

21. A voxel is a 3D unit in medical imaging similar to a pixel.
 22. In supervised learning, the model is trained on labeled input-output pairs.
 23. In data mining, patterns should be valid, novel, useful, and understandable
 24. Deep learning models often consist of many layers and can learn features automatically
25. Explain the difference between data mining, machine learning, and data analytics.
Give one example for each.

Data Mining is the process of discovering patterns and knowledge from large datasets, often using statistical and machine learning tools.

👉 *Example:* Finding customer purchase patterns in a supermarket database.

Machine Learning (ML) is a subset of AI that enables systems to learn from data and make predictions or decisions without being explicitly programmed.

👉 *Example:* A spam filter that learns to classify emails based on labeled data.

Data Analytics involves interpreting, visualizing, and drawing insights from data to support decision-making.

👉 *Example:* Analyzing website traffic to understand user behavior and improve marketing.

26. Describe the steps of building a machine learning model from raw data to testing.

- **Collect data**
- **Clean and preprocess** (handle missing values, normalize, encode)
- **Feature selection/engineering**
- **Split data** into training/test (and optionally validation)
- **Choose model** (e.g., Decision Tree, SVM)
- **Train model** on training data
- **Evaluate** on test data using metrics (accuracy, F1, etc.)
- **Tune hyperparameters** (e.g., via cross-validation)
- **Deploy and monitor**

27. Compare discriminative and generative models. How do they differ in what they learn?

- **Discriminative models** learn the **conditional probability** $P(y|x)$, i.e., how to **directly predict the label** given the input.

👉 Example: Logistic Regression, Support Vector Machines

🧠 Focus: learning the **decision boundary** between classes.

- **Generative models** learn the **joint probability** $P(x,y)$, i.e., how the **data and labels are generated together**. From that, they can compute $P(y|x)P(y|x)P(y|x)$ using Bayes' rule.
👉 Example: Naive Bayes, Gaussian Mixture Models
🧠 Focus: modeling the **full data distribution**, including how inputs are generated.

Key difference:

- Discriminative: Focus on **classification accuracy**.
- Generative: Can also **generate new samples**, useful for **semi-supervised learning, missing data, or simulation**

28. Why is data representation important in ML? Provide an example of how bad representation can affect performance.

Data representation is crucial in machine learning because models learn patterns based on how data is structured. If features don't capture the true relationships in the data, the model will struggle to learn or generalize.

Example:

Suppose we use raw pixel values of images to classify handwritten digits without resizing or centering them.

If the same digit is shifted or scaled differently in various images, the model may fail to recognize it — leading to **poor accuracy**.

✓ Using techniques like normalization, centering, or converting to edge-based features improves performance.

Conclusion:

Good representation = easier learning and better results.

Bad representation = model confusion, overfitting, or underperformance.

29. Define supervised, unsupervised, semi-supervised, and reinforcement learning with one example each.

?

Supervised Learning

The model is trained on **labeled data**, where each input has a known output. The goal is to learn a mapping from inputs to outputs.

👉 *Example:* Predicting house prices using features like size, location, and number of rooms — where each training example includes the actual price.

?

Unsupervised Learning

The model is trained on **unlabeled data**, trying to discover hidden patterns, structure, or groupings in the data.

👉 *Example:* Customer segmentation using clustering (e.g., K-means) to identify distinct user groups without knowing the true categories.

?

Semi-supervised Learning

Combines a **small amount of labeled data** with a **large amount of unlabeled data**.

Useful when labeling is expensive or time-consuming.

👉 *Example:* Classifying medical images when only a few scans are labeled by experts, and the rest are used to guide learning.

?

Reinforcement Learning

An **agent interacts with an environment**, learns from **feedback signals (rewards or penalties)**, and improves its strategy over time.

👉 *Example:* A robot learning to navigate a maze by trial and error, receiving a reward for reaching the goal and penalties for hitting walls.

30. What are some of the main challenges of applying AI in healthcare, and how are they addressed?

?

Challenge: Data Privacy and Security

Medical data is highly sensitive and protected by regulations (e.g., GDPR, HIPAA).

Sharing raw data across institutions is often not allowed.

⌚ *Solution:* **Federated Learning** allows training models across multiple hospitals without sharing patient data. **Encryption techniques** and **differential privacy** also help secure the learning process.

?

Challenge: Lack of Labeled Data

Annotating medical data (e.g., MRI scans, pathology slides) requires expert knowledge and is time-consuming and costly.

✓ *Solution:* **Semi-supervised learning** can leverage large amounts of unlabeled data with a small labeled subset. **Transfer learning** reuses knowledge from models trained on similar tasks or domains to reduce labeling needs.

💡 Challenge: Model Interpretability

In high-stakes decisions (e.g., diagnosis, treatment recommendation), black-box predictions are not acceptable to clinicians.

🔍 **Solution:** Explainable AI (XAI) tools like **SHAP** and **LIME** provide insight into which features influenced a model's prediction, increasing trust and transparency.

💡 Challenge: Clinical Workflow Integration

AI tools must be compatible with real-world hospital systems and support, not replace, medical professionals.

👩‍💻 **Solution:** Human-in-the-loop systems involve doctors in reviewing or validating AI outputs. User-centered design ensures that models fit into clinical routines without disrupting care.

19. The inter-quartile range (IQR) is calculated as Q3 minus Q1
20. Label encoding is best suited for ordinal variables.
21. In DBSCAN, points in low-density areas are labeled as noises or outliers.
22. A value is considered an outlier if it lies beyond $Q3 + 1.5 \times IQR$ or below $Q1 - 1.5 \times IQR$.
23. Describe three main steps in the data cleaning process and give one example for each.

Data Cleaning: A 3-Step Process

- Step 1: Find the Dirt

- Identify issues in your data before starting the cleaning process

Finding outliers in boxplot

- Step 2: Scrub the Dirt

- Use appropriate cleaning techniques based on the type of errors found.

This is the most intensive step.

Using iqr or z-score methods to remove outliers

- Step 3: Rinse and Repeat

- Once cleaned, repeat steps 1 and 2 to ensure data quality

24. Explain two different methods for detecting and handling outliers, and specify when each method is appropriate.

Z-Score method calculates how many standard deviations a data point is from the mean: $Z = (x - \mu) / \sigma$. A threshold is set for Z; if a value exceeds this threshold, it is removed or replaced. This method works well when data is normally distributed and not skewed.

IQR method calculates the interquartile range: $IQR = Q3 - Q1$. It then checks whether a value is below $Q1 - 1.5 \times IQR$ or above $Q3 + 1.5 \times IQR$. Such values are considered outliers and are removed or imputed. This method works well with skewed data too.

25. Compare and contrast One-Hot Encoding and Target Encoding. When should each be used, and what are the risks?

One-hot encoding transforms categorical feature to numeric binary, where 1 is presence, 0 is absence. Works well for nominal data with low cardinality. Drawback is high memory usage.

Target Encoding replaces categorical values with statistical measures derived from target variable. Useful for high cardinality data but prone to overfitting and data leakage if not handled properly

26. What are the differences between IQR and Z-score methods for outlier detection? Which one is more robust to skewed data?

Z-Score method calculates how many standard deviations a data point is from the mean: $Z = (x - \mu) / \sigma$. A threshold is set for Z; if a value exceeds this threshold, it is removed or replaced. This method works well when data is normally distributed and not skewed.

IQR method calculates the interquartile range: $IQR = Q3 - Q1$. It then checks whether a value is below $Q1 - 1.5 \times IQR$ or above $Q3 + 1.5 \times IQR$. Such values are considered outliers and are removed or imputed. This method works well with skewed data too.

27. Describe the impact of poor data quality on real-world ML applications, providing at least two concrete risks or failures.

Poor data quality negatively impacts many data processing efforts. Companies incur in a cost of at least 10% of its revenue due to poor data quality

For example, building a people loan risk classification model using poor data can

- Wrongfully deny loans for some credit-worthy candidates
- Issue loans to individuals who ultimately default

28. Explain why categorical encoding is necessary and give an example where improper encoding leads to model error.

Categorical encoding is necessary because most machine learning algorithms can only work with numerical input. Categorical variables — such as "color", "brand", or "city" — must be converted into numbers in a meaningful way to allow the model to learn patterns from them.

Without proper encoding, the model might misinterpret categorical data as **ordinal** or **continuous**, leading to incorrect assumptions.

Nominal data can be converted by One-hot, ordinal by Label, high dimensional by Target or Embedding encoding

21. The Manhattan distance is also called the 1 norm.

22. The cosine similarity of two orthogonal vectors is 0.

23. The edit distance between “Peter” and “Piotr” using unit cost is 3.

24. In LOOCV, each data point is used once as the test set.

25. Describe how k-NN performs classification. What are the key choices a practitioner must make when applying k-NN?

k-NN assigns labels based on the majority class among the k nearest training points.

It does not require training — it's a lazy learner — and its accuracy can improve with more data.

Key components include the choice of **k**, **distance metric**, and **neighbor weighting**.

It can also be used for regression by averaging the values of neighboring points. However, it is computationally expensive ($O(n)$) because prediction happens at query time, and it is sensitive to the choice of metric and feature scaling.

26. Compare Euclidean, Manhattan, and Minkowski distances. In what scenarios would each be preferred?

Euclidean distance:

Straight-line (L2) distance.

- Preferred when data is continuous and **differences in all dimensions matter equally**.
- Sensitive to outliers.

Manhattan distance:

Grid-based (L1) distance.

- Better for **high-dimensional, sparse** data (e.g., text vectors).
- Less sensitive to outliers.

Minkowski distance:

Generalized distance:

$$d = (\sum |x_i - y_i|^p)^{1/p}$$

- $p=1 \rightarrow$ Manhattan
- $p=2 \rightarrow$ Euclidean
- Flexible — can **tune p** depending on problem.
Used when you want a compromise between L1 and L2.

27. Explain the components and interpretation of the Structural Similarity Index (SSIM). How does it differ from MSE?

IN CV it's essential to measure the similarity between images.

1. Mean Squared Error (MSE)

MSE measures the average squared differences between pixel values in two images. Useful for loss functions in image processing and compression.

2. Structural Similarity Index Measure (SSIM)

SSIM evaluates structural information, luminance, and contrast similarities

between two images. Effective in tasks involving human perception, like image quality assessment.

- MSE: High values indicate significant differences, but it doesn't consider human perception.
- SSIM: Incorporates perceptual aspects, offering a more human-centered measure.

In SSIM there are three components - luminance (brightness), contrast, and structure.

28. What are the advantages and limitations of Leave-One-Out Cross Validation compared to k-fold Cross Validation?

Leave-One-Out Cross Validation (LOOCV):

Each sample is used as a **single test point**, while the rest form the training set.

- For n data points, LOOCV runs **n iterations**.
 **Advantage:** Uses almost all data for training → **low bias**
 **Limitation:** **Very slow** (high computational cost), and **high variance** in estimates.

k-Fold Cross Validation:

The dataset is split into **k equal-sized folds**. Each fold is used once as test, $k-1$ folds as training.

-  **Advantage: Efficient** (only k iterations), good balance between bias and variance
-  **Limitation:** Slightly higher bias than LOOCV (less training data per fold)

When to use:

- Use **LOOCV** when data is **very small** and accuracy matters more than speed
- Use **k-Fold (e.g., $k=5$ or 10)** for **larger datasets** to save time and get stable results

29. Define edit distance. What operations are allowed, and how is it calculated? Provide an example.

Edit distance measures similarity by transforming one object to another. It uses Substitution, Insertion, Deletion. It shows how many transformations needed between objects. Example: Peter -> Piotr. I for e Piter, insert o pioter, deletes e Piotr

30. What is triplet loss? Describe its purpose, how it works, and where it is commonly used.

Triplet loss is a loss function used to train models to learn similarity between data points. It compares three samples at once:

- an **anchor**,
- a **positive** (same class as anchor),
- and a **negative** (different class).

The idea is to make the model learn an embedding space where the anchor is **closer to the positive** than to the negative by some margin.

It works by minimizing the distance between anchor and positive, and maximizing the distance between anchor and negative.

It's commonly used in **face recognition** (like FaceNet), **image retrieval**, and **verification systems**.

21. The cost function used in linear regression is often Mean Square error.

22. Gradient descent updates weights by moving in the direction of the negative gradient.

23. Polynomial regression allows fitting multidimensional relationships.

24. In SGD, parameter updates are performed on each individual sample or mini-batch.

25. Describe the difference between gradient descent and stochastic gradient descent. What are their advantages and disadvantages?

Gradient descent is a powerful optimization algorithm used in machine learning to minimize a function, often a cost or loss function, by iteratively adjusting parameters(bias and slope) in the direction of negative gradient. GD uses entire training set to compute gradients which can be computationally expensive. However it provides stable and accurate direction for minimizing error

While sgd uses single sample or mini-batch of training set. This makes it **faster and more scalable**, especially for large datasets. However, it introduces **noise** in the updates, which can make the path to the minimum **less stable** or "zigzagged"

27. Why does convexity matter in optimization problems? What issues can arise with non-convex surfaces?

Some models have a single minimum on their error surface, which makes the function **convex**. Convexity ensures the presence of a **unique global minimum**, so gradient-based methods can reliably reach it.

In contrast, some problems have **multiple local minima**, which makes the surface **non-convex**. Because of this, the optimizer may get stuck in a **local minimum** and **fail to reach the global optimum**.

Mathematically, in a convex set, the line between any two points is also entirely contained in the set — this helps ensure smooth optimization paths.

28. Explain how polynomial regression extends linear regression and give an example where it might be useful.

Polynomial regression extends linear regression by allowing the model to fit **non-linear relationships** between input variables and the target. It does this by adding **polynomial terms** (e.g., x^2, x^3, x^2, x^3 , etc.) as features.

Example:

Predicting **housing prices** based on size — where very small or very large houses deviate from the linear price trend. A polynomial curve can better fit such data than a straight line.

29. What are the roles of learning rate and momentum in gradient descent? How do they affect convergence?

Learning rate, a critical hyperparameter in gradient descent is the learning rate, which determines the size of each step during the parameter updates.

- too small value may result in slow convergence
- too large value can lead to overshooting the minimum

Momentum is an additional technique that uses information from **previous gradients** to smooth the updates.

- It helps **reduce oscillations**, especially in **stochastic gradient descent**, and can **accelerate convergence** in directions with consistent gradients.
- It works like a moving average — pushing the optimizer forward and dampening zig-zag behavior.

30. Outline the gradient descent algorithm step-by-step, from initialization to convergence, using SSE as a cost function.

Initialization of parameters usually randomly,

defining cost function (SSE). Use **SSE** to measure the error between true and predicted values

Compute Predictions

For all training examples, compute the predicted output using the current parameters.

compute gradients and update rule where new parameter=old-learning rate*derivative of cost function. Repeat it all.

21. In PCA, the transformation matrix consists of top eigenvectors of the covariance matrix.

22. The Discrete Fourier Transform converts data into the frequency domain.

23. A face can be approximated by a weighted combination of eigenfaces.

24. The Hough transform maps image points to the parameter space.
25. Describe the PCA algorithm step by step. Why do we center the data?

Principal Component Analysis (PCA) is a dimensionality reduction technique that transforms data into a new coordinate system, where the axes (principal components) represent directions of maximum variance.

Step-by-step PCA algorithm:

1. Center the data:

Subtract the mean of each feature so that the dataset has zero mean.

This step is important because PCA assumes that the data is centered around the origin. Without centering, the principal components may be biased toward the mean.

2. Compute the covariance matrix:

This matrix captures the variance within features and the correlation between them.

3. Calculate eigenvalues and eigenvectors of the covariance matrix:

- Eigenvectors define the directions (principal components).
- Eigenvalues indicate the amount of variance explained by each component.

4. Sort eigenvectors by descending eigenvalues:

This ensures that the most important components (those that explain the most variance) are selected first.

5. Select the top k components:

Choose the top k eigenvectors to reduce dimensionality while retaining most of the variance.

6. Project the data onto the new feature space:

Multiply the original centered data by the selected eigenvectors to obtain the transformed dataset.

Why do we center the data?

Centering ensures that PCA captures directions of variance rather than absolute position. Without centering, the covariance matrix may not correctly reflect the true relationships in the data, leading to incorrect principal components.

26. What are the main differences between spatial and transform domain representations? Provide one application for each.

The **spatial domain** refers to the representation of an image as a 2D grid of pixel intensities indexed by coordinates (x,y) .

Image processing in this domain involves direct operations on pixel values, such as convolution, filtering, and geometric transformations.

It is intuitive and commonly used for basic image enhancement, denoising, and edge detection.

The **transform domain** is an alternative representation of an image in which it is expressed as a sum of **basis functions**, such as sine and cosine waves in the Discrete Fourier Transform (DFT).

This domain enables efficient image processing operations like **filtering, compression, and denoising**, because image content is organized according to **frequency components**, allowing targeted manipulation.

Crucially, every transform must have an **inverse**, so the modified image can be reconstructed back to the spatial domain without information loss.

For many tasks — especially when dealing with **periodic structures, noise removal, or compression** — the transform domain offers **greater flexibility and interpretability** than direct spatial processing.

Example applications:

- **Spatial domain:** edge detection using Sobel or Laplacian filters
- **Transform domain:** JPEG image compression using Discrete Cosine Transform (DCT)

27. What are eigenfaces, and how are they used in facial recognition?

Eigenfaces is a technique based on **Principal Component Analysis (PCA)** applied to a set of face images. It extracts **eigenvectors** (called **eigenfaces**) from normalized and centered face images.

Each eigenface represents a direction of variation in face appearance. Any face image can be approximately reconstructed as a **weighted sum** of these eigenfaces.

Face recognition is performed by **projecting new faces into the eigenface space** and comparing their weights with those of known faces.

28. Explain the principle of orthogonality in transforms and its benefits.

Orthogonality in transforms (e.g., DFT, DCT, PCA) means that the basis functions (or vectors) used to represent data are **mutually perpendicular** in the mathematical sense — their inner product is zero.

This property ensures that:

- Each transformed coefficient captures **independent information**, with **no redundancy**
- Energy (or variance) is **preserved and separable** across components
- Reconstruction is **exact and lossless** if all components are used

Benefits:

- Simplifies mathematical operations (e.g., projections, inverses)
- Makes compression and noise filtering more efficient — we can discard low-magnitude components with minimal impact
- Enables dimensionality reduction (as in PCA) while retaining key information

29. What are the steps for applying 2D-DFT and inverse DFT to an image, and how is fftshift used in visualization?

Steps to apply 2D Discrete Fourier Transform (DFT):

1. **Start with an image** represented as a 2D array of pixel intensities
2. **Apply 2D-DFT** using `fft2()` (e.g., from NumPy or MATLAB)
→ Converts spatial domain into **frequency domain**
3. **Optionally apply `fftshift()` to center the zero-frequency component**
→ Without this, low frequencies are at the corners, which is hard to interpret
4. **Visualize the magnitude spectrum** using $\log(1 + |F|)$ to enhance contrast
5. (Optional) **Apply filtering** in frequency domain (e.g., high-pass, low-pass)

Inverse DFT:

6. After processing, apply `ifftshift()` (if needed), then `ifft2()` to reconstruct the image
→ This brings the data **back to spatial domain**

30. Describe the voting mechanism in the Hough Transform and how it handles noise and partial shape occlusion.

Hough Transform is a method used to detect **simple geometric shapes** (such as **lines** and **circles**) in images using their **mathematical equations**.

For lines, the equation $y=mx+b$ is not ideal because it is **unbounded** when the line is vertical. Therefore, the **polar form** $\rho = x\cos\theta + y\sin\theta$ is used instead, where ρ is the distance to the origin and $\theta \in [0^\circ, 180^\circ]$ is the angle of the normal vector.

For circles, the equation $(x-a)^2 + (y-b)^2 = r^2$ is used, where (a,b) is the center and r is the radius.

The Hough Transform works by mapping edge points (x,y) into a **parameter space** (also called **accumulator**), where each point votes for all possible parameter combinations that could represent a shape passing through it.

After accumulating votes for all edge points, the algorithm looks for **peaks** in the accumulator array — these peaks represent **intersections of many votes**, indicating likely shape configurations. A **threshold** can be applied to ignore weak detections and keep only the most confident shapes.

Hough Transform is **robust to noise**, since random noise points don't vote consistently. It also **tolerates missing data**, because a shape can still be detected as long as **enough points agree** on its presence.

Hough Transform works best for **simple, parameterizable shapes** (lines, circles). For more complex shapes (like rectangles or stars), other methods like **contour analysis** or **machine learning** are preferred.

23. In PCA, we use the eigenvectors of the covariance matrix to project the data.
24. The frequency component $(0,0)$ in 2D DFT represents the average brightness of the image.
25. In transform domain, an image is represented as a combination of basis functions.
26. Orthogonal basis functions have dot products equal to 0
27. Hough Transform maps points in image space to curves in parameter space.
28. The peak of the accumulator array in Hough space suggests the most likely shape configuration
31. How are basis functions used to represent images in the transform domain?

In the transform domain, an image is represented as a **weighted sum of basis functions**. Each basis function captures a **specific pattern** (e.g., a frequency, edge orientation, or spatial structure), and the weights indicate how much of each pattern is present in the image.

For example:

- In the **Discrete Fourier Transform (DFT)**, the basis functions are **sinusoids** of varying frequencies and orientations.
- In the **Discrete Cosine Transform (DCT)**, the basis functions are cosine waves.
- In **PCA**, the basis functions are the **principal components** (eigenfaces in face recognition).

This approach allows complex image content to be decomposed into **simpler components**, which can be analyzed, modified, or compressed more efficiently than working in the pixel (spatial) domain.

32. Why are orthogonal basis functions advantageous in image transformation and compression?

Orthogonal basis functions have the property that their **inner product is zero**, meaning they are **mutually independent**. This provides several advantages:

1. **No redundancy**: Each coefficient captures **unique information** about the image.
2. **Efficient reconstruction**: The original image can be exactly reconstructed (if all components are used), with **minimal computational overlap**.
3. **Energy compaction**: Most of the important image information (energy) is concentrated in **a few coefficients**, allowing for effective **compression** by discarding small components.
4. **Simplified math**: Operations like projection, filtering, and inverse transforms become simpler and more stable.

These properties are the reason why orthogonal transforms like DFT, DCT, and PCA are widely used in **compression standards** (e.g., **JPEG**) and **image analysis**.

34. How does the notch filter help reduce periodic noise in the frequency domain?

A **notch filter** is used to reduce **periodic noise** in images by suppressing specific frequency components in the frequency domain.

1. **Apply the 2D Fourier Transform** to convert the image into the frequency domain.

2. **Identify bright spots** in the spectrum that correspond to periodic noise frequencies.
3. **Create a mask (notch filter)** that sets those specific frequency components to zero while preserving the rest of the spectrum.
4. **Multiply the frequency-domain image** by the filter mask to remove noise components.
5. **Apply the inverse 2D Fourier Transform** to reconstruct the denoised image in the spatial domain.

This process effectively removes **periodic patterns** while preserving the original image content.

The **frequency domain** is an alternative way of representing an image (or any signal), where instead of describing it in terms of pixel intensities at spatial locations, it is described in terms of **how those intensities vary**, i.e., in terms of **frequencies**.

In the frequency domain, an image is represented as a **combination of sinusoidal patterns** (basis functions) of different frequencies and directions. These components are obtained using a **2D Discrete Fourier Transform (2D DFT)**, which converts the image into a complex matrix of frequency coefficients. Each coefficient represents the **amplitude and phase** of a particular frequency in the image.

The **magnitude** of a frequency coefficient tells us how much that frequency contributes to the overall image structure, while the **phase** determines the spatial position of those features. The resulting representation provides a powerful way to understand and manipulate images.

For example, low frequencies correspond to smooth variations and general shapes, while high frequencies represent sharp edges, textures, and noise. This allows targeted operations such as **low-pass filtering** (to smooth the image), **high-pass filtering** (to enhance edges), or even **notch filtering** (to remove periodic noise).

The frequency domain also simplifies many image processing tasks. For instance, convolution in the spatial domain becomes simple multiplication in the frequency domain, significantly reducing computation. Additionally, since many frequency components may be negligible, the frequency domain is widely used in **image compression** (e.g., JPEG uses DCT) and **denoising**.

Overall, the frequency domain is an essential concept in image processing, offering better interpretability and control over image structure than the spatial domain in many applications.

The **2D Discrete Fourier Transform (2D DFT)** is a mathematical tool used to convert an image from the **spatial domain** (where it is represented as pixel intensities) into the **frequency domain**, where it is represented as a set of frequency components.

The transform expresses the image as a **weighted sum of sinusoidal basis functions** with different frequencies and directions. Each frequency component has both **magnitude** (indicating how much it contributes to the image) and **phase** (indicating spatial alignment).

The formal definition of the 2D DFT involves a double summation over all pixels, producing a complex-valued matrix that contains all the frequency information of the image. The **central frequency** (DC component) typically represents the average brightness, while frequencies farther from the center correspond to faster variations like edges or noise.

2D DFT is widely used in image processing for **filtering, denoising, and compression**. Many periodic disturbances such as vertical striping can be easily detected and suppressed in the frequency domain using masks (e.g., notch filters).

Moreover, thanks to the **separability** of the DFT, it can be computed efficiently by applying 1D DFT first to the rows, and then to the columns. This is what makes it feasible to apply DFT to large images using the **Fast Fourier Transform (FFT)**.

The **separability** property of the 2D Discrete Fourier Transform (DFT) allows it to be computed by applying **1D DFT** first along the rows and then along the columns (or vice versa).

This not only simplifies implementation, but also **enables the use of the Fast Fourier Transform (FFT)** — a highly efficient algorithm that reduces computational complexity from $O(N^2)$ to $O(N \log N)$.

The inverse DFT (IDFT) allows reconstruction of the original image, making the process fully reversible. This makes 2D DFT a fundamental tool in modern image analysis and signal processing.

26. Clustering is an example of unsupervised learning.
27. K-means minimizes a cost function based on euclidean distance.
28. In hierarchical clustering, a dendrogram visualizes how clusters merge.
29. Single linkage considers the closest distance between two clusters.
30. The elbow point is found by plotting SSE against number of clusters.
31. A silhouette score close to 1 indicates a point is well matched to its cluster.
32. Describe the K-means clustering algorithm step-by-step. Include initialization, assignment, update, and termination.

K-Means Clustering Algorithm – Step-by-Step:

K-Means is a **centroid-based** clustering algorithm that partitions a dataset into **k clusters** by minimizing the **within-cluster variance** (inertia). It iteratively updates cluster assignments based on **proximity to centroids**.

Steps:

1. Initialization:

Randomly choose **k initial centroids** (either randomly from the data points or using smarter methods like k-means++).

2. Assignment step:

Assign each data point to the **nearest centroid** using a distance metric (typically

Euclidean distance).

This creates **k clusters** based on similarity.

3. Update step:

Recalculate the centroid of each cluster by taking the **mean of all points** assigned to that cluster.

4. Termination condition:

Repeat steps 2 and 3 until one of the following:

- **Centroids do not change** (convergence),
 - or a **maximum number of iterations** is reached.
-

 **Applications:**

- Image compression and segmentation
- Customer segmentation in marketing
- Anomaly detection
- Document or gene clustering

33. Explain at least three limitations of the K-means algorithm and how they can be addressed.

Limitations of the K-Means Algorithm and How to Address Them:

1. Varying cluster shapes, sizes, or densities

K-Means assumes that clusters are **spherical and equally sized**, which makes it ineffective for detecting **non-globular** or **unevenly distributed** clusters.

 **Solution:** Use algorithms like **DBSCAN** or **Gaussian Mixture Models (GMM)** that can model clusters of different shapes and densities.

2. Sensitivity to outliers

K-Means minimizes **squared Euclidean distances**, so **outliers** can strongly influence the position of centroids, leading to distorted clusters.

 **Solution:** Use **K-Medoids**, which selects actual data points as cluster centers, or apply **outlier removal** techniques before clustering.

3. Poor initialization of centroids

The final clustering result can vary significantly depending on the **initial centroids**. Poor initialization can lead to **local minima** and unstable results.

 **Solutions:**

- Use **K-Means++** for smarter initial centroid selection
- Run the algorithm **multiple times** with different initializations and pick the best result
- Optionally, use **hierarchical clustering** to estimate better starting points

34. Compare centroid-based, density-based, and hierarchical clustering. Include when each is preferred.

Comparison of Clustering Methods:

1. Centroid-based clustering (e.g., K-Means):

- **Idea:** Partitions data into k clusters by minimizing the distance between points and the cluster **centroid** (mean).
- **Assumptions:** Clusters are **spherical**, of **similar size and density**.
- **Advantages:** Fast, efficient for large datasets.
- **Limitations:** Fails with **non-globular** clusters and **outliers**.
- *Preferred when:* You know the number of clusters in advance and clusters are roughly spherical (e.g., customer segmentation by purchase volume).

2. Density-based clustering (e.g., DBSCAN):

- **Idea:** Forms clusters as **dense regions** separated by low-density areas.
- **Advantages:** Can detect **arbitrary-shaped clusters** and **handle noise/outliers** well.
- **Limitations:** Struggles with varying density; requires tuning of ϵ (neighborhood size).
- *Preferred when:* Data contains **non-linear structures**, **noise**, or **unknown number of clusters** (e.g., spatial data, anomaly detection).

3. Hierarchical clustering (e.g., Agglomerative):

- **Idea:** Builds a **tree (dendrogram)** by recursively merging or splitting clusters based on distance.

- **Advantages:** No need to specify number of clusters initially; provides **multi-level hierarchy**.
- **Limitations:** Computationally expensive for large datasets; sensitive to noise.
- *Preferred when:* You need a **hierarchical structure** of data or want to **visualize cluster merging** (e.g., gene expression analysis).

35. What are linkage methods in hierarchical clustering? Compare single, complete, average, and ward linkage.

Hierarchical is a clustering algorithm or technique used to group data points into a hierarchical tree-like structure based on their similarity or dissimilarity. It creates a hierarchy of clusters where data points are gradually merged into larger clusters.

In hierarchical clustering, **linkage methods** define how the distance between two clusters is calculated during the merging process. The choice of linkage affects the **shape, size, and structure** of the resulting clusters.

1. Single Linkage

- Computes the distance between the **closest pair** of points (one from each cluster).
- Tends to form **elongated, chain-like clusters**.
-  Sensitive to **noise and outliers**.

2. Complete Linkage

- Uses the distance between the **farthest pair** of points in two clusters.
- Produces **compact and spherical** clusters.
- Less sensitive to outliers than single linkage.

3. Average Linkage

- Calculates the **average distance** between all pairs of points across clusters.
- Balances the behavior of single and complete linkage.
- Often results in **well-formed and balanced** clusters.

4. Centroid Linkage

- Measures distance between the **centroids** (mean vectors) of clusters.

- May lead to undesirable results such as **inversion** (non-monotonic merging), but can be efficient for certain types of data.

5. Ward Linkage

- Minimizes the **increase in total within-cluster variance** when merging two clusters.
- Favors merging clusters that result in **minimum loss of information**.
- Produces **compact, spherical clusters** and is less sensitive to noise.
- Often considered the **most effective** for quantitative data.

36. How is the silhouette score calculated and interpreted? Provide an example scenario.

Measures similarity within clusters compared to similarity between clusters. Ranges from -1 to 1, with higher scores indicating better clustering. Higher average scores imply well-defined clusters, while lower scores suggest clustering issues. Relatively easy to interpret and robust to outliers.

Can be found by formula $s(i) = (b(i) - a(i)) / \max(a(i), b(i))$ where $a(i)$ is average distance of data point i to all other data points in the same cluster (represents cohesion). $b(i)$ shows smallest distance to another cluster (separation).

The Silhouette Score for the entire dataset is the mean of the silhouette scores for all data points.

Example:

Suppose you cluster customer data into 3 segments. After computing the silhouette score, you get an average value of **0.72**.

This indicates that customers within each cluster are similar to each other and well-separated from other clusters — meaning **clustering is likely good**.

37. Describe the Elbow Method. How does it help choose the number of clusters in K-means?

Method to determine the optimal number of clusters (k) for a dataset. Involves plotting the sum of squared distances (SSD) against different k values. The "elbow point" on the plot, where the SSD rate of decrease significantly slows down, is used to estimate the optimal k value. Balances clustering quality with avoiding overfitting by selecting an appropriate number of clusters.

38. What preprocessing and postprocessing techniques can improve clustering quality?

- Pre-processing:

- Data Normalization: Standardize the data to ensure all features have the same scale.
- Outlier Removal: Identify and remove outliers from the dataset.

- Post-processing:

- Eliminate Small Clusters: Remove clusters that contain a minimal number of data points.
- Split 'Loose' Clusters: Identify clusters with relatively high SSE.
- Split these clusters to create more well-defined sub-clusters.
- Merge Close Clusters: Identify clusters that are close to each other in feature space. Merge such clusters, particularly when they have relatively low SSE.

39. Explain how dendograms can be used to determine the number of clusters and detect outliers.

A **dendrogram** is a tree-like diagram that visualizes the process of hierarchical clustering. It shows how data points or clusters are progressively merged based on similarity.

 **Determining the number of clusters:**

- By drawing a **horizontal cut** across the dendrogram at a certain height (distance threshold), you can determine how many **clusters** exist.
- The number of vertical lines intersected by the cut equals the number of clusters.
- A good cut is often made where the vertical distance between merged clusters is **large** (i.e., a big jump in linkage distance), indicating well-separated groups.

 **Detecting outliers:**

- Outliers appear as **points that join clusters at a very high linkage distance** — they remain isolated for most of the merging process.
- These are seen as **branches that extend far** before connecting to the rest of the tree, indicating that the point is not similar to any existing cluster.

Example:

In a dendrogram of customer data, if one point connects at the very top of the tree while others form compact groups below, that point is likely an **anomaly or outlier**.

36. The activation function that returns $\max(0, x)$ is called ReLU.
37. A neuron 'fires' when the output of its activation function is non-zero.
38. The process of adjusting weights using loss gradients is called backpropagation
39. The problem of very small gradients in early layers is known as the vanishing gradient problem.
40. A technique that randomly disables neurons during training is called dropout
41. The derivative of ReLU for positive input is 1
42. In PyTorch, the shape of an MNIST input image after flattening is **(1, 784)**
43. The optimizer that combines momentum and adaptive learning is Adam.
44. CrossEntropyLoss in PyTorch combines softmax and negative log-likelihood loss.
45. The architecture where multiple layers are stacked one after another is called a feedforward neural network.
46. Explain the difference between ReLU, sigmoid, and tanh activation functions. When should each be used?

1. ReLU (Rectified Linear Unit):

$$f(x) = \max(0, x)$$

- **Behavior:** Outputs 0 for negative inputs, linear for positive.
- **Advantages:** Fast to compute, helps avoid vanishing gradient.
- **Used in:** Hidden layers of **deep networks**, especially in CNNs and modern architectures.

- **Limitation:** Can cause “dead neurons” (outputs stuck at 0).

2. Sigmoid:

- **Behavior:** Maps input to range (0, 1).
- **Used in:** **Binary classification output layer** (e.g., logistic regression).
- **Limitation:** **Vanishing gradients**, especially for large/small inputs.

3. Tanh (Hyperbolic Tangent):

- **Behavior:** Maps input to (-1, 1).
- **Used in:** Hidden layers in **shallow networks**; better than sigmoid since it's zero-centered.
- **Limitation:** Still suffers from vanishing gradient in deep networks.

Summary:

- Use **ReLU** in hidden layers for deep nets.
- Use **sigmoid** for binary outputs.
- Use **tanh** when zero-centered activations are preferred, and model is not very deep.

47. Describe the complete process of training a neural network using backpropagation and gradient descent.

The training process involves minimizing a **loss function** by adjusting the network's weights using **backpropagation** and **gradient descent**.

Step-by-step Process:

1. Forward Pass:

- Input data passes through each layer.
- Activations are computed using weights and activation functions.
- Output is generated and compared to the target using a **loss function**.

2. Compute Loss:

- Measure how far the predicted output is from the true label (e.g., using MSE or cross-entropy).

3. Backward Pass (Backpropagation):

- Compute **gradients of the loss** with respect to each weight using the **chain rule**.
- Error is propagated **from output to input layer**, layer by layer.

4. Gradient Descent Update

5. Repeat:

- Iterate over the training data in **epochs**.
- Optionally use **mini-batches** (mini-batch gradient descent) and apply **regularization**.

Over time, the model learns weights that **minimize the loss** and improve predictions on unseen data.

48. What causes the vanishing and exploding gradient problems? How can they be mitigated?

Vanishing Gradient Problem

Causes:

- Occurs with activation functions like sigmoid and tanh due to small derivatives
- Gradients shrink exponentially during backpropagation in deep networks

Consequences:

- Early layers doesn't learn, resulting in poor model performance.

Solutions:

- Use ReLU or Leaky ReLU (derivatives = 0/1 or small constant)

Exploding Gradient Problem

Causes:

- Large weight values (>1) amplify gradients during backpropagation
- Common in deep networks with unbounded activations (e.g., linear)

Consequences:

- Numerical instability (NaN/Inf loss)
- Model parameters diverge, leading to erratic predictions

Solutions:

- Gradient clipping: Limit gradient magnitude during updates
- Weight regularization (e.g., L2) to constrain magnitudes

49. Describe the role of the loss function in neural networks. How does it influence training?

The **loss function** quantifies how far the model's predictions are from the true values. It serves as the **objective function** that the training process seeks to **minimize**.

◆ **Role in training:**

- After each forward pass, the loss function computes the **error** between predicted and actual outputs.
- During backpropagation, the **gradients of the loss** with respect to each weight are calculated.
- These gradients guide **weight updates** via an optimization algorithm (e.g., SGD).

📌 **Influence on training:**

- The choice of loss function affects **what the model learns**:
 - **Cross-entropy** for classification
 - **Mean Squared Error (MSE)** for regression
- A poorly chosen loss function can lead to **slow convergence** or **suboptimal learning**.
- Properly defined loss ensures the model learns **meaningful patterns** from the data.

50. Compare SGD, batch gradient descent, and Adam optimizer. What are the pros and cons of each?

Optimizer	Pros	Cons
Batch Gradient Descent	Uses the entire dataset for each update. <input checked="" type="checkbox"/> Stable convergence <input checked="" type="checkbox"/> Exact gradients	<input checked="" type="checkbox"/> Slow on large datasets <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> Not scalable
Stochastic Gradient Descent (SGD)	Updates weights using one data point at a time . <input checked="" type="checkbox"/> Fast updates <input checked="" type="checkbox"/> Can escape local minima	<input checked="" type="checkbox"/> High variance in updates <input checked="" type="checkbox"/> Noisy convergence
Adam (Adaptive Moment Estimation)	Combines momentum + adaptive learning rates . <input checked="" type="checkbox"/> Fast convergence <input checked="" type="checkbox"/> Works well out of the box	<input checked="" type="checkbox"/> May converge to suboptimal solutions <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> Sensitive to learning rate sometimes

Summary:

- Use **Batch GD** when the dataset is small and you want stability.
- Use **SGD** for online learning or streaming data.
- Use **Adam** for deep learning tasks and large-scale problems — it's often the default choice due to adaptive updates.

51. What is the function of dropout during training? Why does it help generalization?

Dropout is a regularization technique used during the training of neural networks to reduce **overfitting** and improve **generalization**.

How Dropout Works (with p)

During training, each neuron in a layer is **randomly deactivated (dropped)** with probability $1-p$ and **kept active** with probability p .

This means:

- If $p=0.5$, then 50% of neurons are randomly "turned off" during each forward pass.
- This dropout is applied **independently** for each training example and each iteration.

- At **test time**, all neurons are kept active, but their outputs are typically **scaled by p** to compensate for the reduced activation during training.
-

✓ Why Dropout Helps Generalization

- It prevents the network from **relying too heavily on any single neuron**.
 - Encourages the network to learn **redundant, distributed representations**.
 - Dropout effectively trains an **ensemble of subnetworks** and averages their predictions.
 - This makes the model more **robust** and reduces the risk of overfitting.
-

📌 Example

In a dense layer with dropout $p=0.5$, the network learns to perform well even when **half of the neurons are missing**.

As a result, it generalizes better to new, unseen data, even if training accuracy slightly decreases.

52. Describe the structure and forward pass in a multilayer perceptron (MLP).

A **Multilayer Perceptron (MLP)** is a **fully connected feedforward neural network**, composed of:

◆ Structure:

1. **Input layer:** Receives the input features.
2. **One or more hidden layers:**
 - Each neuron is connected to all neurons in the previous layer.
 - Uses activation functions (e.g., ReLU, tanh) to introduce non-linearity.
3. **Output layer:**
 - Produces the final prediction.
 - Activation depends on task (e.g., softmax for classification, none or sigmoid for regression).

Forward pass:

1. Inputs are multiplied by weights and summed with biases.
2. Result is passed through an activation function to produce output for the next layer.
3. This process repeats through all layers until the final prediction is produced.

 The MLP learns to model complex relationships between inputs and outputs by adjusting weights to minimize the loss during training.

53. How does momentum help in accelerating gradient descent? Illustrate with an example.

Momentum is an optimization technique that improves gradient descent by **accumulating past gradients** to build up speed in consistent directions and dampen oscillations in noisy or curved regions.

Instead of updating weights based only on the current gradient, momentum **adds a fraction of the previous update**, creating a “rolling effect.”

Why It Helps:

- Speeds up convergence in directions of consistent gradients.
- **Reduces oscillation** in narrow valleys or along steep curves.
- Helps escape **local minima** and plateaus.

Example:

In a loss surface shaped like a ravine, standard gradient descent may oscillate left-right. Momentum smooths the updates and **accelerates in the correct direction** (like pushing a ball downhill with increasing speed).

54. What is overfitting in neural networks, and what techniques are commonly used to prevent it?

Overfitting occurs when a neural network learns not only the underlying patterns in the training data but also **noise and random fluctuations**.

As a result, the model performs well on training data but **poorly on unseen test data**.

▶ **Symptoms:**

- **High training accuracy**, but **low validation/test accuracy**
 - Large gap between training and validation loss
-

✓ **Common Techniques to Prevent Overfitting:**

1. **Regularization** (L1/L2): Adds penalties for large weights.
 2. **Dropout**: Randomly disables neurons during training.
 3. **Early stopping**: Stops training when validation loss stops improving.
 4. **Data augmentation**: Increases training diversity (e.g., flips, rotations).
 5. **Reduce model complexity**: Use fewer layers/neurons.
 6. **Cross-validation**: Helps tune hyperparameters effectively.
-

📌 **Goal:** Improve **generalization** by forcing the model to learn robust patterns that work beyond the training set.

57. What are the benefits of using a GPU for training deep neural networks?

GPUs (Graphics Processing Units) are widely used to accelerate training of deep neural networks due to their **parallel processing capabilities**.

✓ **Key Benefits:**

1. **Massive parallelism**
 - GPUs have thousands of cores that can perform **many operations simultaneously**, ideal for matrix multiplications in neural nets.
2. **Faster training**

- Training deep models on large datasets can be **10–100× faster** compared to CPUs.

3. Better memory bandwidth

- GPUs handle **large amounts of data** more efficiently, speeding up data loading and computation.

4. Optimized deep learning libraries

- Libraries like **TensorFlow**, **PyTorch**, and **CUDA** are designed to leverage GPU hardware for maximum performance.
-

📌 **Conclusion:** GPUs enable the practical training of complex deep learning models that would otherwise take days or weeks on CPUs.

58. Explain how the architecture of MNIST digit classification is structured. Include layers and activation functions.

MNIST is a classic dataset of **28×28 grayscale images** of handwritten digits (0–9). The architecture for classifying MNIST digits is typically a **convolutional neural network (CNN)**, but can also be a fully connected network (MLP).

💡 Example 1: Multilayer Perceptron (MLP) Architecture

1. Input layer:

- 784 neurons (28×28 flattened pixels)

2. Hidden layer 1:

- 128 neurons, **ReLU activation**

3. Hidden layer 2:

- 64 neurons, **ReLU activation**

4. Output layer:

- 10 neurons (one per digit), **Softmax activation**
-

Example 2: CNN Architecture

1. **Input:** $1 \times 28 \times 28$ grayscale image
 2. **Conv layer 1:** 32 filters (3×3), **ReLU** \rightarrow **MaxPooling (2x2)**
 3. **Conv layer 2:** 64 filters (3×3), **ReLU** \rightarrow **MaxPooling (2x2)**
 4. **Flatten**
 5. **Fully connected (Dense) layer:** 128 units, **ReLU**
 6. **Output layer:** 10 units, **Softmax**
-

 The model learns to extract spatial features (edges, shapes) from the input and classify them into one of 10 digit classes.

Both versions achieve high accuracy (>98%) on the MNIST test set.

60. Describe a real-world application where deep neural networks outperform traditional ML methods, and explain why.

Application: Image recognition (e.g., facial recognition, object detection in photos)

Why DNNs outperform traditional ML here:

1. **Automatic feature extraction**
 - o Traditional ML (like SVM, logistic regression) requires **manual feature engineering** (e.g., edge detectors, HOG descriptors).
 - o Deep neural networks, especially **CNNs**, learn **hierarchical features** (edges \rightarrow shapes \rightarrow objects) **directly from raw pixels**.
 2. **Scalability**
 - o DNNs can handle **large-scale datasets** (like ImageNet) and high-dimensional input, where traditional methods fail or plateau.
 3. **Transfer learning**
 - o Pretrained deep models can be **fine-tuned** on new image tasks with minimal data, improving accuracy and reducing development time.
-

❖ **Example:** In facial recognition (e.g., FaceNet, DeepFace), deep networks achieve **human-level accuracy**, while traditional ML methods struggle with lighting, pose, and alignment variability.

36. In CNNs, the operation of sliding a filter over input is called convolution.
37. A filter that reduces spatial size without learning parameters is called pooling.
38. The (height, width) of a 3x3 kernel is (3, 3)
39. A ReLU activation returns 0 for negative inputs.
40. The layer that maps each feature map to one value is **Global** pooling.
41. The architecture that introduced skip connections is called ResNet
42. CNNs are primarily used in grid data.
43. In transfer learning, we reuse weights from a pretrained model.
44. Data augmentation increases the effective size of the training data.
45. Autoencoders are trained to reproduce their outputs
46. Describe the complete architecture of a typical CNN. Mention the role of convolution, activation, pooling, and fully connected layers.

A **Convolutional Neural Network (CNN)** is a deep learning architecture especially effective for image and spatial data. It processes input through **convolutional**, **activation**, **pooling**, and **fully connected layers** to learn hierarchical features.

❖ 1. Convolutional Layers

- Apply **filters (kernels)** that slide over the input image to detect local patterns like edges, corners, or textures.
- Each filter generates a **feature map** that highlights the presence of specific features.

- These layers preserve **spatial relationships** and drastically reduce the number of parameters compared to dense layers.
-

◆ 2. Activation Functions

- After convolution, a **non-linear activation** function is applied, usually **ReLU (Rectified Linear Unit)**:

$f(x) = \max(0, x)$ Introduces **non-linearity** into the model so it can learn complex patterns.

◆ 3. Pooling Layers

- Reduce the **spatial dimensions** (width and height) of feature maps while keeping the most important information.
 - Common types: **Max pooling** (keeps max value in a window), **Average pooling**.
 - Helps reduce computation, prevent overfitting, and improve translation invariance.
-

◆ 4. Fully Connected (Dense) Layers

- ❑ After the convolutional and pooling layers have extracted spatial features, the output (i.e., feature maps) is **flattened** into a 1D vector.
 - ❑ This vector is passed through one or more **dense (fully connected) layers**, where **each neuron is connected to all neurons of the previous layer**.
 - ❑ These layers act as the "**decision-making**" part of the network, performing **high-level reasoning** based on the extracted features.
 - ❑ The **final dense layer** produces the output:
 - Uses **Softmax** activation for **multi-class classification**.
 - Uses **Sigmoid** activation for **binary classification** or **multi-label problems**.
 - ❑ Fully connected layers introduce many parameters and are often regularized (e.g., with **dropout**) to prevent overfitting.
-

✓ Typical CNN Flow:

1. **Input image** (e.g., 32×32×3)

2. → Conv layer + ReLU
 3. → Pooling layer
 4. → (repeat Conv + Pool multiple times)
 5. → Flatten
 6. → Fully connected layer(s)
 7. → Output layer (with softmax/sigmoid)
-

👉 This layered architecture allows CNNs to learn from **low-level to high-level features**, making them extremely powerful for tasks like image classification, object detection, and face recognition.

47. Explain how convolution operation works with an example of a 3x3 filter sliding over an image.

The **convolution operation** in CNNs extracts local features from input data (e.g., images) by sliding a **filter** (also called kernel) across the input and computing a **dot product** between the filter and the input region.

Each filter learns to detect a **specific pattern** like an edge, texture, or shape.

💡 **Basic Idea of Convolution (3×3 filter):**

At each position:

1. Select a **3×3 region** from the input.
 2. Multiply each input value by the corresponding filter weight.
 3. Sum the results to get **one number** (activation).
 4. Move the filter and repeat.
-

📐 **Formal Output Size Formula:**

Output Size=($N+2P-F$)/Stride+1

🧠 **Practical Example:**

- **Input volume:** $32 \times 32 \times 3$ (height \times width \times channels)
- **Filter:** 10 filters of size $5 \times 5 \times 3$
- **Stride:** 1
- **Padding:** 2

$$32 + 2 \times 2 - 5 + 1 = 32$$

 So the output volume is:

$32 \times 32 \times 10$ \rightarrow 10 feature maps of size 32×32 (one per filter)

Number of Parameters in This Layer:

- Each filter has:

$$5 \times 5 \times 3 = 75 \text{ weights} + 1 \text{ bias} = 76 \text{ parameters}$$

With 10 filters:

$$76 \times 10 = 760 \text{ parameters total}$$

 **Conclusion:** Convolution allows CNNs to scan the input for local patterns using small, trainable filters. The number of filters defines the **depth of the output volume**, while stride and padding control **spatial resolution**.

48. Compare max pooling and average pooling. When is each preferable?

Pooling layers reduce the spatial dimensions of feature maps, helping with generalization and computational efficiency.

Max Pooling

- Takes the **maximum value** in each window (e.g., 2×2).
- Captures the **strongest activation** (e.g., edge, texture).
- Helps with **edge detection** and preserving **salient features**.

Preferred when:

- You want to retain **dominant patterns**.
- Common in **image classification** (e.g., CNNs for visual recognition).

◆ **Average Pooling**

- Computes the **mean** value of each window.
- Produces **smoother** feature maps.
- Less aggressive than max pooling.

✓ **Preferred when:**

- You want to retain **overall texture** or background information.
 - Sometimes used in **signal smoothing** or **compressed representations**.
-

📌 **Summary:**

- **Max pooling** = sharper, more abstract features → better for detection.
- **Average pooling** = smoother, more general info → better for preserving background or reducing noise.

49. What are the main causes and consequences of vanishing gradients in deep CNNs?
How does ResNet solve this?

▼ **Causes of Vanishing Gradients in Deep CNNs:**

- Use of activation functions like **sigmoid** or **tanh**, whose derivatives are **very small** for large or small inputs.
 - In very deep networks, gradients are **multiplied through many layers**, which causes them to **shrink exponentially**.
 - As a result, **early layers receive almost no gradient signal** during backpropagation.
-

⚠ **Consequences:**

- **Early layers fail to learn**, leading to poor feature extraction.
- The network **trains very slowly** or fails to converge.
- Model **underfits**, even if deep.

How ResNet Solves This (Residual Networks):

- Introduces **skip connections (shortcut paths)** that bypass one or more layers:

$$y=F(x)+x$$

- Instead of learning the full mapping, the network learns the **residual** $F(x)=H(x)-x$
- Gradients can **flow directly through the skip paths**, which helps **preserve the signal** even in very deep networks.
- This enables training of **extremely deep CNNs** (e.g., 50, 101, 152 layers) with improved convergence and accuracy.

50. Describe the architecture and novelty of LeNet. How did it influence later CNN designs?

LeNet-5 is one of the earliest Convolutional Neural Network architectures, developed by **Yann LeCun** in 1998 for **handwritten digit recognition** (e.g., MNIST).

Architecture Overview:

1. **Input:** 32×32 grayscale image
2. **Conv layer 1:** 6 filters (5×5), **tanh**, output: $28 \times 28 \times 6$
3. **Avg pooling:** 2×2 , output: $14 \times 14 \times 6$
4. **Conv layer 2:** 16 filters (5×5), **tanh**, output: $10 \times 10 \times 16$
5. **Avg pooling:** 2×2 , output: $5 \times 5 \times 16$
6. **Fully connected layers:** $120 \rightarrow 84 \rightarrow 10$ (softmax)

Key Innovations:

- Introduced the idea of using **convolution + pooling + fully connected layers**.
- Used **parameter sharing** to reduce complexity.
- Demonstrated the power of **end-to-end training** of features and classifier.

Influence on Later CNNs:

- Inspired modern architectures like **AlexNet**, **VGG**, and **ResNet**.
 - Validated the idea that deep networks can learn **hierarchical representations**.
 - Established core building blocks (conv → pool → dense) that are still used today.
-

 LeNet was a **foundational breakthrough** that made modern deep learning for images possible.

51. What innovations did AlexNet introduce to CNN training and performance?

AlexNet (2012) was a landmark CNN architecture that **won the ImageNet challenge** and triggered the modern deep learning revolution in computer vision.

Key Innovations:

1. **Use of ReLU activations**
 - Replaced sigmoid/tanh with **ReLU**:
 - Allowed for **faster training** and mitigated vanishing gradient.
2. **Dropout for regularization**
 - Prevented **overfitting** by randomly deactivating neurons during training.
3. **Data augmentation**
 - Applied transformations (e.g., flips, crops) to synthetically enlarge the training set.
4. **GPU acceleration**
 - Trained on two GPUs in parallel, showing that **deep networks could scale** with hardware.
5. **Large-scale deep architecture**
 - Used 8 layers (5 convolutional, 3 fully connected) — deeper than LeNet.
 - **First to show deep CNNs could outperform traditional methods on large datasets.**

📌 **Impact:** AlexNet proved that deep CNNs trained with enough data and GPU power could achieve **state-of-the-art results**, setting the foundation for modern architectures.

52. How does VGGNet achieve deeper architectures using small kernels? What are the advantages?

VGGNet (2014) is a deep CNN that emphasizes simplicity and depth by stacking **small 3×3 convolutional filters**.

💡 **How It Builds Depth:**

- Instead of large filters (e.g., 7×7), VGG uses **multiple 3×3 filters** in sequence.
- Example: Two 3×3 conv layers approximate the receptive field of one 5×5 filter.
- Stacking **3 conv layers (3×3)** yields an effective **7×7 receptive field**, but with **fewer parameters** and more non-linearities.

✓ **Advantages of Small Kernels:**

1. **Deeper architectures**
 - Easier to build very deep networks (VGG-16, VGG-19).
2. **More non-linearity**
 - More activation functions per receptive field → **richer representations**.
3. **Fewer parameters**
 - A stack of 3×3 filters is **more efficient** than a single large filter.
4. **Better generalization**
 - Small filters + depth + simplicity improves performance on large-scale datasets (e.g., ImageNet).

📌 VGGNet influenced many later architectures by showing that **depth + small filters = strong performance**, forming the basis for models like ResNet and ZFNet

53. Describe the structure and idea behind inception modules in GoogLeNet. Why use multiple kernel sizes?

GoogLeNet (2014) introduced the **Inception module**, which allowed the network to **capture features at multiple spatial scales** in parallel, making it both deep and computationally efficient.

Inception Module Structure:

Each module performs **parallel operations** on the same input:

- **1×1 convolution** (for dimensionality reduction)
- **3×3 convolution**
- **5×5 convolution**
- **3×3 max pooling**

All outputs are **concatenated** along the depth axis to form the final output of the module.

Why Use Multiple Kernel Sizes?

- Different kernel sizes **capture different types of features**:
 - **1×1**: reduces dimensionality, adds non-linearity.
 - **3×3**: captures **mid-level patterns**.
 - **5×5**: captures **larger, more global features**.
 - **Pooling**: adds **robustness and abstraction**.

 Instead of choosing one kernel size, the network **learns which scale works best** at each layer.

Benefits:

- Encourages **multi-scale learning**.
- Reduces computation by using **1×1 convolutions** before expensive ones.
- Enables very **deep networks** without excessive computational cost.

54. How do skip connections in ResNet help in training very deep networks?

ResNet (Residual Networks) introduced **skip connections**, which directly pass the input of a layer to a deeper layer.

How Skip Connections Work:

Instead of learning $H(x)$, the layer learns the **residual**:

$$y = F(x) + x$$

Where:

- $F(x)$: the transformation (convolution, BN, ReLU, etc.)
 - x : input passed through unchanged
-

Why This Helps Training:

1. **Improved gradient flow**

- Gradients can flow **directly through skip paths**, mitigating **vanishing gradients** in deep networks.

2. **Easier optimization**

- Learning residuals (differences) is often easier than learning full transformations.

3. **Network behaves like an ensemble**

- Skip connections allow multiple paths for signal propagation, leading to **better generalization**.
-

 **Result:** ResNet enables training of **very deep models** (e.g., 50, 101, 152 layers) with **better accuracy and faster convergence** compared to plain deep networks.

55. What is data augmentation, and what are common techniques used? Why is it important?

Data augmentation is the process of artificially increasing the size and diversity of a training dataset by applying **transformations** to the original data.

✓ **Common Techniques (especially in image data):**

- **Flipping** (horizontal/vertical)
 - **Rotation** (e.g., $\pm 15^\circ$)
 - **Cropping and resizing**
 - **Zooming in/out**
 - **Color jittering** (brightness, contrast, saturation)
 - **Translation** (shifting left/right/up/down)
 - **Adding noise or blurring**
-

🎯 **Why It's Important:**

1. **Prevents overfitting:** Augmented data exposes the model to more variation, making it generalize better.
 2. **Improves robustness:** The model learns to recognize objects in different positions, angles, or lighting conditions.
 3. **Reduces need for large labeled datasets:** Especially useful when collecting labeled data is expensive.
-

📌 Widely used in computer vision tasks like classification, detection, and segmentation to improve performance without collecting new data.

56. Define transfer learning. Give an example of how it can reduce training time in image classification.

Transfer learning is a technique where a model trained on one task (usually on a **large dataset**) is **reused or fine-tuned** for another related task.

🧠 **Key Idea:**

- Early layers in deep neural networks learn **generic features** (edges, textures), which are transferable across tasks.
 - Instead of training from scratch, we **reuse a pretrained model** and only **fine-tune the top layers** for a new task.
-

✓ Example in Image Classification:

- Use a model like **ResNet-50 pretrained on ImageNet**.
- For a custom task (e.g., classifying medical X-rays with 10 categories), you:
 - . Replace the final classification layer.
 - . Freeze the early layers (optional).
 - . Fine-tune the top few layers on your smaller dataset.

⚡ Benefits:

- **Significantly reduces training time**.
- Requires **less data** and **computational resources**.
- Achieves **high accuracy** even with limited labeled data.

📌 Transfer learning is especially effective when your dataset is small, but similar in nature to a large public dataset.

57. Explain what happens in the bottleneck layer of an autoencoder. Why is this representation useful?

The **bottleneck layer** in an autoencoder is the **central (most compressed) layer** of the network. It has **fewer neurons** than the input, forcing the network to **compress the input into a low-dimensional representation**.

🔧 What happens:

- The encoder network transforms the input into this **compressed latent space**.
- The decoder then reconstructs the input from this compressed form.

- No labels are used — the model is trained to **minimize reconstruction loss** (e.g., MSE between input and output).
-

 **Why it's useful:**

- The bottleneck forces the network to **learn the most relevant features**, removing noise and redundancy.
- This representation is useful for:
 - **Dimensionality reduction**
 - **Data compression**
 - **Feature extraction**
 - **Anomaly detection**
 - **Pretraining for downstream tasks**

 Essentially, the bottleneck captures the **essence** of the data in a compact form.

58. Compare PCA and autoencoders for dimensionality reduction. Include differences in linearity and training.

Both PCA and autoencoders aim to reduce the dimensionality of data while preserving important structure — but they differ in **architecture, linearity, and learning process**.

Aspect	PCA	Autoencoder
Type	Linear method	Can be non-linear (with activations)
Learning	Analytical (eigen decomposition)	Trained using backpropagation
Features	Projects data onto orthogonal axes (principal components)	Learns nonlinear embeddings from data
Interpretability	Eigenvectors are interpretable	Latent features are often abstract
Training Data	Needs centered data; no training loop	Requires more data and epochs

Aspect	PCA	Autoencoder
Expressiveness	Limited to linear relationships	Can model complex nonlinear patterns

When to use:

- Use **PCA** for **quick, interpretable** linear dimensionality reduction.
- Use **autoencoders** when:
 - You need to capture **nonlinear structure**
 - You're working with large, complex datasets (e.g., images, audio)

59. What are common regularization methods used in CNNs to prevent overfitting?

Overfitting occurs when a CNN learns to **memorize training data**, failing to generalize to unseen examples. Regularization techniques help improve **generalization** by reducing model complexity or increasing robustness.

Common Regularization Methods in CNNs:

1. Dropout

- Randomly disables a fraction of neurons during training (e.g., 0.5), forcing the network to learn redundant representations.

2. L2 Regularization (Weight Decay)

- Adds a penalty proportional to the square of the weights to the loss function.
- Encourages smaller weights, reducing model complexity.

3. Data Augmentation

- Generates new training examples using transformations (rotation, flipping, cropping), exposing the model to more variation.

4. Batch Normalization

- Normalizes layer inputs to reduce internal covariate shift and stabilize learning, which can also have a regularizing effect.

5. Early Stopping

- Stops training when validation performance stops improving, preventing overfitting on training data.
-

📌 These methods are often used **together** for best results.

60. How does early stopping work, and what does it monitor during CNN training?

Early stopping is a regularization technique used to prevent overfitting by **halting training when the model stops improving** on a **validation set**.

🔍 How It Works:

1. Split the data into **training** and **validation** sets.
 2. During training, monitor a performance metric on the **validation set** (typically **validation loss** or **validation accuracy**).
 3. If the metric does not improve for a defined number of epochs (called **patience**), stop training.
 4. Optionally, **restore the weights** from the epoch with the best validation performance.
-

✅ Why It Helps:

- Prevents the model from learning noise in the training set.
- Ensures training **stops before overfitting begins**.
- Saves computation time by avoiding unnecessary epochs.

📌 Early stopping is widely supported in deep learning frameworks like **Keras**, **PyTorch**, and **TensorFlow**.

21. Generalization is the ability of a model to perform well on unseen data.
 22. The two components of generalization error are bias and variance.
 23. L1 regularization promotes model sparsity by eliminating features.
 24. Dropout helps reduce overfitting by randomly disabling neurons during training.
 25. Ridge regression is also known as L2 regularization.
27. Explain the difference between bias and variance in a machine learning model. How do they affect generalization?

Bias and **variance** are two fundamental sources of error in machine learning models that impact how well a model generalizes to unseen data.

◆ Bias

- Bias refers to the **error due to overly simplistic assumptions** in the model.
- It leads to **underfitting**, where the model fails to capture the true patterns in the data.
- High bias results in **high training and test error**.

📌 *Example:* A linear model trying to fit non-linear data.

◆ Variance

- Variance measures the model's **sensitivity to changes in the training data**.
- High variance leads to **overfitting**, where the model memorizes the training data but performs poorly on new data.
- High variance results in **low training error but high test error**.

📌 *Example:* A deep decision tree that changes drastically with small data changes.

⚖️ Bias-Variance Tradeoff

- The goal is to find a **balance**:
 - **Low bias** → accurate model
 - **Low variance** → stable, generalizable model
 - Too simple → high bias (underfit)
 - Too complex → high variance (overfit)
-

✓ **Generalization** is best achieved when both bias and variance are **reasonably low**, allowing the model to capture patterns without being overly sensitive to noise.

29. What is the purpose of the bias-variance tradeoff? How can a balance be achieved?

Purpose:

To achieve the **lowest possible generalization error** by finding the optimal balance between:

- **Bias** — error due to overly simple models that underfit the data.
 - **Variance** — error due to overly complex models that overfit the training data.
-

⚖️ **How to achieve balance:**

1. **Model selection:** Choose a model with the right complexity for your data (e.g., regularized linear model vs. deep neural network).
2. **Cross-validation:** Helps estimate generalization performance and detect overfitting/underfitting.
3. **Regularization** (e.g., L1/L2): Penalizes overly complex models to reduce variance.
4. **Training with more data:** Can reduce variance and stabilize the model.

30. How do L1 and L2 regularization differ in terms of their effect on feature selection and model complexity?

L1 Regularization (Lasso):

- Adds the **absolute value** of the weights to the loss function:

$$\text{Loss} + \lambda \sum |\theta_i|$$

- Encourages **sparsity**: many weights become exactly **zero**, effectively performing **feature selection**.
- Useful when you expect **only a few features to be relevant**.

L2 Regularization (Ridge):

- Adds the **squared value** of the weights:

$$\text{Loss} + \lambda \sum \theta_i^2$$

- Encourages **small weights**, but **rarely zero**. It **shrinks** all coefficients, reducing model **complexity**, but **keeps all features**.
- Useful when all features are expected to contribute a little.

31. Explain how Elastic Net regularization works. Why might it be preferred over using L1 or L2 alone?

Elastic Net combines both **L1 and L2** penalties in the loss function. It balances **sparsity** (from L1) and **stability** (from L2).

Why prefer Elastic Net?

- **L1 alone** may arbitrarily select one feature from a group of correlated features.
- **L2 alone** keeps all features, even irrelevant ones.
- **Elastic Net** handles **multicollinearity better**, selecting **groups of correlated features** and promoting generalization.
- It is especially useful when:
 - You have **many correlated features**
 - You expect **some feature selection** but also want **model stability**

32. What are three practical techniques (other than L1/L2) to prevent overfitting in deep learning?

?

Dropout

- Randomly deactivates a fraction of neurons during training (e.g., 20–50%).
- Prevents the network from becoming too reliant on specific neurons, promoting redundancy and robustness.
- At test time, all neurons are active but scaled appropriately.

?

Early Stopping

- Monitors performance on a **validation set** during training.
- Stops training when validation loss **starts to increase**, even if training loss is still decreasing.
- Prevents the model from memorizing noise in the training data.

?

Data Augmentation

- Artificially increases the size and diversity of the training dataset by applying transformations (e.g., rotation, flipping, cropping, color shift).
- Helps the model generalize better by seeing more varied examples without collecting new data.

33. Why is early stopping considered a form of regularization? How is it implemented?

Early stopping is considered a form of **regularization** because it helps prevent **overfitting** by controlling how long the model is allowed to learn. Instead of modifying the loss function like L1 or L2, it limits the model's **training time** to prevent it from fitting noise in the training data.

?

How It Works:

1. During training, the model's performance is evaluated on a **validation set** after each epoch.
 2. If the **validation loss stops improving** for a certain number of epochs (called **patience**), training is halted.
 3. The model parameters from the **best-performing epoch** are restored.
-

✓

Why it's regularization:

- It implicitly limits model complexity by **not allowing too many training iterations**.
 - Helps maintain **generalization** by stopping before the model starts memorizing the training set.
-

📌 **Implementation Tip:** Most deep learning frameworks (like Keras, PyTorch, TensorFlow) have built-in callbacks for early stopping.

34. What makes L1 regularization especially useful in high-dimensional datasets?

L1 regularization (Lasso) is particularly useful in high-dimensional datasets because it performs **automatic feature selection**.

It encourages **sparsity** in the model by driving many feature coefficients to exactly **zero**, effectively removing irrelevant or redundant features.

📌 This helps simplify the model, reduce overfitting, and improve interpretability — especially when the number of features is much larger than the number of samples.

35. Describe a scenario where increasing the number of training samples improves model performance.

When training a model on a **small dataset**, it may easily **overfit** — learning noise or specific patterns that don't generalize.

For example, in **image classification**, a CNN trained on only 500 images per class may perform poorly on unseen data.

✅ **Increasing the number of training samples** exposes the model to more **variation**, helping it learn **general patterns** rather than memorizing. As a result, both **training stability** and **test accuracy** improve.

📌 This is especially effective for **high-capacity models** (e.g., deep neural networks) that require a large amount of data to avoid overfitting.

36. Explain how regularization helps in achieving better generalization, especially when training data is limited.

When training data is limited, models can easily **overfit**, memorizing the noise rather than learning the underlying patterns.

Regularization helps prevent this by **constraining the model's complexity**, making it less sensitive to fluctuations in the training data.

Examples:

- **L1/L2 regularization** penalizes large weights, promoting simpler models.
- **Dropout** prevents co-adaptation by randomly deactivating neurons.
- **Early stopping** halts training before the model starts overfitting.

These techniques **reduce variance** and help the model generalize better, even when **data is scarce**.

ResNet, Transfer Learning

Backpropagation is a method for optimizing the weights in a neural network.

It works by **applying the chain rule** to compute the **gradient of the loss function with respect to each weight**, starting from the output layer and moving backward through the network.

These gradients are then used to **adjust the weights in the direction that minimizes the loss**.

Transfer learning is a technique where a model trained on a large source dataset (e.g., ImageNet) is reused for a different but related target task.

Instead of training a new model from scratch, we use the **generic features** learned in the early layers of the pre-trained model and **freeze them**.

Then we **replace and train the final layers** (e.g., classification head) using data from the new task.

This approach is especially useful when the target dataset is small.

Fine-tuning is an extension of transfer learning.

After using a pre-trained model, we **unfreeze some of its layers** (usually the higher ones) and **continue training them** on the new task.

This allows the model to adjust both generic and task-specific features to better fit the new data.

Fine-tuning requires more data and computing but often results in higher accuracy.

ResNet (Residual Network) is a deep neural network architecture that introduces **skip connections** to make training of very deep networks more stable and effective.

Instead of learning the full transformation $H(x)$, each block learns a **residual function** $F(x)$, which represents the difference between the desired output and the input.

The final output is computed as:

$$H(x) = F(x) + x$$

This residual connection allows the network to **preserve useful information from earlier layers**, improves **gradient flow**, and helps avoid problems like **vanishing gradients**.

As a result, ResNet can be trained with **many layers (e.g., 50, 101, 152)** while maintaining or even improving accuracy.