat{\mathbf{x}}_{k|k}$$

where $\hat{\mathbf{x}}_{k|k}$ is the state estimate at time k given all observations up to k.

• Covariance Prediction: Predict the error covariance for the next state:

$$\mathbf{P}_{k+1|k} = \mathbf{F}\mathbf{P}_{k|k}\mathbf{F}^T + \mathbf{Q}$$

where $\mathbf{P}_{k|k}$ is the error covariance matrix at time k and \mathbf{Q} is the process noise covariance.

Update Step:

• Measurement Residual: Calculate the difference between the actual measurement and the predicted measurement:

$$\mathbf{y}_k = \mathbf{z}_k - \mathbf{H}\hat{\mathbf{x}}_{k|k-1}$$

• Kalman Gain: Compute the Kalman Gain, which balances the importance of the prediction and the measurement:

$$\mathbf{K}_k = \mathbf{P}_{k|k-1}\mathbf{H}^T(\mathbf{H}\mathbf{P}_{k|k-1}\mathbf{H}^T + \mathbf{R})^{-1}$$

• State Update: Update the state estimate with the measurement residual:

$$\hat{\mathbf{x}}_{k|k} = \hat{\mathbf{x}}_{k|k-1} + \mathbf{K}_k \mathbf{y}_k$$

• Covariance Update: Update the error covariance:

$$\mathbf{P}_{k|k} = (\mathbf{I} - \mathbf{K}_k \mathbf{H}) \mathbf{P}_{k|k-1}$$

3. Applications

- Object Tracking: Predict and track the position and velocity of objects in video sequences or real-time systems.
- Navigation Systems: Estimate the position and velocity of vehicles or aircraft.
- Control Systems: Apply in robotics for trajectory estimation and control.

4. Challenges and Considerations

- Linear Assumption: The Kalman Filter assumes linear dynamics and Gaussian noise. If the system or noise is highly non-linear, extensions like the Extended Kalman Filter (EKF) or Unscented Kalman Filter (UKF) may be required.
- **Model Accuracy**: Accurate modeling of the state transition and measurement models is crucial for good performance. Incorrect models can lead to poor tracking results.

Dept. of CSE AIML/B.Tech AIML

- **Computational Complexity**: While the Kalman Filter is computationally efficient, real-time implementations in complex scenarios might require optimizations.
- The Kalman Filter is an efficient recursive algorithm used to estimate the state of a linear dynamical system from noisy measurements.
- By predicting the state and updating it with new observations, it provides a robust method for tracking objects, estimating positions, and controlling systems.

Unit-IV

HIGH-LEVEL VISION: Registration, Registering Rigid and Deformable Objects, Smooth Surfaces and Their Outlines - Contour Geometry, Koenderink's Theorem, The Bitangent Ray Manifold, Object Matching using Interpretation Trees and Spin Images, Classification, Error, and Loss.

High Level Vision

- High-level vision in computer vision refers to the interpretation of complex visual information and understanding scenes or objects in a way that is meaningful for tasks like object recognition, scene understanding, and activity recognition.
- It contrasts with low-level vision, which focuses on basic image processing tasks such as edge detection, color analysis, and texture extraction.
- Registering rigid and deformable objects in computer vision involves aligning images or data sets where the objects of interest can either be assumed to be rigid (i.e., having fixed shape and structure) or deformable (i.e., able to change shape).
- The techniques for each type of registration differ due to the nature of the objects and the transformations involved.

Rigid Registration

- Rigid registration deals with aligning objects where the transformations are constrained to rotations and translations.
- This means the object maintains its shape and size, and the only changes between images are due to the viewpoint or position of the object.

Registering Rigid and Deformable Objects

• In computer vision, registering rigid and deformable objects involves aligning multiple views or images of these objects to achieve a common reference frame.

Rigid Object Registration

- Rigid objects do not change their shape or size during registration, only their position and orientation may change. The goal is to align these objects such that their geometric shapes match perfectly.
- 1. Feature-Based Methods:
 - a. Key point Detection: Extract features like corners or edges from the images using algorithms such as SIFT (Scale-Invariant Feature Transform), SURF (Speeded-Up Robust Features), or ORB (Oriented FAST and Rotated BRIEF).
 - b. Feature Matching: Match these features across different images to establish correspondences.
 - c. Transformation Estimation: Compute the rigid transformation (rotation and translation) that aligns the matched features. Methods such as the RANSAC (Random Sample Consensus) algorithm are often used to handle outliers.

2. Direct Methods:

- Iterative Closest Point (ICP): An optimization algorithm used to minimize the distance between points on the surfaces of the rigid objects. It iteratively refines the alignment of two-point clouds.
- Least Squares Optimization: Compute the rigid transformation that minimizes the sum of squared distances between corresponding points or features.

3. Affine Transformation:

- For scenarios where a slight scale change or skew might be acceptable, affine transformations can be used. They involve rotation, translation, scaling (Changes the size of an object by multiplying its original coordinates by a scaling factor. Scaling can be used to expand or compress the dimensions of an object) and shearing (Slants or distorts the shape of an object by adding a multiple of one coordinate to the other).
- Scaling & shearing are linear transformations that can be applied to 2D or 3D point coordinates.

Deformable Object Registration

• **Deformable objects** can change their shape and size. Registering such objects involves more complex methods to account for these variations.

1. Non-Rigid Deformation Models:

- 1. **Thin-Plate Splines (TPS)**: A method used for smooth and flexible deformations. TPS can model complex warps by minimizing bending energy.
- 2. B-splines: These allow for flexible and local adjustments to model deformations smoothly.

2. Variational Methods:

- Active Contour Models (Snakes): These models evolve to fit the contours of the object by minimizing an energy function that incorporates image gradients and smoothness constraints.
- Level Set Methods: These methods evolve a contour or surface to capture the shape of the object while accounting for changes in topology.

3. Machine Learning Approaches:

• **Deep Learning**: Convolutional Neural Networks (CNNs) and other deep learning models can learn features and deformations directly from data. Models such as deformable object detection networks can be trained to handle a wide variety of deformations.

4. Registration Algorithms:

- **Demons Algorithm**: A technique for non-rigid registration that uses a transformation field to warp one image to align with another.
- **Normal Mixture Models**: For statistical modeling of deformations, these models can capture variations across multiple instances of deformations.

5. Elastic Registration:

• **Elastic Deformation Models**: Use spring-like forces to model deformation, applying elasticity principles to align images or shapes.

Smooth Surfaces and Their Outlines

- Smooth surfaces and their outlines are crucial concepts in computer vision, computer graphics, and geometric modeling.
- Understanding and processing these surfaces involves several key aspects:

Smooth Surfaces

• Smooth surfaces are those without abrupt changes in curvature. They are often modeled mathematically and visualized using various techniques.

1. Mathematical Representation:

- Parametric Surfaces: Represented by equations involving parameters. Common examples include Bézier surfaces, B-splines, and NURBS (Non-Uniform Rational B-Splines). These are used to create smooth and flexible shapes.
- Implicit Surfaces: Defined by an equation F(x, y, z) = 0F(x, y, z) = 0F(x, y, z) = 0. Examples include spheres and toruses. Implicit surfaces are particularly useful for modeling complex shapes where explicit parameterization is challenging.
- Subdivision Surfaces: Created by recursively refining a base mesh (e.g., Catmull-Clark subdivision). These surfaces smooth out as the mesh is refined, making them suitable for high-quality modeling.

2. Surface Smoothing:

- Gaussian Smoothing: Applies a Gaussian filter to smooth the surface, often used in image processing to reduce noise.
- Laplacian Smoothing: Adjusts vertices of a mesh based on the average of neighboring vertices, useful for mesh refinement and noise reduction.
- **Bilateral Filtering**: Preserves edges while smoothing the surface, balancing noise reduction and edge preservation.

3. Surface Fitting:

- Least Squares Fitting: Minimizes the difference between the observed data points and the fitted surface.
- **RANSAC (Random Sample Consensus)**: Robustly fits surfaces to noisy data by iteratively selecting subsets of data points.

Outlines of Smooth Surfaces

Outlines or contours of smooth surfaces are the curves or lines that define the boundary or significant features of a surface.

1. Edge Detection:

- 1. Gradient-Based Methods: Techniques like the Sobel or Canny edge detectors find edges by detecting significant changes in intensity or gradient magnitude.
- 2. Zero-Crossing Methods: Detect edges by identifying points where the second derivative of the image intensity function crosses zero.

2. Contour Extraction:

- Level Set Methods: Evolve contours to match the edges of the object by minimizing an energy function that incorporates edge information and smoothness constraints.
- Active Contours (Snakes): Refine contour shapes to fit the object boundaries by minimizing an energy function that balances edge attraction and smoothness.

3. Shape Analysis:

- Fourier Descriptors: Represent contours using frequency components, useful for shape recognition and analysis.
- **Curvature Analysis**: Analyze the curvature of contours to understand shape properties and identify significant features.

4. Boundary Extraction:

- **Region Growing**: Start from seed points and grow the boundary by adding neighboring pixels that satisfy certain criteria (e.g., similarity in intensity or color).
- Watershed Transform: Segment the image by treating it like a topographic surface, where boundaries correspond to watershed lines.

Applications

- 1. Computer Vision:
 - a. **Object Detection and Recognition**: Smooth surfaces and their outlines help in detecting and recognizing objects by identifying their shapes and boundaries.
 - b. **3D Reconstruction**: Smooth surface modeling is essential for reconstructing 3D shapes from multiple 2D images.

2. Computer Graphics:

- a. **Surface Rendering**: Smooth surfaces are used to create realistic visual representations of objects in graphics and animation.
- b. **Mesh Refinement**: Smooth surfaces improve the visual quality of 3D models by refining their mesh structures.

3. Medical Imaging:

a. **Organ Segmentation**: Outlines and smooth surface models help in segmenting and analyzing organs and other structures from medical scans.

4. Robotics:

a. **Path Planning**: Understanding smooth surfaces and their contours aids in planning robot trajectories and avoiding obstacles.

Koenderink's Theorem

- Koenderink's Theorem is a significant result in computer vision and differential geometry, particularly in the study of shape perception and the structure of visual information.
- It deals with the properties of the visual system in relation to the shapes and contours of objects.
- Koenderink's Theorem, formulated by Jan Koenderink, addresses the following key insights:
- 1. **Curvature Properties:** The theorem provides a fundamental relationship between the curvatures of a 3D surface and its projection onto a 2D image. Specifically, it describes how the intrinsic curvature of a surface (the way it bends in space) affects the way that surface appears in a 2D image.
- 2. **Shape from Shading**: One of the practical applications of Koenderink's Theorem is in shape-from-shading problems.
- It helps in understanding how variations in shading can be used to infer surface geometry.
- By analyzing the curvature information encoded in the shading of an image, one can reconstruct the 3D shape of the surface.

3. Local Curvature Estimates: The theorem also addresses how local curvature estimates can be derived from the image's intensity variations. This is important for algorithms that aim to recover the shape of an object from its 2D projection.

Mathematical Insight

- Koenderink's Theorem is rooted in the concept of **differential geometry**, which deals with the properties of surfaces using calculus.
- The theorem provides a link between the **Gaussian curvature** and **mean curvature** of a surface and the **image curvature** observed in a 2D projection.

Key Components

- Gaussian Curvature: Measures how the surface curves in different directions at a point. It's the product of the principal curvatures.
- Mean Curvature: Measures the average rate of surface bending and is the average of the principal curvatures.

- Koenderink's results show that these intrinsic properties of a surface influence how it appears in a 2D image.
- For example, if you have an image of a curved surface, the curvatures can be analyzed to deduce the surface's shape.

Implications in Computer Vision

- 1. **Surface Reconstruction**: By understanding how curvatures affect the appearance in 2D images, computer vision algorithms can reconstruct 3D surfaces from their 2D projections.
- 2. **Shape-from-Shading**: The theorem is particularly relevant in algorithms that use shading information to infer the 3D structure of a surface. It helps in understanding how variations in light intensity correspond to surface geometry.
- 3. **Object Recognition**: Accurate recognition of objects can benefit from understanding the relationship between surface curvature and image appearance, aiding in more robust recognition systems.
- 4. **Visual Perception**: The theorem also provides insights into how the human visual system interprets shapes and surfaces based on curvature information.

Koenderink's Theorem is a powerful tool for linking the geometric properties of surfaces to their visual representation, making it a cornerstone in understanding and developing computer vision techniques for 3D shape reconstruction and recognition.

The Bitangent Ray Manifold

- The concept of the Bitangent Ray Manifold is an advanced topic in computer vision and computer graphics, related to the study of how light interacts with surfaces and how these interactions can be used for various visual tasks.
- It is particularly relevant in the context of shape-from-shading, photometric stereo, and 3D reconstruction.

Understanding the Bitangent Ray Manifold

1. Surface Reflection and Rays

- <u>Surface Normals</u>: In computer vision, the normal vector at a point on a surface is crucial for understanding how light interacts with the surface. The normal vector is perpendicular to the surface.
- <u>Bitangent Rays</u>: These are rays that lie in the plane formed by the surface normal and the view direction. They are useful in understanding how light that is reflected off a surface from different angles contributes to the visual appearance of that surface.

2. Ray Manifolds

- A ray manifold represents the set of all possible rays that could be emitted from or received by a surface. In the context of bitangent rays:
- Bitangent Ray Manifold is the collection of rays that lie in the bitangent plane, which is defined by the bitangent vector and the surface normal.

Applications in Computer Vision

1. Shape-from-Shading

- Reconstructing 3D Shapes: By analyzing how the intensity of light changes across an image, algorithms can infer the surface shape.
- The bitangent ray manifold provides a geometric framework for understanding how different light paths contribute to shading variations, aiding in the reconstruction of the 3D shape from a single image.

2. Photometric Stereo

- Surface Geometry Recovery: Photometric stereo uses multiple images of a surface under different lighting conditions to recover its 3D shape.
- The bitangent ray manifold helps in understanding how different lighting directions interact with the surface and how these interactions are captured in the images.

3. Reflectance Models

• Surface Appearance: By modeling how light reflects off a surface, bitangent rays help in creating accurate reflectance models. These models are used for realistic rendering and for understanding how surfaces appear under different lighting conditions.

Mathematical Perspective

- The bitangent ray manifold can be mathematically represented as follows:
- **Parameterization**: The manifold can be parameterized using the surface normal and bitangent directions. This helps in mapping out the set of possible bitangent rays for a given surface.
- **Ray-Surface Interaction**: Equations describing the interaction between rays in the bitangent manifold and the surface can be used to derive constraints and solve for unknown surface properties.

Practical Considerations

- Algorithm Design: Designing algorithms to exploit the bitangent ray manifold requires a good understanding of both the mathematical properties of ray manifolds and the physical properties of light reflection.
- **Computational Efficiency**: Efficiently computing and utilizing the bitangent ray manifold involves careful consideration of computational resources, especially when dealing with high-resolution images and complex scenes.
- In summary, the bitangent ray manifold is a sophisticated concept in computer vision that provides insights into how surface geometry and light interaction affect image appearance.
- It plays a crucial role in various applications involving shape reconstruction, surface analysis, and realistic rendering.

Object Matching using Interpretation Trees and Spin Images

- Object matching is a fundamental problem in computer vision, where the goal is to identify and match objects between different images or scenes.
- Two techniques that are particularly relevant for this task are Interpretation Trees and Spin Images.
- Each of these approaches addresses the problem from different angles, and they can be used either separately or in combination to improve matching accuracy.

1. Interpretation Trees

• Interpretation Trees are used for object recognition and matching by representing and comparing the spatial relationships between features in a scene. Here's a detailed breakdown of how they work:

Concept

- **Feature Representation**: Interpretation Trees represent objects based on the spatial relationships between features. Features are typically points, edges, or regions extracted from the object in the image.
- **Tree Structure**: Each node in the tree represents a feature, and the branches represent the spatial relationships (e.g., distance, angle) between these features.

Process

- 1. Feature Extraction: Extract key features from the images of the objects to be matched.
- 2. **Tree Construction**: Build an interpretation tree for each object based on the extracted features and their spatial relationships.
- 3. **Matching**: Compare the interpretation trees of the objects. This involves checking the similarity between the trees based on feature correspondences and spatial relationships.

Advantages

- **Robustness to Transformations**: Interpretation Trees are robust to changes in scale, rotation, and partial occlusions because they focus on relative spatial relationships rather than absolute positions.
- **Structured Representation**: They provide a structured way to represent and compare objects, making it easier to handle complex shapes and relationships.

2. Spin Images

• **Spin Images** are a shape descriptor used for 3D object recognition and matching. They capture the shape of an object in a way that is invariant to rotation and translation. Here's how Spin Images work:

Concept

• **3D Shape Representation**: Spin Images represent the local shape around a point on the object's surface by projecting it onto a 2D histogram. This histogram captures the distribution of points around the reference point in terms of distance and angle.

Process

- 1. Feature Point Selection: Select key points on the surface of the 3D object.
- 2. **Spin Image Calculation**: For each feature point, compute the Spin Image by projecting the surface points around the feature point onto a 2D plane. The resulting image is a histogram of point distributions relative to the feature point.
- 3. **Matching**: Compare Spin Images from different objects. This involves computing similarities between the histograms to find correspondences between the objects.

Advantages

- **Invariance to Orientation**: Spin Images are invariant to rotation and translation, making them effective for matching objects in different orientations.
- Local Shape Descriptor: They provide a robust local shape descriptor that captures the geometry around a specific point on the surface.

Combining Both Techniques

- In practice, combining Interpretation Trees and Spin Images can leverage the strengths of both methods:
- 1. Feature Extraction with Spin Images: Use Spin Images to describe local surface patches or features.
- 2. Hierarchical Matching with Interpretation Trees: Build an interpretation tree based on these features

and their spatial relationships. Use this tree to handle global object structure and relationships.

- This combined approach can improve matching accuracy by using detailed local shape information from Spin Images and structured spatial relationships from Interpretation Trees. This is particularly useful in scenarios with complex shapes, varying orientations, and partial occlusions.
- **Interpretation Trees** focus on the hierarchical and spatial relationships between features, providing a robust structure for object recognition and matching.
- **Spin Images** capture local shape information around key points and are invariant to transformations, making them useful for recognizing and matching 3D shapes.
- By integrating both methods, you can achieve more accurate and robust object matching, handling a wider range of variations in object appearance and positioning.

Classification

• In computer vision, classification, error, and loss are fundamental concepts used to evaluate and improve the performance of models, particularly in tasks like image recognition, object detection, and segmentation. Here's an overview of these concepts:

Classification

• Classification in computer vision refers to the task of assigning an input image (or a part of an image) to one of several predefined categories or classes. This is typically achieved using machine learning algorithms, particularly deep learning models such as Convolutional Neural Networks (CNNs).

Key Aspects

- Class Labels: The predefined categories that the model can predict. For example, in an image classification task, class labels might include "cat," "dog," "car," etc.
- Training Data: Images that are labeled with their respective classes, used to train the model.
- Predictions: The output of the model, which assigns a class label (or a probability distribution over class labels) to an input image.

Examples

- Image Classification: Classifying an image as "dog" or "cat."
- Object Detection: Identifying objects within an image and classifying each object (e.g., detecting and classifying multiple objects like "car" and "pedestrian" in a street scene).
- Semantic Segmentation: Classifying each pixel in an image into predefined categories (e.g., "road," "sky," "building").

Error

• Error quantifies the discrepancy between the model's predictions and the actual ground truth labels. It's a measure of how well or poorly the model performs.

Types of Error

• <u>Classification Error</u>: The proportion of incorrect predictions out of the total predictions. For example, if a model incorrectly classifies 3 out of 100 images, the classification error rate is 3%.

• <u>Confusion Matrix</u>: A table that summarizes the performance of a classification algorithm by showing the true positives, false positives, true negatives, and false negatives. It helps in understanding the types of errors made by the model.

Examples

- True Positive (TP): Correctly predicted instances of a class.
- False Positive (FP): Incorrectly predicted instances where the model predicted a class that was not actually present.
- False Negative (FN): Instances where the model failed to identify the class that was actually present.
- True Negative (TN): Correctly predicted instances where the class was not present.

Loss

• **Loss** is a measure of how well the model's predictions match the ground truth, expressed as a numerical value. It quantifies the penalty for errors and guides the optimization process during training.

Types of Loss Functions

• **Cross-Entropy Loss**: Commonly used for classification tasks, it measures the difference between the true class distribution and the predicted probability distribution. For binary classification, it's often referred to as Binary Cross-Entropy, and for multi-class classification, it's referred to as Categorical Cross-Entropy.



Where y_i is the ground truth label and \hat{y}_i is the predicted probability for class i.

Mean Squared Error (MSE): Used in regression tasks and measures the average squared difference between predicted and actual values. While not typical for classification, it's used in tasks involving continuous outputs.

$$L=rac{1}{N}\sum_{i=1}^N(y_i-\hat{y}_i)^2$$

Hinge Loss: Often used for Support Vector Machines (SVMs), it penalizes predictions based on the margin between classes. It's suitable for binary classification tasks.

$$L = \max(0, 1 - y_i \cdot \hat{y}_i)$$

Optimization

- **Gradient Descent**: A common optimization algorithm used to minimize the loss function by iteratively updating the model parameters in the direction that reduces the loss.
- **Backpropagation**: In deep learning, backpropagation computes gradients of the loss function with respect to the model parameters, which are then used to update the weights during training.

Summary

- Classification: Assigning labels to images based on learned features and patterns.
- Error: A measure of the discrepancy between predictions and true labels, often summarized by a confusion matrix or error rate.
- Loss: A numerical measure of how well the model's predictions match the true labels, used to guide the optimization and training process.

Unit-V

OBJECT DETECTION AND RECOGNITION: Detecting Objects in Images - The Sliding Window Method, Face Detection, Detecting Humans, Boundaries and Deformable Objects, Object Recognition – Categorization, Selection, Applications – Tracking People, Activity Recognition.

Detecting Objects in Images

• Detecting objects in images is a fundamental task in computer vision, crucial for various applications such as autonomous driving, surveillance, and robotics.

Key Concepts

- 1. Object Detection vs. Object Classification:
 - **Object Detection**: Identifies and localizes multiple objects within an image, providing bounding boxes and class labels.
 - **Object Classification**: Determines the presence of a single object class in the image without localization.

2. Bounding Boxes:

• Rectangular areas that enclose detected objects, defined by coordinates (x, y) for the top-left corner and width and height.

3. Class Labels:

• Categories assigned to detected objects, such as "car," "person," or "dog".

Techniques and Methods

1. Traditional Approaches:

- Haar Cascades: A machine learning object detection method based on feature extraction and classifiers, often used for face detection.
- HOG (Histogram of Oriented Gradients): Extracts features based on the gradients of image intensity, commonly used with classifiers like SVM.

2. Sliding Window Method:

• A classical method that involves moving a window across the image and classifying each segment. It can be computationally intensive but is foundational for understanding object detection.

3. Region Proposal Networks (RPN):

• Part of the Faster R-CNN architecture, RPNs generate proposals (potential bounding boxes) which are then refined and classified.

4. Deep Learning Approaches:

- Convolutional Neural Networks (CNNs): Used to automatically learn features from images. CNNs are the backbone of many modern detection systems.
- YOLO (You Only Look Once): A real-time object detection system that frames detection as a single regression problem, predicting bounding boxes and class probabilities simultaneously.

- SSD (Single Shot Multi Box Detector): Similar to YOLO but uses multiple feature maps at different scales to detect objects, balancing speed and accuracy.
- 5. Transformers for Object Detection:
- **DETR (DEtection TRansformer)**: A new approach that combines transformers with CNNs for end-toend object detection, achieving high accuracy with fewer post-processing steps.

Evaluation Metrics

- Mean Average Precision (mAP): Measures the precision and recall across different intersection-overunion (IoU) thresholds, commonly used to evaluate the performance of detection algorithms.
- **IoU** (**Intersection over Union**): A metric used to determine the overlap between the predicted bounding box and the ground truth box.

Challenges

- **Occlusion**: Objects may be partially blocked, complicating detection.
- Variability: Changes in lighting, viewpoint, and scale can affect detection accuracy.
- **Class Imbalance**: Some object classes may be underrepresented in training data, leading to biased models.

The Sliding Window Method

The sliding window method is a classical technique for object detection in images. It involves moving a fixedsize window across an image and classifying the content within that window.

Steps in the Sliding Window Method

1. Define the Window Size:

a. Choose the dimensions of the sliding window based on the size of the objects you want to detect.

4. Slide the Window:

a. Move the window across the image in a step-wise fashion, typically both horizontally and vertically. The step size can vary; smaller steps provide more detailed coverage but increase computational cost.

3. Extract Features:

For each position of the window, extract features that can help in classifying the content. Common techniques include:

- Histogram of Oriented Gradients (HOG)
- Color histograms
- Texture descriptors

4. Classify the Window:

- Use a classifier (like SVM, decision trees, or neural networks) to determine if the content within the window corresponds to the object of interest. This can be a binary classification (object present vs. not present).
- 5. Repeat for Different Scales:
- To detect objects of varying sizes, the window can be resized, and the process is repeated for each scale.

6. Non-Maximum Suppression:

• After classifying multiple overlapping windows, apply non-maximum suppression to eliminate duplicate detections and keep the most confident ones.

Advantages

- **Simplicity**: The method is straightforward and easy to implement.
- Flexibility: It can be adapted to different object sizes by using multiple scales.

Disadvantages

- **Computationally Expensive**: Sliding the window over every possible position and scale can be very slow, especially for high-resolution images.
- Redundant Computation: Many windows might contain similar information, leading to inefficiencies.

Improvements

- Using Region Proposal Networks (RPNs): Modern approaches often combine the sliding window method with more advanced techniques like RPNs in convolutional neural networks (CNNs) to improve efficiency and accuracy.
- **Integrating Deep Learning:** Utilizing deep learning models (like YOLO, SSD) can significantly enhance performance by learning features directly from the data.

Applications

- Face Detection: Locating faces in images.
- Vehicle Detection: Identifying vehicles in traffic images.
- **Object Tracking**: Following specific objects across frames in video analysis.
- The sliding window method laid the groundwork for many advances in computer vision and continues to be relevant, especially as part of hybrid approaches in modern object detection systems.

Face Detection

- Face detection, also called facial detection, is an artificial intelligence (AI)-based computer technology used to find and identify human faces in digital images and video.
- Face detection technology is often used for surveillance and tracking of people in real time. It is used in various fields including security, biometrics, law enforcement, entertainment and social media.
- In face analysis, face detection uses facial expressions to identify which parts of an image or video should be focused on to determine age, gender and emotions.
- In a facial recognition system, face detection data is required to generate a faceprint and match it with other stored faceprints.

How face detection works

- Face detection applications use AI algorithms, ML, statistical analysis and image processing to find human faces within larger images and distinguish them from nonface objects such as landscapes, buildings and other human body parts.
- Before face detection, the analyzed media is preprocessed to improve its quality & remove images that might interfere with detection.

- Face detection algorithms typically start by searching for human eyes, one of the easiest features to detect.
- They then try to detect facial landmarks, such as eyebrows, mouth, nose, nostrils and irises.
- Once the algorithm concludes that it has found a facial region, it does additional tests to confirm that it has detected a face.
- To ensure accuracy, the algorithms are trained on large data_sets that incorporate hundreds of thousands of positive and negative images.
- The training improves the algorithms' ability to determine whether there are faces in an image and where they are.



Face detection methods

- Face detection software uses several different methods, each with advantages and disadvantages:
- Viola-Jones algorithm: This method is based on training a model to understand what is and isn't a face. Although the framework is still popular for recognizing faces in real-time applications, it has problems identifying faces that are covered or not properly oriented.
- **Knowledge- or rule-based:** These approaches describe a face based on rules. Establishing well-defined, knowledge-based rules can be a challenge, however.
- **Feature-based or feature-invariant.** These methods use features such as a person's eyes or nose to detect a face. They can be negatively affected by noise and light.
- **Template matching.** This method is based on comparing images with previously stored standard face patterns or features and correlating the two to detect a face. However, this approach struggles to address

variations in pose, scale and shape.

- **Appearance-based.** This method uses statistical analysis and ML to find the relevant characteristics of face images. The appearance-based method can struggle with changes in lighting and orientation.
- **Convolutional neural network-based.** A Convolutional Neural Network (CNN) is a type of deep learning ANN used in image recognition and processing that's designed to process pixel data.
- A region-based CNN, also called an R-CNN, generates proposals on a CNN framework that localizes and classifies objects in images.
- These proposals focus on areas, or regions, in a photo that are similar to other areas, such as the pixelated region of an eye.
- If this region of the eye matches up with other regions of the eye, then the R-CNN knows it has found a match.
- However, CNNs can become so complex that they "overfit," which means they match regions of noise in the training data and not the intended patterns of facial features.

Single shot detector (SSD)

- While region proposal network-based approaches such as R-CNN need two camera shots -- one for generating region proposals and one for detecting the object of each proposal -- SSDs only require one shot to detect multiple objects within the image.
- Therefore, SSDs are faster than R-CNN. However, SSDs have difficulty detecting small faces or faces farther away from the camera.
- Some techniques used in face detection applications include the following:
- **Background removal.** If an image has a plain, mono-color background or a predefined, static one, removing the background can reveal the face boundaries.
- **Skin color.** In color images, skin color can sometimes be used to find faces; however, this might not work with all complexions.
- **Motion.** Using motion to find faces is another option. In real-time video, a face is almost always moving, so users of this method must calculate the moving area. One drawback of this approach is the risk of confusion with other objects moving in the background.

Advantages of Face Detection

- As a key element in facial imaging applications, such as facial recognition and face analysis, face detection creates various advantages for users, including the following:
- **Improved security.** Face detection improves surveillance efforts and helps track down criminals and terrorists. Personal security is enhanced when users use their faces in place of passwords, because there's nothing for hackers to steal or change.
- **Easy to integrate.** Face detection and facial recognition technology is easy to integrate, and most applications are compatible with the majority of cybersecurity software.
- Automated identification. In the past, identification was manually performed by a person; this was

inefficient and frequently inaccurate. Face detection allows the identification process to be automated, saving time and increasing accuracy.

Disadvantages of Face Detection

- Face detection also holds various disadvantages, including the following:
- **Massive data storage burden.** The ML technology used in face detection requires a lot of data storage that might not be available to all users.
- **Inaccuracy.** Face detection provides more accurate results than manual identification processes, but it can also be thrown off by changes in appearance, camera angles, expression, position, orientation, skin color, pixel values, glasses, facial hair, and differences in camera gain, lighting conditions and image resolution.
- A potential breach of privacy. Face detection's ability to help the government track down criminals creates huge benefits. However, the same surveillance can let the government observe private citizens. Strict regulations must be set to ensure the technology is used fairly and in compliance with human privacy rights.
- **Discrimination.** Experts have raised_concerns_about face detection's inaccuracy in recognizing people of color, mostly women, and how that issue could result in falsely connecting people of color with crimes they didn't commit. These worries are part of a broader concern about racial biases in machine learning algorithms.

Detecting Humans

• Human detection is a subfield of object detection. An object detection system that can identify humans in an image is called a human detection system. Let us look at some of the approaches through which we can do human detection.



Early approaches for Human Detection

- The early approaches used for human detection require less computing power than modern techniques and are readily available in computer vision libraries such as OpenCV.
- However, they are not very high in accuracy.

Haar Cascades

- Haar cascades is a feature-based object recognition system.
- OpenCV provides Haar Cascade-based object detection and contains pre-trained models for face detection, full-body detection, upper body detection, and lower body detection.
- It is most commonly used for face detection.

```
Import cv2
human_body_classifier = cv2.HOGDescriptor()
human_body_classifier.setSVMDetector(cv2.HOGDescriptor_getDefaultPeopleDetector())
```

Modern approaches for Human Detection

- Modern approaches for human detection mainly consist of deep convolution neural networks. Modern object detection systems based on CNN are very accurate and can detect objects belonging to multiple classes.
- Deep neural networks tasks can be solved by using TensorFlow. TensorFlow is an open-source API provided by Google. We will talk about three models that can be implemented by using TensorFlow, which can detect humans. Let us look at some of the popular deep neural networks.

SSD MobileNet V1

- MobileNets are light convolutional neural networks that can be implemented on mobile applications using TensorFlow.
- SSD stands for single shot detector. It can learn to predict the bounding boxes and classify them in one go. SSD can be trained end-to-end.



Faster RCNN Inception V2

- ٠ This model is used for object detection and is faster than RCNN and Fast RCNN.
- It is used in many real-life object detection tasks and is even used in self-driving cars. •
- RCNN stands for region-based convolution neural network and works by dividing the image into different features or regions and classifying them.
- Faster RCNN model can provide good accuracy in less time if GPU acceleration is enabled. •



Dept. of CSE AIML/B.Tech AIML

RCNN NASNet

- NASNet stands for Neural Architecture Search Network.
- NASNet is one of the most accurate models available, and it can detect humans easily. Below are some of the results shown by NASNet.



accurate results by NASNet

YOLO

- YOLO (You only look once) is a real-time object detection system. YOLO passes the image through a network only once and detects the object, unlike the RCNN discussed above that performs object detection on the image regionally.
- YOLO is very fast and hence, very popular for object detection. YOLO splits the input image into maximum grids and generates two bounding boxes and class probabilities for each grid.



Figure 2: The Model. Our system models detection as a regression problem. It divides the image into an $S \times S$ grid and for each grid cell predicts B bounding boxes, confidence for those boxes, and C class probabilities. These predictions are encoded as an $S \times S \times (B * 5 + C)$ tensor.
Comparison of the models:

Model name	Accuracy	Speed
Haar Cascades	Low	High
Histograms of oriented gradients (HOGs)	Low	High
SSD MobileNet V1	High	Low
Faster RCNN Inception V2	High	Low
RCNN NasNet	Very High	Very Low
YOLO	High	High

Object recognition

• Object recognition in computer vision involves two main components: categorization and selection. Here's a deeper dive into each aspect:

1. Categorization

• Definition: Categorization refers to the process of identifying and classifying objects in an image based on their visual features. The goal is to assign a label or class to each detected object.

Techniques:

> Convolutional Neural Networks (CNNs):

• CNNs are the cornerstone of modern image classification. They consist of multiple layers that learn hierarchical features from images, enabling them to classify objects accurately.

> Transfer Learning:

• This involves using pre-trained models (like VGG, ResNet, or Inception) on large datasets (like ImageNet) and fine-tuning them on a specific task, which saves time and resources.

> Multi-Class Classification:

• The model predicts one class from multiple categories (e.g., distinguishing between a cat, dog, or car).

> Multi-Label Classification:

- The model can predict multiple labels for a single image (e.g., an image containing both a dog and a cat).
- > Feature Extraction:
- Techniques like SIFT (Scale-Invariant Feature Transform) or HOG (Histogram of Oriented Gradients) can be used to extract important features that assist in categorization.

2. Selection

Definition: Selection refers to identifying the location of objects within an image, typically by drawing bounding boxes or creating masks around them. It's about localizing the objects.

Techniques:

Object Detection:

- **Bounding Box Detection**: This involves identifying the coordinates of a box around each detected object. Techniques include:
- **a. YOLO** (**You Only Look Once**): A real-time object detection system that predicts bounding boxes and class probabilities simultaneously.
- **b.** SSD (Single Shot Multi Box Detector): Similar to YOLO but uses a different approach for detection.

> Region-Based CNNs (R-CNN):

• Combines region proposal with CNNs to detect objects. Variants like Fast R-CNN and Faster R-CNN improve efficiency by streamlining the process of generating region proposals.

> Instance Segmentation:

• This technique goes further than bounding boxes by providing pixel-wise masks for each object instance. **Mask R-CNN** is a popular model that can perform both object detection and instance segmentation.

> Semantic Segmentation:

• Classifies every pixel in an image into categories, treating all instances of a category the same (e.g., all pixels labeled as "car" regardless of the number of cars).

Applications

- Autonomous Vehicles: Detecting and classifying pedestrians, vehicles, and road signs.
- Retail: Analyzing customer behavior and inventory management through object recognition.
- Healthcare: Assisting in medical image analysis (e.g., identifying tumors in radiology images).

Challenges

- Variability in Appearance: Differences in lighting, angles, and occlusions can significantly affect detection accuracy.
- **Computational Complexity**: Balancing the trade-off between accuracy and processing speed, especially in real-time applications.
- **Dataset Bias**: Ensuring the training data is diverse to improve the model's ability to generalize across different environments.

Tools and Frameworks

- TensorFlow and PyTorch: Widely used for building and training deep learning models.
- **OpenCV**: Provides tools for image processing and traditional computer vision tasks.
- **Detectron2**: A framework for object detection and segmentation developed by Facebook AI Research.
- By effectively combining categorization and selection, object recognition systems can robustly identify and localize objects, enabling a wide range of applications across different fields.

Applications

Tracking People

• Tracking people in computer vision is a fascinating area that involves several techniques and algorithms to identify and follow individuals across frames in video sequences.

1. Detection vs. Tracking

- Detection: Identifying objects (like people) in individual frames using methods like Convolutional Neural Networks (CNNs).
- Tracking: Continuously following detected objects across frames, often using information from previous frames.
- •

2. Tracking Algorithms

- Kalman Filter: A probabilistic algorithm used for estimating the state of a moving object. It predicts the object's future position based on its current state and motion model.
- Particle Filter: Similar to Kalman filters but can handle non-linearities and non-Gaussian noise better.
- Optical Flow: Analyzes motion between two image frames based on pixel intensity changes to track objects.
- Mean Shift and CAMshift: Techniques that locate the peak of a probability distribution, commonly used for tracking objects based on color histograms.

3. Deep Learning Approaches

- YOLO (You Only Look Once): A real-time object detection system that can be combined with tracking algorithms.
- **SORT (Simple Online and Realtime Tracking)**: Utilizes object detections from YOLO or other detectors and applies Kalman filtering for tracking.
- **DeepSORT**: An extension of SORT that incorporates appearance features (using deep learning) to improve tracking accuracy.

4. Multi-Object Tracking (MOT)

• This involves tracking multiple individuals simultaneously, requiring robust data association methods to match detections to existing tracks.

5. Challenges

- Occlusions: When people overlap or are temporarily blocked, tracking becomes difficult.
- Appearance Changes: Variations in clothing, pose, or lighting can hinder tracking accuracy.
- **Real-Time Processing**: Achieving high accuracy while maintaining a fast-processing speed is crucial for applications like surveillance.

6. Applications

- Surveillance Systems: Monitoring public spaces for security purposes.
- Human-Computer Interaction: Enabling gesture recognition and other interactive applications.
- Sports Analytics: Tracking players in games for performance analysis.

7. Evaluation Metrics

• Metrics like Multiple Object Tracking Accuracy (MOTA) and Multiple Object Tracking Precision (MOTP) are used to assess the performance of tracking systems.

Tools and Libraries

- **OpenCV**: A popular library with functions for both detection and tracking.
- **TensorFlow and PyTorch**: Frameworks for implementing deep learning-based detection and tracking models.
- By combining these techniques and continuously improving algorithms, the accuracy and robustness of person tracking in computer vision continue to advance.

Activity recognition

- Activity recognition in computer vision involves identifying and classifying specific actions or behaviors performed by individuals in images or video sequences.
- It plays a crucial role in various applications, from surveillance and human-computer interaction to healthcare and sports analytics.

Types of Activity Recognition

- Static Activity Recognition: Identifying actions from still images (e.g., identifying a person reading).
- Dynamic Activity Recognition: Classifying actions in video sequences (e.g., walking, running, jumping).

2. Techniques Used

- **Feature Extraction**: Techniques such as Histogram of Oriented Gradients (HOG), Scale-Invariant Feature Transform (SIFT), and Optical Flow to extract relevant features from images or frames.
- Machine Learning Models:
 - **Traditional Models**: Support Vector Machines (SVM), Random Forests, and Hidden Markov Models (HMM) were commonly used for activity recognition before deep learning became prevalent.
 - **Deep Learning**: Convolutional Neural Networks (CNNs) for spatial feature extraction, combined with Recurrent Neural Networks (RNNs) or Long Short-Term Memory networks (LSTMs) for temporal sequence analysis.

3. Datasets

- Several publicly available datasets are used for training and benchmarking activity recognition systems:
- UCF101: A dataset with 101 action categories captured in various scenarios.
- HMDB51: A dataset with 51 action categories, featuring clips from movies and daily life.
- Kinetics: A large-scale dataset containing diverse human actions collected from YouTube.

4. Challenges

- Variability in Actions: Differences in how people perform the same action (e.g., running vs. jogging).
- **Occlusion**: Objects or people being partially blocked can hinder recognition accuracy.
- **Context Dependence**: The same action might be interpreted differently depending on the context (e.g., a person waving can mean greeting or signaling).

5. Applications

- Surveillance: Monitoring activities in public spaces to detect suspicious or unusual behavior.
- Healthcare: Monitoring patients' movements and activities for rehabilitation or elderly care.
- **Human-Computer Interaction**: Enhancing user experience by recognizing gestures and actions for control and interaction.
- Sports Analytics: Analyzing player movements and actions for performance evaluation and coaching.
- Smart Homes: Detecting activities within a home to assist with elderly care or provide security alerts.

6. Evaluation Metrics

• To assess the performance of activity recognition systems, metrics such as accuracy, precision, recall, and F1 score are commonly used. Additionally, confusion matrices can help visualize misclassifications among different activities.

7. Future Directions

- **Multi-Modal Approaches**: Combining video data with other modalities like audio or sensor data (e.g., wearables) for improved accuracy.
- **Real-Time Processing**: Enhancing algorithms to process video streams in real-time for applications like surveillance and interactive systems.
- **Transfer Learning**: Utilizing pre-trained models on large datasets to improve performance on smaller, domain-specific datasets.