

ANNAMACHARYA INSTITUTE OF TECHNOLOGY & SCIENCES
RAJAMPET
(Autonomous)

Department of ARTIFICIAL INTELLIGENCE AND DATA SCIENCE
Lecture Notes



Name of the Faculty: B.SIREESHA
Class: III B.Tech II Sem.
Branch: AI&DS
Name of the Course: COMPUTER VISION
Subject Code: 23A306FT
Academic Year: 2025-2026

UNIT-I:LINEARFILTERS

IntroductiontoComputerVision

ComputerVisionisthefieldthatenablescomputerstosee,**analyze,andunderstand images and videos** similar to human vision.

Computervisionisaninterdisciplinaryfieldofartificial intelligencethatenablesmachinesto interpret and gain high-level understanding from digital images or videos, essentially mimickinghumansight.

- **CoreConcepts:**Thefieldinvolvestaskssuchasobjectdetection(locatingobjects), objectrecognition(identifyingobjects),imageclassification(categorizingimages), and scene understanding (interpreting relationships between objects).
- **Process:** The typical workflow involves image acquisition, pre-processing (noise reduction,contrastenhancement),featureextraction(identifyingkeyelementslike edges, shapes), pattern recognition (comparing features to known examples), and decision-making.
- **Applications:**Itisusedacrossmanyindustriesincludinghealthcare(medicalimage analysis), autonomous vehicles (navigation, obstacle avoidance), manufacturing (quality control), and security (surveillance, face recognition).

GoalsofComputerVision

- Extractusefulinformationfromimages
- Detectobjects,edges,textures
- Recognizepatternsandshapes
- Understandscenestructure

Applications

- Facerecognition
- Medicalimageanalysis
- Autonomousvehicles
- Surveillancesystems
- Industrialinspection

DigitalImageRepresentation

A digital image is represented as a 2D function:

$$f(x,y)$$

Where:

- x, y → spatial coordinates
- $f(x,y)$ → intensity (gray value)

GrayLevelImage

- Values range from **0 to 255**
- 0 → Black
- 255 → White

ColorImage

- Represented using **RGB**
- Each pixel has 3 values: (R, G, B)

LinearFilters

A linear filter modifies an image using a linear operation.

Linear filters are transformations applied to an image to modify signals, often for noise reduction, sharpening, or blurring.

Properties of Linear Systems

1. Additivity

$$T(f_1 + f_2) = T(f_1) + T(f_2) \quad T(f_1 + f_2) = T(f_1) + T(f_2)$$

2. Homogeneity

$$T(af) = aT(f) \quad T(af) = aT(f)$$

If both properties are satisfied → system is **linear**. **Convolution in**

Computer Vision

Convolution is the **core operation** in image filtering.

- **Convolution:** This is a mathematical operation central to linear filtering. It combines an input image with a smaller matrix of weights, called a **kernel** or mask, to produce

anewimage.Theoutputpixel'svalueisaweightedsumofitsneighboringpixelsin the original image.

- **Properties:** Linear filters are characterized by the properties of superposition and shift-invariance (moving the input shifts the output by the same amount).
- **Types:** Common types include low-pass filters (for smoothing/blurring) and high-pass filters (for sharpening/edge detection).

Mathematical Definition

$$g(x,y)=\sum_m\sum_n f(m,n)\cdot h(x-m,y-n) \quad g(x,y)=\sum_m\sum_n f(m,n)\cdot h(x-m,y-n)$$

Where:

- f → input image
- h → kernel (filter mask)
- g → output image

Kernel

- Small matrix (e.g., 3×3, 5×5)
- Slides over the image
- Determines type of filtering

Example: Averaging Filter

$$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

Purpose → **Smoothing/Noise removal**

Shift Invariant Linear Systems

A system is **shift invariant** if shifting the input image causes the output to shift by the same amount.

A system is shift-invariant if a shift in the input signal causes an identical shift in the output signal, without changing the shape of the output. Linear, shift-invariant (LSI) systems are fundamental in computer vision and can be entirely characterized by their impulse response (the kernel).

$$T[f(x-x_0, y-y_0)] = g(x-x_0, y-y_0) \quad T[f(x-x_0, y-y_0)] = g(x-x_0, y-y_0)$$

Importance

- Convolution assumes shift invariance
- Most image filters satisfy this property

Spatial Frequency

Spatial frequency refers to how rapidly intensity values change in an image.

- **Spatial Frequency:** This describes how quickly pixel values change across an image. Low frequencies represent smooth, gradual changes (like large shapes), while high frequencies capture rapid changes (like edges and fine details).

Frequency Type Meaning

Low frequency Smooth regions

High frequency Edges, noise

Fourier Transform (FT)

Fourier Transform converts an image from **spatial domain** → **frequency domain**.

- **Fourier Transform (FT):** The FT is a crucial tool that decomposes an image from its original spatial domain into the frequency domain, representing it as a sum of sine and cosine waves of different frequencies.
- **Convolution Theorem:** A key property is that convolution in the spatial domain is equivalent to simple multiplication in the frequency domain, which can be computationally efficient, especially for large kernels.

2D Fourier Transform

$$F(u,v) = \sum_x \sum_y f(x,y) e^{-j2\pi(ux+vy)} \quad F(u,v) = \sum_x \sum_y f(x,y) e^{-j2\pi(ux+vy)}$$

Meaning

- DC Component (0,0) → Average brightness
- High frequencies → Edges
- Low frequencies → Smooth areas

SamplingandAliasing

- **SamplingandAliasing:** Sampling is the process of converting a continuous signal (like light entering a camera) into a discrete digital image. Aliasing occurs when a signal is sampled at a rate too low to accurately represent its highest frequency components (below the Nyquist rate). This results in visual distortions like "jaggies" (staircase effect on edges) in images. Anti-aliasing filters (typically low-pass filters) are used before sampling to remove high frequencies that would otherwise cause this distortion.

FiltersasTemplates

Filters act as **templates** to detect specific patterns.

- **Filters as Templates:** Linear filters, specifically kernels, act as templates. By convolving an image with a specific kernel, you can find patterns or features that match that kernel's design, such as using an edge-detector kernel to find edges

Examples

- Edgedetection kernels
- Linedetection kernels
- Cornerdetection kernels

NormalizedCorrelation

Used for **pattern matching**.

- **Normalized Cross-Correlation (NCC):** This technique is widely used for pattern matching and template matching. It measures the similarity between a template and different regions of a larger image.
- **Function:** NCC produces a value between -1 (perfect mismatch) and 1 (perfect match) and is robust to linear changes in brightness and contrast. A peak in the correlation output indicates the location where the pattern is most likely present.

Formula

$$C = \frac{\sum (f - \bar{f})(t - \bar{t})}{\sqrt{\sum (f - \bar{f})^2 \sum (t - \bar{t})^2}} \quad C = \frac{\sum (f - \bar{f})(t - \bar{t})}{\sqrt{\sum (f - \bar{f})^2 \sum (t - \bar{t})^2}}$$

Range

- -1to+1
- +1→Perfectmatch

ScaleandImagePyramids

- **Image Pyramids:** This is a multi-scale representation of an image, created by repeatedly smoothing (usually with a Gaussian filter) and subsampling (downscaling) the image. The original, full-size image is at the bottom, and progressively smaller, blurrier versions stack on top.
- **Purpose:** Pyramids are useful for detecting objects that appear at different sizes or scales within an image, as a fixed-size template can find large objects in lower-resolution layers and small objects in higher-resolution layers. They are also used in image blending and compression.

ImagePyramid

- Multi-resolution representation
- Same image at different scales

Types

1. **Gaussian Pyramid**—smoothing+subsampling
2. **Laplacian Pyramid**—difference between levels

Applications

- Object detection
- Image compression
- Feature extraction

UNIT-II:EDGE DETECTION

Noise in Images

Additive Stationary Gaussian Noise (AWGN) is a fundamental noise model used in image processing and communication systems. It represents random disturbances that affect pixel values.

- **Additive:** The noise is simply added to the original pixel value; it is independent of the image content.
- **Stationary (or White):** The noise has a uniform power spectral density across all frequencies, meaning all frequencies are equally present, analogous to white light.

- Gaussian: The amplitude of the noise follows a normal (bell-curve) probability distribution, with a mean of zero and a specific variance (σ^2 sigma squared

σ^2) that determines the noise strength.

- Source: In practical terms, it often models thermal noise in electronics sensors and circuits used for image acquisition.

Additive Stationary Gaussian Noise

- Noise is **unwanted random variation** in image intensity.
- **Additive noise:**

$$g(x,y) = f(x,y) + n(x,y) \quad g(x,y) = f(x,y) + n(x,y) \quad g(x,y) = f(x,y) + n(x,y)$$

- Gaussian noise follows **normal distribution**.
- Stationary \rightarrow statistical properties do not change over space.

Effect of Noise

- Creates false edges
- Reduces accuracy of edge detection

Why Finite Differences Respond to Noise

Finite difference operators approximate the derivative of an image by calculating the difference between adjacent pixel values.

- Noise Sensitivity: Because noise is often characterized by high-frequency random variations between neighboring pixels, these simple difference operations amplify the noise significantly, creating many false edges.
- Second Derivatives: Second-order derivatives (like the basic Laplacian) are even more sensitive to noise than first-order methods (like the basic Sobel operator)
- Finite differences approximate derivatives:

$$\frac{\partial f}{\partial x} \approx \frac{f(x+1) - f(x)}{\Delta x} \quad \frac{\partial f}{\partial x} \approx \frac{f(x+1) - f(x)}{\Delta x}$$

- Noise has **high-frequency components**
- Derivatives amplify high frequencies \rightarrow **noise increases**

Estimating Derivatives

To combat noise sensitivity, smoothing is combined with differentiation.

- Smoothing: A Gaussian filter is applied first to the image to suppress high-frequency noise components.
- Derivative of Gaussian (DoG): Instead of applying a Gaussian filter and then a derivative operator separately, the operations can be combined due to the associative property of convolution. We can pre-calculate a filter kernel that is the derivative of a Gaussian function and convolve the image with this single kernel. This provides a smoothed derivative estimate, reducing the noise impact while still highlighting intensity changes

First-order Derivative

- Detects **edges**
- Sensitive to noise

Second-order Derivative

- Zero-crossings indicate edges
- More sensitive to noise

Derivative of Gaussian Filters

- Noise Reduction: Smoothing works as a low-pass filter, attenuating the high-frequency components that typically correspond to noise, while preserving important structural information (edges are lower frequency than random noise).
- Gaussian's Optimality: The Gaussian filter is widely used because it provides an optimal trade-off between spatial localization (accuracy of edge position) and frequency localization (effective noise suppression). It minimizes the product of spatial and frequency uncertainty, as described by the uncertainty principle.
- Separability: The 2D Gaussian filter is separable into two 1D filters (horizontal and vertical), which makes the computation much more efficient
- Combines **smoothing+differentiation**
- Reduces noise sensitivity

$$\frac{\partial}{\partial x}(G * f) = (\frac{\partial G}{\partial x}) * f \quad \frac{\partial}{\partial x}(G * f) = (\frac{\partial G}{\partial x}) * f \quad \frac{\partial}{\partial x}(G * f) = (\frac{\partial G}{\partial x}) * f$$

Advantage

- Smooth image first

- Thencomputederivative

WhySmoothingHelps

- Removeshigh-frequencynoise
- Makes edgedetectionreliable

CommonSmoothingFilter

- Gaussianfilter

ChoosingaSmoothingFilter

- Small $\sigma \rightarrow$ detects fine edges but noisy
- Large $\sigma \rightarrow$ smooth but may miss details
- Trade-off between **accuracy and noise suppression**

LaplacianforEdgeDetection

The Laplacian operator is a second-order derivative operator (

$\nabla^2 f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$ that detects edges in all directions.

- Zero Crossings: The Laplacian response is zero in areas of constant intensity and changes sign (crosses zero) exactly at the center of an ideal edge. Edges are detected by finding these "zero crossings."
- Laplacian of Gaussian (LoG): The practical implementation combines Gaussian smoothing with the Laplacian operator (LoG or "Mexican hat" kernel) to make it robust to noise.

Laplacian Operator

$\nabla^2 f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$

Laplacian of Gaussian (LoG)

- Smoothimage first
- ThenapplyLaplacian
- Edgedetectedusingzerocrossings

Gradient-BasedEdgeDetectors Gradient

Magnitude

$$|\nabla f| = \sqrt{f_x^2 + f_y^2}$$

GradientDirection

$$\theta = \tan^{-1}\left(\frac{f_y}{f_x}\right)$$

Examples

- Sobeloperator
- Prewittoperator
- Robertsoperator

OrientationRepresentationsandCorners

Gradient-basedmethodsworkbyfindinglocationswherethefirstderivative (gradient magnitude) of the image intensity is maximum.

- SobelOperator:Acommonexampleuses

3x3cross3

3x3kernelstoapproximatethegradientsinthehorizontal (G_x) and vertical (G_y) directions.

- GradientMagnitude&Orientation:Theedgestrength(magnitude)iscalculatedas

$G = \sqrt{G_x^2 + G_y^2}$

$$G = \sqrt{G_x^2 + G_y^2}$$

,andthedirection(orientation)is

$\theta = \arctan(G_y/G_x)$ theta equals arctangent of G_y/G_x

$\theta = \arctan(G_y/G_x)$

- Canny Edge Detector: This robust, multi-stage algorithm uses gradient information, non-maximum suppression (to thin edges), and hysteresis thresholding (to link strong edges to nearby weak edges).
- Corners: Corners are points where the image gradient changes significantly in two or more directions simultaneously. Detectors like the Harris corner detector analyze the gradient information to identify these key feature points
- Orientation = direction of edge
- Corners occur when **intensity changes in multiple directions**
- Important for feature detection

UNIT-III: TEXTURE

What is Texture?

Texture is an important visual cue that describes the local characteristics of surfaces in an image.

Filter banks are a primary tool for texture analysis.

- Filter Bank: This is a collection of filters designed to analyze an image at multiple scales and orientations. Common filter banks use Gabor filters, Laplacian of Gaussian, or wavelet transforms to capture different frequency and orientation information.
- Process: Convolution of an image with these filters helps reveal local structures because a strong filter response occurs where the image pattern matches the filter kernel's design. The collection of responses across the entire bank provides a rich description of the local texture

Texture describes **surface properties** like:

- Smoothness
- Roughness
- Regularity

Representing Texture Using Filter Banks

The raw output from a large filter bank can be redundant, so statistical features are extracted to create a compact and meaningful texture representation.

- **Statistical Features:** For each filtered image (channel), statistical measures are calculated, such as the mean, variance, skewness, and energy.
- **Feature Vectors and Textons:** These statistics are concatenated into a single feature vector that represents the texture of an area. In some approaches, filter responses are clustered into prototypical response vectors called "textons," which act as a dictionary for characterizing different texture classes
- Apply multiple filters at different orientations & scales
- Each filter captures specific patterns

Example

- Gabor filters
- Oriented Gaussian filters

Statistics of Filter Outputs

- Mean
- Variance
- Energy
- Correlation

Used as texture features.

Analysis and Synthesis Using Oriented Pyramids

- **Image Pyramids:** These multi-scale representations (Gaussian and Laplacian) are key to texture analysis and synthesis.
- **Laplacian Pyramid (LP):** This structure captures the "detail images" or high-frequency components at various resolutions by subtracting blurred (low-pass) versions from the original images. The LP is crucial for applications like image compression and blending because it separates details at different scales.
- **Oriented Pyramids:** These are similar to Laplacian pyramids but incorporate orientation-specific filters (like Gabor filters) to analyze directional structures in the

texture. They are highly effective for capturing complex patterns and directional information.

- **Frequency Domain:** Filters can be designed and analyzed in the spatial frequency domain using the Fourier Transform, which helps in understanding their effect on different texture frequencies

Analysis

- Decompose image into orientations and scales

Synthesis

- **Texture Synthesis:** This is the process of generating a large digital image from a small sample image such that the new image has the same visual texture appearance. It is widely used in computer graphics for tasks like texture mapping in 3D rendering to avoid visible tiling artifacts and seams.
- **Homogeneity:** In image processing, homogeneity is a measure of how similar the elements (pixels or texels) of an image region are. A homogeneous region might consist of the same grey level or a single, uniform texture pattern.
- **Synthesis by Sampling Local Models:** Many modern texture synthesis algorithms are non-parametric and rely on sampling existing neighborhoods from a sample image to create new pixels or patches, often using Markov Random Field models
- Reconstruct texture from filter responses

Laplacian Pyramid

- Multi-scale image representation
- Capture texture details at different resolutions

Spatial Frequency Domain Filters

- Low-pass → smooth textures
- High-pass → rough textures
- Band-pass → periodic textures

OrientedPyramids

- Capturedirectionalinformation
- Usefulfortextureclassification

TextureSynthesis

Goal

- Generatenewtexturesimilartooriginal

Methods

- Samplinglocalmodels
- Statisticalmatching

Homogeneity

- Texture consistencyacrossregions
- Usedinsegmentation

ShapefromTexture

ShapefromTexture(SfT)isacomputervisiontechniquewherea3Dshapeorsurface orientation is recovered from a single 2D image using texture as a primary cue.

- Mechanism:Humanscaneasilyperceivedepthandshapefromtexturechanges(e.g., repeatingpatternsappearingsmalleranddenserfurtheraway).SfTalgorithmsmimic this by analyzing the distortion or perspective changes of texture elements (texels) acrossasurface.
- For Planes: The analysis is often simplified by assuming the surface can be approximatedaslocallyplanar,whichallowsforthederivationofrelationships between texture distortions and surface orientation parameters.
- Texturedistortion givesdepth information
- Usedfor3Dsurfaceestimation

UNIT-IV:SEGMENTATIONBYCLUSTERING

What is Image Segmentation?

Image segmentation is the process of partitioning a digital image into multiple segments (sets of pixels, also known as superpixels). The goal is to simplify and/or change the representation of an image into something that is more meaningful and easier to analyze.

- **Purpose:** The result of image segmentation is a set of segments that collectively cover the entire image, or a set of contours extracted from the image. Each pixel in a region is similar with respect to some characteristic or computed property, such as color, intensity, or texture.
- **Applications:** Segmentation is a foundational step in most computer vision applications, including medical imaging (tumor boundary detection), autonomous driving (identifying road from pavement), and object recognition.

Segmentation divides an image into **meaningful regions**.

Human Vision: Grouping and Gestalt

The human visual system automatically groups visual elements into meaningful wholes. Computer vision algorithms often draw inspiration from the Gestalt principles of perception.

- **Gestalt Principles:** These include proximity (elements close to each other are grouped), similarity (similar elements are grouped), continuity (lines/shapes are seen as continuous), closure (incomplete shapes are perceived as complete), and symmetry.
- **Relevance:** These principles inform how algorithms define "meaningful" regions during segmentation.

Gestalt principles:

- Proximity
- Similarity
- Continuity
- Closure

Applications

- Shotboundarydetection
- Backgroundsubtraction
- Objectdetection

ImageSegmentationbyClusteringPixels

Clusteringalgorithmsgrouppixelsbasedontheirproperties(e.g.,intensity,colorvalues,or texture features) in feature space.

- K-MeansClustering: Apopulariterativealgorithmthatpartitionsdatapoints(pixels) into a predefined number (K) of clusters. It is simple and fast but requires knowing the number of clusters beforehand.

K-MeansClustering

- Grouppixelsbasedonintensity/color
- Iterativeprocess

SegmentationbyGraph-TheoreticClustering

Graph-basedmethodsmodeltheimageasagraph,wherepixelsarenodesandedgeweights represent the similarity between adjacent pixels. Segmentation then becomes a problem of cutting the graph into meaningful partitions.

- Graph Cuts: Algorithms like normalized cuts find a set of edges to remove (cut) that minimizeacostfunctionrelatedtothedissimilaritybetweensegmentsandmaximize the similarity within segments. These methods often produce high-quality, globally optimal segmentations.
- Pixels=nodes
- Similarity=edgeweights
- Partitiongraphtoformsegments

HoughTransform

The Hough Transform is a feature extraction technique used in image analysis, computer vision, and digital image processing. Its primary purpose is to find imperfect instances of objects within a certain class of shapes (e.g., lines, circles, ellipses) by a voting procedure.

- **Fitting Lines:** It transforms image points (x, y) into parameter space (e.g., slope m and intercept b , or polar coordinates ρ and θ). Peaks in the accumulator (parameter space) indicate the presence and parameters of lines in the original image.
- **Fitting Curves:** The concept extends to more complex shapes like circles or ellipses, though the parameter space becomes larger and computationally more expensive. It is a robust method because it can handle gaps in the curves or lines and is relatively noise-resistant.

Used to detect lines and curves.

Line Equation

$$\rho = x \cos \theta + y \sin \theta$$

Fitting Lines and Curves

- Robust to noise
- Used in lane detection, shape detection

UNIT-V: RECOGNITION BY RELATIONS BETWEEN TEMPLATES

Object Recognition

This topic addresses how to move beyond simple template matching (like normalized correlation in Unit I) to more complex object recognition by considering how parts of an object relate to each other.

- **Part-Based Models:** Objects are often represented as a collection of features or "templates" (e.g., a wheel, a window, a door) rather than a single holistic template.
- **Voting Schemes:** Methods like the Generalized Hough Transform use local features to vote for the probable position and orientation of the entire object's center. This approach is robust to occlusion (when parts of the object are hidden) because even a

few visible parts can still cast votes for the correct object location. The method aggregates evidence from different parts to find the best match.

Goal: Identify objects in images.

Voting on Relations Between Templates

- Match parts instead of whole object
- Voting accumulates evidence

Relational Reasoning Using Probabilistic Models

To improve the accuracy and robustness of recognition, especially with variations in viewpoint and articulation, we incorporate probabilistic models and search strategies.

- Probabilistic Models: Models like Bayesian networks or Markov Random Fields define the probability of finding a part at a specific location given the location of other parts. This captures the spatial relationships (e.g., a car's left headlight is typically to the left and slightly above the bumper).
- Search: Recognition often involves searching through a vast parameter space (possible locations, scales, orientations of all parts). Search strategies, such as branch and bound or greedy algorithms, are used to efficiently find the configuration of parts that best match the probabilistic model of the object.
- Handle uncertainty
- Use probabilities to model relationships

Search and Classifier Pruning

This search space can be extremely large. Classifiers help reduce this search effort by quickly discarding unlikely hypotheses.

- Hypothesis Generation & Verification: Initial steps might generate many potential object locations (hypotheses).
- Pruning: Classifiers (e.g., support vector machines, simple neural networks) can be used to rapidly assess the quality of a hypothesis based on local visual features. Hypotheses deemed highly unlikely by the classifier are pruned (removed) from the search space, allowing the more complex, computationally expensive relational reasoning steps to focus only on promising candidates.

- Reducesearchspace
- Useclassifierstoeliminateunlikelymatches

HiddenMarkovModels(HMM)

HiddenMarkovModelsarepowerfulstatisticalmodelsusedtomodelsystemsthatcanbein certain "hidden" states (whicharenotdirectly observable) and produce observableoutputs based on those states.

- Sequential Data: HMMs are particularly effective for sequential or temporal data, suchasspeech,handwriting,orgesturesinsignlanguage.Thesequenceofobserved images(features)correspondstoasequenceof underlyinghiddenstates(theaction being performed).
- Training & Recognition: During training, the model learns the probabilities of transitioning between states and the probabilities of observing specific features in eachstate.Duringrecognition,itusesalgorithmsliketheViterbialgorithmtofindthe most likely sequence of hidden states (e.g., the specific sign being performed) given an input sequence of images.
- Statisticalmodelforsequentialdata

Components

- States
- Observations
- Transitionprobabilities
- Emissionprobabilities

ApplicationsofHMM

- **SignLanguageUnderstanding:**HMMsareaclassicapproachforrecognizingdynamic gestures and signs. The observed visual features (hand position, joint angles, movementtrajectoryovertime)aremodeledasoutputsfromhiddenstates representing different phonemes or signs.
- **Finding People with HMM:** HMMs can be adapted for spatial tasks too, such as modeling the typical visual sequence found when scanning vertically for a person (e.g.,headfeatures->torsofeatures->legfeatures).Themodelcanslideawindow across an image and use the HMM to evaluate if the vertical sequence of features matches a typical human profile structure.
- Speechrecognition

- Gesturerecognition
- Signlanguageunderstanding
- Findingpeopleinvideos