

**Question ID:**	M01-Q01
**Topic:**	M1a – What is Data Science?
**Difficulty:**	Easy
**Blooms:**	Understand

## What is Data Science mainly about?

1. Using data to identify patterns and make informed decisions
2. Creating as much data as possible
3. Replacing scientists with AI
4. Making charts without analyzing meaning
5. Searching the internet for interesting pictures

**Question ID:**	M01-Q02
**Topic:**	M1a – Flood of Data
**Difficulty:**	Easy
**Blooms:**	Understand

## Why has Data Science become more important recently?

1. Because we now collect and store enormous amounts of data
2. Because statistics were just invented recently
3. Because computers are disappearing
4. Because hand-drawn measurements are more popular
5. Because people have shorter attention spans

**Question ID:**	M01-Q03
**Topic:**	M1a – Data Visualization
**Difficulty:**	Medium
**Blooms:**	Apply

## What role does visualization play in Data Science?

1. Helps humans understand patterns that exist in large datasets
2. It eliminates the need for data cleaning
3. It guarantees conclusions are correct
4. It replaces statistics entirely
5. It is only useful for small datasets

**Question ID:**	M01-Q04
**Topic:**	M1a – Statistics: Nonsense Protection
**Difficulty:**	Medium
**Blooms:**	Analyze

## What does statistics help prevent in analysis?

1. False conclusions caused by random patterns
2. Collecting too much data
3. Good visualizations
4. The need for computers
5. Faster machine learning training

**Question ID:**	M01-Q05
**Topic:**	M1a – Coping with Data
**Difficulty:**	Medium
**Blooms:**	Remember

## Which is NOT one of the core parts of data science discussed?

1. Handling data scale
2. Visualization
3. Algorithms
4. Astrology
5. Computation

**Question ID:**	M01-Q06
**Topic:**	M1b – Tools
**Difficulty:**	Easy
**Blooms:**	Understand

## Why is Python widely used in Data Science?

1. It has libraries for math, visualization, and machine learning
2. It is the only programming language for data analysis
3. It was designed to replace image editors
4. It is extremely fast for everything
5. It can only be used in Jupyter notebooks

**Question ID:**	M01-Q07
**Topic:**	M1b – Google Colab
**Difficulty:**	Easy
**Blooms:**	Remember

**Which tool provides a Python environment in the browser without local installation?**

1. Anaconda
2. Google Colab
3. Visual Basic
4. GitHub Desktop
5. Local terminal only

**Question ID:**	M01-Q08
**Topic:**	M1c – Arrays
**Difficulty:**	Medium
**Blooms:**	Understand

**Which statement about arrays is TRUE?**

1. Easy to insert in the middle
2. Always sorted
3. Fast access by index
4. Automatically distributed across different computers
5. Removing an element is always O(1)

**Question ID:**	M01-Q09
**Topic:**	M1c – Stacks
**Difficulty:**	Medium
**Blooms:**	Analyze

**Which is a disadvantage of stacks?**

1. Hard to access arbitrary elements
2. Easy to remove from the top
3. Very slow at resizing
4. They don't allow push operations
5. They cannot be implemented in Python

**Question ID:**	M01-Q10
**Topic:**	M1c – Queue vs Stack
**Difficulty:**	Easy
**Blooms:**	Understand

## What is the key difference between a Queue and a Stack?

1. Queue removes the earliest inserted item first; Stack removes the most recently inserted
2. Queue is stored in trees
3. Stack elements must be unique
4. Queue requires a GPU
5. Stack never allows push operations

**Question ID:**	M01-Q11
**Topic:**	M1c – Linked Lists
**Difficulty:**	Medium
**Blooms:**	Apply

## Linked Lists are particularly good when...

1. We need constant-time random indexing
2. Frequent insertions or deletions at ends or arbitrary positions
3. Data never changes
4. We want built-in sorting
5. We only store one element

**Question ID:**	M01-Q12
**Topic:**	M1c – Hash Tables
**Difficulty:**	Medium
**Blooms:**	Understand

## A key-value store such as a Python dictionary...

1. Allows fast lookup by key
2. Preserves items in sorted order
3. Can only store numbers
4. Requires sequential search for access
5. Only allows string keys

**Question ID:**	M01-Q13
**Topic:**	M1c – Trees
**Difficulty:**	Medium
**Blooms:**	Understand

## What is a Tree most useful for compared to linear structures?

1. Representing hierarchical relationships
2. Guaranteeing O(1) for all operations
3. Storing only numbers
4. Using no pointers internally
5. Always being binary

**Question ID:**	M01-Q14
**Topic:**	M1c – Broadcasting
**Difficulty:**	Medium
**Blooms:**	Apply

## In NumPy, what is “broadcasting”?

1. Automatically expanding arrays to compatible shapes during arithmetic
2. Sending arrays to a TV channel
3. Only used when adding identical shapes
4. A memory compression method
5. Copying data into SQL databases

**Question ID:**	M01-Q15
**Topic:**	M1c – Shape compatibility
**Difficulty:**	Medium
**Blooms:**	Analyze

## Why can't you always add two NumPy arrays of different shapes?

1. The shapes must be compatible so elements correspond
2. NumPy refuses to add integers
3. Arrays must always be  $2 \times 2$
4. The larger array deletes itself
5. Only works for prime-number shapes

**Question ID:**	M01-Q16
**Topic:**	M1c — Array memory organization
**Difficulty:**	Medium
**Blooms:**	Understand

**One major reason arrays allow fast random access is:**

1. The array keeps a hash table of all indexes
2. Each element stores the address of the next
3. Elements are stored contiguously so index lookup is constant time
4. The OS automatically accelerates access for arrays
5. Arrays are always stored in CPU cache

**Question ID:**	M01-Q17
**Topic:**	M1c — Array disadvantages
**Difficulty:**	Medium
**Blooms:**	Analyze

**Why is deleting an element from the middle of a Python list inefficient?**

1. All later elements must shift to fill the gap
2. Python must rebuild the entire interpreter
3. Removal causes permanent fragmentation
4. Python sorts the list after every deletion
5. The list converts to a linked list internally

**Question ID:**	M01-Q18
**Topic:**	M1c — FILO semantics
**Difficulty:**	Medium
**Blooms:**	Apply

**Which scenario best leverages a stack data structure?**

1. Tracking nested function calls
2. Processing real-time streaming logs
3. Storing only unique items
4. Implementing alphabetical search
5. Indexing stock-market time series

**Question ID:**	M01-Q19
**Topic:**	M1c — FIFO
**Difficulty:**	Easy
**Blooms:**	Understand

## A queue supports which behavior?

1. Requires O(1) random access
2. Only works if all items are numbers
3. Removes the earliest inserted item first
4. Can only grow, never shrink
5. Automatically sorts items by priority

**Question ID:**	M01-Q20
**Topic:**	M1c — Linked lists
**Difficulty:**	Medium
**Blooms:**	Analyze

## Linked lists scale well for which workload?

1. Fast binary search
2. Frequent insertions at arbitrary locations
3. Automatic rebalancing of data
4. Multi-dimensional indexing
5. Constant-time random access

**Question ID:**	M01-Q21
**Topic:**	M1c — Hash functions
**Difficulty:**	Medium
**Blooms:**	Understand

## Why are dictionaries (hash tables) efficient for lookup by key?

1. Keys are stored in sorted order
2. A hash function maps keys to near-constant-time index access
3. Values are duplicated in multiple locations
4. Every operation scans the entire table
5. The data is stored on the GPU

**Question ID:**	M01-Q22
**Topic:**	M1c — Trees
**Difficulty:**	Medium
**Blooms:**	Understand

**A tree differs from a stack or queue primarily because it:**

1. Represents hierarchical relationships
2. Stores only numbers
3. Guarantees all operations are O(1)
4. Uses no pointers internally
5. Is always binary

**Question ID:**	M01-Q23
**Topic:**	M1c — Composite
**Difficulty:**	Hard
**Blooms:**	Create

**Which situation is best for a composite data structure (e.g., tree of arrays)?**

1. You only need sequential iteration
2. You must store data that never changes
3. You need to mix hierarchical lookup with fast local indexing
4. You want to minimize implementation complexity
5. You require strict alphabetical ordering

**Question ID:**	M01-Q24
**Topic:**	M1c — Hash table disadvantages
**Difficulty:**	Medium
**Blooms:**	Analyze

**Which is a key drawback of hash tables?**

1. They require keys to be sorted
2. They cannot delete items
3. They do not preserve meaningful order
4. Searching is always slower than a linked list
5. They can only store fixed-size elements

**Question ID:**	M01-Q25
**Topic:**	M1c — Broadcasting rules
**Difficulty:**	Medium
**Blooms:**	Apply

## In NumPy, why does adding a $1 \times 3$ vector to a $3 \times 3$ matrix work?

1. The vector is broadcast across rows to match shape
2. The matrix is flattened automatically
3. NumPy guesses the user's intention
4. The vector overwrites the diagonal
5. Because 3 is a special broadcast-safe number

**Question ID:**	M01-Q26
**Topic:**	M1c — Shape mismatch
**Difficulty:**	Medium
**Blooms:**	Analyze

## Which NumPy operation fails without broadcasting compatibility?

1. Adding arrays with mismatched dimensions that cannot expand
2. Adding a scalar to any array
3. Element-wise multiplication of equal shapes
4. Index slicing
5. Dot product of same-length vectors

**Question ID:**	M01-Q27
**Topic:**	M1b — Tools justification
**Difficulty:**	Easy
**Blooms:**	Understand

## Why is Python frequently chosen for data science workflows?

1. The language forces all variables to be floats
2. Python code always runs faster than C
3. Strong ecosystem of numerical and ML libraries
4. It only works in notebook environments
5. It replaces the need for distributed computing

**Question ID:**	M01-Q28
**Topic:**	M1a — Human cognition limits
**Difficulty:**	Medium
**Blooms:**	Evaluate

## Why is data visualization essential in data science?

1. Visualizations remove all noise from data
2. Humans understand visual patterns better than raw numbers
3. Plots guarantee a correct conclusion
4. Visualizations replace statistical inference
5. Visualization is only useful on labeled data

**Question ID:**	M01-Q29
**Topic:**	M1a — Big data scale challenges
**Difficulty:**	Medium
**Blooms:**	Apply

## Handling big data requires more than spreadsheets because:

1. Distributed tools allow scalable computation
2. Spreadsheets always corrupt files over 10 MB
3. Databases compute without hardware
4. Large datasets require no cleaning
5. CPUs cannot process tabular data

**Question ID:**	M01-Q30
**Topic:**	M1a — Cleaning/processing effort
**Difficulty:**	Medium
**Blooms:**	Analyze

## Why will data engineering remain a major effort in DS pipelines?

1. Real-world data is messy and must be cleaned before modeling
2. Machine learning automatically repairs all errors
3. Storage formats determine the model accuracy
4. Sensor data arrives perfectly structured
5. Data rarely changes once collected

**Question ID:**	M03-Q1
**Topic:**	Exploratory Data Analysis
**Difficulty:**	Easy
**Blooms:**	Understand

**Which part of the data science workflow primarily focuses on understanding the structure, patterns, and anomalies present in the data?**

1. Data Collection
2. **Exploratory Data Analysis (EDA)**
3. Confirmatory Data Analysis
4. Feature Deployment
5. GPU Optimization

**Question ID:**	M03-Q2
**Topic:**	Histogram
**Difficulty:**	Easy
**Blooms:**	Remember

**Which visualization is most appropriate for examining the distribution of a single quantitative variable?**

1. Scatter plot
2. **Histogram**
3. Network graph
4. Choropleth map
5. Box-and-whisker + KDE overlay

**Question ID:**	M03-Q3
**Topic:**	Boxplot
**Difficulty:**	Easy
**Blooms:**	Understand

**Box-and-whisker plots are especially useful for:**

1. Showing data on a map
2. **Comparing medians and detecting outliers**
3. Displaying 3D surfaces
4. Tracking time series index returns
5. Showing relationships between nominal variables

**Question ID:**	M03-Q4
**Topic:**	Nominal vs Nominal
**Difficulty:**	Medium
**Blooms:**	Apply

**A bar chart comparing 'Plant Type' and 'Fruit Variety' represents what type of data?**

1. Ordinal vs Quantitative
2. Nominal vs Nominal
3. Ordinal vs Ordinal
4. Quantitative vs Quantitative
5. Geospatial vs Quantitative

**Question ID:**	M03-Q5
**Topic:**	Scatter Plot
**Difficulty:**	Easy
**Blooms:**	Remember

**Scatter plots are primarily used to visualize:**

1. Relationships between two quantitative variables
2. Distribution of a single variable only
3. Hierarchical relationships
4. Statistical inference and p-values
5. Survey proportions

**Question ID:**	M03-Q6
**Topic:**	Time Series
**Difficulty:**	Medium
**Blooms:**	Apply

**Which visualization best compares multiple time-dependent quantities simultaneously?**

1. Stacked area chart
2. Single-variable histogram
3. Parallel coordinates
4. Radar chart
5. Candlestick chart

**Question ID:**	M03-Q7
**Topic:**	Candlestick Chart
**Difficulty:**	Easy
**Blooms:**	Remember

**Which graph type is explicitly linked to financial market data in the lecture?**

1. Horizon graph
2. Candlestick chart
3. Contour map
4. Boxplot
5. Tree map

**Question ID:**	M03-Q8
**Topic:**	Matplotlib API
**Difficulty:**	Medium
**Blooms:**	Understand

**In matplotlib, why is using ‘fig, ax = plt.subplots()‘ preferred over calling ‘plt.plot()‘ directly?**

1. It uses more memory so it's faster
2. It gives explicit references to figure and axes objects for better control
3. It enables automatic machine learning integration
4. It prevents adding labels and legends
5. It is required by NumPy

**Question ID:**	M03-Q9
**Topic:**	KDE
**Difficulty:**	Medium
**Blooms:**	Analyze

**A KDE (Kernel Density Estimate) is used to:**

1. Visualize category labels
2. Smooth the distribution of sampled data
3. Display geospatial movement
4. Simulate stock trading
5. Normalize missing values

**Question ID:**	M03-Q10
**Topic:**	Outliers
**Difficulty:**	Medium
**Blooms:**	Evaluate

**The lecture suggests that before trusting an outlier, you should:**

1. Remove all outliers automatically
2. Trace back to the original data and verify context
3. Replace with the mean
4. Convert to categorical encoding
5. Report it as a major scientific discovery

**Question ID:**	M03-Q11
**Topic:**	High-dimensional visualization
**Difficulty:**	Medium
**Blooms:**	Analyze

**Parallel coordinate plots are best suited for:**

1. Two quantitative variables
2. Nominal-only comparisons
3. High-dimensional numeric data
4. Exact probability estimation
5. Animated game graphics

**Question ID:**	M03-Q12
**Topic:**	Q-Q Plot
**Difficulty:**	Hard
**Blooms:**	Analyze

**Which visualization is used to examine whether data follows a theoretical distribution?**

1. Q-Q plot
2. Box plot
3. Stacked bar chart
4. Flow map
5. 3D volume rendering

**Question ID:**	M03-Q13
**Topic:**	Composition over Time
**Difficulty:**	Easy
**Blooms:**	Remember

**Which visualization technique was illustrated using time-series of Panda Hats and Underpants?**

1. Index chart
2. **Stacked time-series visualization**
3. Heat map
4. Network diagram
5. Histogram with multiple bins

**Question ID:**	M03-Q14
**Topic:**	Quantitative Relationship
**Difficulty:**	Medium
**Blooms:**	Remember

**What is added to the mother's height in calculating Galton's 'midparent' height?**

1.  $0.98 \times$  mother height
2.  **$1.08 \times$  mother height**
3.  $1.50 \times$  mother height
4. Subtract father height
5. It uses only father height

**Question ID:**	M03-Q15
**Topic:**	3-variable visualization
**Difficulty:**	Medium
**Blooms:**	Apply

**Which chart type allows three variables to be shown simultaneously using two axes and color/size encoding?**

1. Bar chart with sublevels
2. **Scatter plot with a third encoding**
3. Pie chart slices only
4. Table layout
5. Q-Q plot

**Question ID:**	M03-Q16
**Topic:**	Summary Statistics
**Difficulty:**	Easy
**Blooms:**	Remember

Which of the following appears in the slides' list of summary statistics for EDA?

1. Skewness
2. Standard deviation
3. Fourier coefficients
4. Mode only
5. Z-score thresholds

**Question ID:**	M03-Q17
**Topic:**	Missing Values
**Difficulty:**	Easy
**Blooms:**	Remember

Which package was highlighted for visualizing patterns of missing data?

1. missingno
2. seaborn.gridplot
3. plotly.missmap
4. statsmodels.na\_viz
5. opencv.impute

**Question ID:**	M03-Q18
**Topic:**	Pairwise Relationships
**Difficulty:**	Medium
**Blooms:**	Understand

A matrix of scatter plots helps you:

1. Visualize pairwise relationships among many quantitative variables
2. Render 3D surfaces of a function
3. Encode hierarchical trees
4. Compute p-values for regression
5. Normalize geospatial coordinates

**Question ID:**	M03-Q19
**Topic:**	Correlation Heat Map
**Difficulty:**	Easy
**Blooms:**	Understand

## What does a correlation heat map display?

1. The sign and magnitude of linear relationships between variables
2. Raw counts of categories
3. Geographic elevation
4. Financial OHLC patterns
5. Kernel bandwidths

**Question ID:**	M03-Q20
**Topic:**	Co-occurrence
**Difficulty:**	Medium
**Blooms:**	Apply

Pair-wise co-occurrence (joint probabilities) in EDA are most appropriate for:

1. Relationships between categorical variables
2. Optimizing hyperparameters
3. 3D surface estimation
4. Fourier spectral analysis
5. GPU memory profiling

**Question ID:**	M03-Q21
**Topic:**	Automated EDA Tools
**Difficulty:**	Easy
**Blooms:**	Remember

Which tool is NOT listed among the automated EDA tools in the slides?

1. Pandas Profiling
2. Sweetviz
3. Autoviz
4. D-Tale
5. PowerBI AutoInspect

**Question ID:**	M03-Q22
**Topic:**	Index Charts
**Difficulty:**	Medium
**Blooms:**	Understand

**In an index chart for time series, what does 'indexing' typically do?**

1. Normalizes each series to a common baseline (e.g., 100) to compare relative change
2. Sorts categories alphabetically
3. Automatically removes outliers
4. Interpolates missing geolocations
5. Computes PCA scores

**Question ID:**	M03-Q23
**Topic:**	Small Multiples
**Difficulty:**	Medium
**Blooms:**	Apply

**Small multiples are especially helpful when:**

1. Comparing trends across many categories using the same axes and scales
2. Rendering a single 3D surface
3. Computing z-scores
4. Encoding packet network flows
5. Estimating kernel bandwidth

**Question ID:**	M03-Q24
**Topic:**	Horizon Graph
**Difficulty:**	Hard
**Blooms:**	Analyze

**Which statement best describes a horizon graph?**

1. A layered time-series display that folds bands of values to save vertical space
2. A 3D bar chart for horizons and altitudes
3. A map projection for polar regions
4. A network centrality diagram
5. A chord diagram for gene expression

**Question ID:**	M03-Q25
**Topic:**	Radar Chart
**Difficulty:**	Medium
**Blooms:**	Apply

Radar charts are typically appropriate when:

1. Variables are positive and you want to compare multivariate profiles
2. You need to visualize negative-only values
3. You must encode geographic directions
4. You want to estimate probability density
5. You need exact correlation coefficients

**Question ID:**	M03-Q26
**Topic:**	Maps
**Difficulty:**	Easy
**Blooms:**	Remember

Which of the following is NOT listed among the map types in the lecture?

1. Flow map
2. Graduated symbol map
3. Choropleth
4. Cartogram
5. Sankey map

**Question ID:**	M03-Q27
**Topic:**	Matplotlib Save
**Difficulty:**	Easy
**Blooms:**	Remember

Which command saves a matplotlib figure to PDF as shown in the slides?

1. plt.savefig('figure.pdf')
2. fig.savefig('figure.pdf')
3. ax.savefig('figure.pdf')
4. plt.writepdf(fig)
5. np.savetxt('figure.pdf')

**Question ID:**	M03-Q28
**Topic:**	Matplotlib Styling
**Difficulty:**	Easy
**Blooms:**	Remember

**In the matplotlib example, which keyword argument sets the line's color?**

1. linewidth
2. color
3. alpha
4. style
5. markerface

**Question ID:**	M03-Q29
**Topic:**	Matplotlib API Style
**Difficulty:**	Medium
**Blooms:**	Understand

**Which coding style does the lecture label as 'Lazy (try to avoid)'?**

1. Using the stateful pyplot interface without figure/axes objects
2. Using the object-oriented fig/ax API
3. Saving figures to files
4. Calling NumPy linspace
5. Adding labels and titles

**Question ID:**	M03-Q30
**Topic:**	Seaborn Facets
**Difficulty:**	Medium
**Blooms:**	Apply

**For faceted comparisons across categories, which approach is highlighted in the visualization lecture?**

1. Seaborn's FacetGrid
2. Matplotlib's imshow
3. D3's force simulation
4. NetworkX spring layout
5. OpenCV's cvtColor

**Question ID:**	M06-Q1
**Topic:**	ML Process: Problem Framing
**Difficulty:**	Medium
**Blooms:**	Analyze

**Which step most directly ensures that the problem you're solving is actually a machine-learning problem and not a simple rules or query task?**

1. Check whether a labeled target exists and if patterns must generalize to new data
2. Visualize predictions with a confusion matrix
3. Collect more data first
4. Tune hyperparameters with cross-validation
5. Train a baseline linear model

**Question ID:**	M06-Q2
**Topic:**	ML Process: Data Splits
**Difficulty:**	Easy
**Blooms:**	Understand

**In a standard ML workflow, a hold-out test set is primarily used to...**

1. Select the best model during tuning
2. Fit the feature scaler and imputer
3. Estimate the final generalization performance after all modeling choices are frozen
4. Balance class labels in the training data
5. Visualize learning curves

**Question ID:**	M06-Q3
**Topic:**	ML Process: Leakage Prevention
**Difficulty:**	Medium
**Blooms:**	Apply

**Which practice best prevents \*\*data leakage\*\* when scaling features?**

1. Fit the scaler only on the training data, then apply the fitted transform to validation/test
2. Use MinMax scaling instead of standardization
3. Use more features so leakage averages out
4. Shuffle the rows before scaling
5. Fit the scaler on all available data, then transform splits

**Question ID:**	M06-Q4
**Topic:**	ML Process: Baselines
**Difficulty:**	Easy
**Blooms:**	Understand

**A \*\*baseline\*\* model in the ML process is best described as...**

1. A model with zero variance predictions
2. The final, most complex model
3. **A simple, often naive model used to set a minimum performance bar**
4. Any linear model
5. A model trained without regularization

**Question ID:**	M06-Q5
**Topic:**	ML Process: Feature Engineering
**Difficulty:**	Hard
**Blooms:**	Evaluate

**During feature engineering, creating target-aware features (e.g., averaging the target by category using all rows) mainly risks...**

1. Worse calibration but correct ranking
2. High bias
3. **Data leakage that inflates validation scores**
4. Underfitting
5. Improved interpretability with no downside

**Question ID:**	M06-Q6
**Topic:**	ML Process: Pipelines
**Difficulty:**	Medium
**Blooms:**	Analyze

**Which statement about \*\*pipelines\*\* is most accurate?**

1. Pipelines only chain models, not transforms
2. Pipelines are slower but identical to manual code
3. **Pipelines ensure that cross-validation folds fit transforms (impute/scale) using only the training fold each time**
4. Pipelines prevent overfitting by adding noise
5. Pipelines eliminate the need for a test set

**Question ID:**	M06-Q7
**Topic:**	ML Process: Splits for Time Series
**Difficulty:**	Medium
**Blooms:**	Apply

**Which split strategy is most appropriate for \*\*time-series\*\* forecasting?**

1. Random k-fold cross-validation
2. Bootstrap resampling only
3. Stratified shuffle split
4. Leave-one-out cross-validation
5. **Forward-chaining (expanding window) validation**

<b>Question ID:</b>	M06-Q8
<b>Topic:</b>	ML Process: Lifecycle
<b>Difficulty:</b>	Easy
<b>Blooms:</b>	Remember

**In CRISP-DM-like lifecycles, the step after 'Modeling' that checks fitness-for-use against stakeholder goals is...**

1. Data understanding
2. **Evaluation**
3. Business understanding
4. Deployment
5. Data preparation

<b>Question ID:</b>	M06-Q9
<b>Topic:</b>	ML Process: Monitoring
<b>Difficulty:</b>	Medium
<b>Blooms:</b>	Analyze

**Which warning most strongly indicates \*\*concept drift\*\* after deployment?**

1. Feature distributions in production shift relative to training and error rises on recent labeled samples
2. Stable calibration curves
3. Steady validation accuracy
4. Lower variance in predictions
5. Slightly higher training loss

<b>Question ID:</b>	M06-Q10
<b>Topic:</b>	ML Process: Metric Fit
<b>Difficulty:</b>	Medium
<b>Blooms:</b>	Analyze

**A confusion matrix is \*\*not\*\* the right tool to evaluate a model when...**

1. You need class-wise errors
2. You want false positive rate
3. **You must compare ranking quality at varying thresholds**
4. You want recall per class
5. You are solving classification

<b>Question ID:</b>	M06-Q11
<b>Topic:</b>	Regression: OLS Objective
<b>Difficulty:</b>	Easy
<b>Blooms:</b>	Remember

**In simple linear regression, the \*\*least squares\*\* estimator chooses coefficients that...**

1. Equalize residuals across x
2. Minimize mean absolute error of residuals
3. Maximize R<sup>2</sup> directly
4. Minimize classification error
5. **Minimize the sum of squared residuals**

<b>Question ID:</b>	M06-Q12
<b>Topic:</b>	Regression: Assumptions
<b>Difficulty:</b>	Medium
<b>Blooms:</b>	Evaluate

**Which situation most clearly violates a key linear regression assumption?**

1. n < p
2. Residuals are approximately normal
3. **Residual variance increases with x (fan-shaped residual plot)**
4. There are categorical predictors encoded with dummies
5. Predictors are scaled to zero-mean

<b>Question ID:</b>	M06-Q13
<b>Topic:</b>	Regression: Multicollinearity
<b>Difficulty:</b>	Medium
<b>Blooms:</b>	Analyze

Multicollinearity primarily affects which aspect of a linear model?

1. Feasibility of predictions
2. Training MSE only
3. Interpretation of intercept only
4. Variance/stability of coefficient estimates
5. Unbiasedness of OLS coefficients

**Question ID:**	M06-Q14
**Topic:**	Regression: Overfitting
**Difficulty:**	Easy
**Blooms:**	Understand

You add a perfectly predictive but noisy feature. Training MSE drops sharply; validation MSE rises. The model most likely...

1. Underfit
2. Overfit due to high variance
3. Suffers from label leakage
4. Is unbiased
5. Has perfect calibration

**Question ID:**	M06-Q15
**Topic:**	Regression: Model Comparison
**Difficulty:**	Medium
**Blooms:**	Analyze

Compared to \*\*R<sup>2</sup>\*\*, \*\*Adjusted R<sup>2</sup>\*\* is preferred for model comparison because it...

1. Is invariant to scaling of y
2. Penalizes added predictors that don't improve fit enough
3. Always increases when you add predictors
4. Is threshold-independent
5. Equals correlation squared between y and  $\hat{y}$  always

**Question ID:**	M06-Q16
**Topic:**	Regression: Ridge
**Difficulty:**	Medium
**Blooms:**	Understand

## Ridge regression differs from OLS by...

1. Guaranteeing sparsity
2. Adding an L1 penalty on coefficients
3. Optimizing MAE instead of MSE
4. Adding interaction terms automatically
5. Adding an L2 penalty that shrinks coefficients toward zero

**Question ID:**	M06-Q17
**Topic:**	Regression: Lasso
**Difficulty:**	Medium
**Blooms:**	Analyze

## The \*\*Lasso\*\* is especially useful when...

1. You only have categorical variables
2. All predictors are essential
3. You need unbiased estimates regardless of  $p \gg n$
4. You want feature selection via many coefficients exactly zero
5. You need grouped shrinkage

**Question ID:**	M06-Q18
**Topic:**	Regression: Model Complexity
**Difficulty:**	Easy
**Blooms:**	Understand

## Polynomial regression of degree 10 on small n most likely increases...

1. Sample size
2. Linearity
3. Bias
4. Variance
5. Homoskedasticity

**Question ID:**	M06-Q19
**Topic:**	Regression: Preprocessing
**Difficulty:**	Easy
**Blooms:**	Apply

**Centering and scaling predictors before regularized regression mainly...**

1. Ensures the penalty treats coefficients comparably across features
2. Makes intercept exactly zero
3. Changes predictions drastically
4. Improves RMSE regardless of data
5. Eliminates multicollinearity

**Question ID:**	M06-Q20
**Topic:**	Regression: Dummy Variables
**Difficulty:**	Medium
**Blooms:**	Apply

**In multiple regression with a categorical variable of k levels, correct dummy encoding requires...**

1. Dropping the intercept and using k-1 dummies and one interaction
2. k-1 dummy columns with an intercept (reference level)
3. Only one dummy column regardless of k
4. k dummy columns plus intercept
5. Target encoding by default

**Question ID:**	M06-Q21
**Topic:**	Evaluation: Cross-Validation
**Difficulty:**	Easy
**Blooms:**	Understand

**Which statement about \*\*k-fold cross-validation\*\* is correct?**

1. It uses the test set k times
2. It provides an estimate of generalization by repeatedly training on k-1 folds and validating on the remaining fold
3. It cannot be used with pipelines
4. Stratification is only for regression
5. It eliminates the need for a final test set

**Question ID:**	M06-Q22
**Topic:**	Evaluation: Metrics for Imbalance
**Difficulty:**	Medium
**Blooms:**	Analyze

For \*\*imbalanced\*\* binary classification, which metric is most informative across thresholds?

1. Accuracy
2. Precision-Recall (PR) curve / Average Precision
3. ROC AUC only
4. Explained variance
5.  $R^2$

**Question ID:**	M06-Q23
**Topic:**	Evaluation: Metric Matching
**Difficulty:**	Easy
**Blooms:**	Remember

Which pair is correctly matched?

1. Calibration  $\leftrightarrow$  Adjusted  $R^2$
2. Ranking  $\leftrightarrow$  ROC AUC
3. Classification  $\leftrightarrow$  RMSE
4. Clustering  $\leftrightarrow$  MAE
5. Regression  $\leftrightarrow$  F1-score

**Question ID:**	M06-Q24
**Topic:**	Evaluation: Confusion Matrix
**Difficulty:**	Easy
**Blooms:**	Remember

Given a confusion matrix, \*\*recall\*\* for the positive class equals...

1.  $FP / (FP + TN)$
2.  $TP / (TP + FP)$
3.  $TP / (TP + TN)$
4.  $TN / (TN + FP)$
5. TP / (TP + FN)

**Question ID:**	M06-Q25
**Topic:**	Evaluation: Calibration
**Difficulty:**	Medium
**Blooms:**	Apply

**When calibration matters (e.g., risk estimation), which is best practice?**

1. Optimize only for accuracy
2. Replace the loss with hinge loss
3. Threshold at 0.5 for all datasets
4. Use isotonic or Platt scaling on a validation set and check calibration curves
5. Prefer hard labels to probabilities

**Question ID:**	M06-Q26
**Topic:**	Evaluation: Regression Metrics
**Difficulty:**	Easy
**Blooms:**	Understand

**Which metric is \*\*scale-sensitive\*\* and penalizes larger errors more strongly?**

1.  $R^2$
2. MAE
3. F1-score
4. Accuracy
5. RMSE

**Question ID:**	M06-Q27
**Topic:**	Evaluation: Threshold Metrics
**Difficulty:**	Easy
**Blooms:**	Understand

**You care equally about precision and recall at a single operating point. Which metric best summarizes this?**

1. Matthews correlation
2. Average precision
3. Balanced accuracy
4. F1-score
5. ROC AUC

**Question ID:**	M06-Q28
**Topic:**	Evaluation: Nested CV
**Difficulty:**	Hard
**Blooms:**	Analyze

During hyperparameter tuning, **\*\*nested cross-validation\*\*** is used to...

1. Speed up grid search
2. Balance classes automatically
3. Provide an outer loop for unbiased performance estimation while inner CV selects hyperparameters
4. Reduce data leakage from scaling
5. Avoid creating a test set

<b>**Question ID:**</b>	M06-Q29
<b>**Topic:**</b>	Evaluation: Averaging
<b>**Difficulty:**</b>	Medium
<b>**Blooms:**</b>	Analyze

Which scenario calls for **\*\*macro-averaged\*\* F1 instead of micro-averaged F1?**

1. You want to give equal weight to each class regardless of frequency
2. You're evaluating regression models
3. You want to weight larger classes more
4. Binary classification only
5. Classes are balanced and you want an overall rate

<b>**Question ID:**</b>	M06-Q30
<b>**Topic:**</b>	Evaluation: ROC Interpretation
<b>**Difficulty:**</b>	Medium
<b>**Blooms:**</b>	Evaluate

An ROC curve that lies **\*\*below\*\*** the diagonal indicates...

1. Random performance
2. High calibration error only
3. A perfect classifier
4. Systematic reversal of labels; flipping the score sign would perform above chance
5. Data leakage

<b>**Question ID:**</b>	M07-Q1
<b>**Topic:**</b>	Classification overview
<b>**Difficulty:**</b>	Easy
<b>**Blooms:**</b>	Understand

**In supervised classification, what is the primary goal when some rows have known labels?**

1. To delete rows without labels so training is cleaner
2. To reduce dimensionality for visualization only
3. To generate synthetic labels for all rows from noise
4. **To label previously unlabeled rows using a learned model**
5. To cluster the data into groups without labels

**Question ID:**	M07-Q2
**Topic:**	Applications of classification
**Difficulty:**	Easy
**Blooms:**	Remember

**Which example task is most clearly a multi-class classification problem?**

1. Predicting house prices in dollars
2. Grouping unlabeled news articles by topic
3. Estimating the mean of a Gaussian distribution
4. Finding a low-dimensional embedding of images
5. **Assigning each handwritten digit image a class from 0–9**

**Question ID:**	M07-Q3
**Topic:**	Usual steps / preprocessing
**Difficulty:**	Medium
**Blooms:**	Analyze

**Which sequence best reflects the 'Usual Steps' emphasized for building a classifier in these slides?**

1. **Preprocessing & data cleaning → Feature finding/selection → Train & evaluate**
2. Hyperparameter search only → Deploy
3. Visualization only → Deploy model
4. Train model → Collect data → Deploy → Clean data
5. Feature engineering → Ignore preprocessing → Train

**Question ID:**	M07-Q4
**Topic:**	KNN and effect of k
**Difficulty:**	Medium
**Blooms:**	Understand

**For k-Nearest Neighbors (kNN), increasing k typically has what effect on the decision boundary?**

1. Has no effect on the boundary; only runtime changes
2. Converts the classifier into a linear separator
3. Smooths the boundary by averaging over more neighbors
4. Forces the classifier to overfit training data
5. Makes the boundary more jagged and sensitive to noise

**Question ID:**	M07-Q5
**Topic:**	Evaluation protocol
**Difficulty:**	Medium
**Blooms:**	Evaluate

**What principle about test data is stressed on the evaluation slide?**

1. Any data used to help prediction must not later be used to change the model
2. Only accuracy matters when evaluating
3. Test sets should be larger than training sets
4. You may tune the model after peeking at the test set once
5. It's fine to mix training and test if you cross-validate

**Question ID:**	M07-Q6
**Topic:**	Cross-validation
**Difficulty:**	Easy
**Blooms:**	Understand

**Cross-validation is introduced primarily to combat which issue?**

1. Overfitting and unreliable estimates from a single split
2. Class imbalance exclusively
3. Label noise only
4. GPU memory limits
5. Too many features

**Question ID:**	M07-Q7
**Topic:**	1-D decision threshold
**Difficulty:**	Medium
**Blooms:**	Apply

In one-dimensional two-class separation, the 'overlap region' between class distributions most directly corresponds to:

1. Training time
2. Regularization strength
3. Bayesian prior probability
4. Model bias
5. Classification error rate

**Question ID:**	M07-Q8
**Topic:**	Otsu thresholding
**Difficulty:**	Easy
**Blooms:**	Remember

Otsu's method chooses a threshold by optimizing which quantity?

1. Equalizing class priors
2. Minimizing within-class variance only
3. Maximizing inter-class variance (between-class)
4. Minimizing training loss
5. Maximizing overall accuracy on test data

**Question ID:**	M07-Q9
**Topic:**	LDA intuition
**Difficulty:**	Medium
**Blooms:**	Analyze

Linear Discriminant Analysis (LDA) can be viewed as finding:

1. A nonlinear kernel mapping to infinite dimensions
2. A decision tree that partitions space by axis-aligned cuts
3. A clustering of unlabeled data by k-means
4. A random forest that averages many trees
5. A projection that maximizes between-class separation relative to within-class scatter

**Question ID:**	M07-Q10
**Topic:**	LDA pros/cons
**Difficulty:**	Easy
**Blooms:**	Remember

**Which is listed as a \*pro\* of LDA in the deck?**

1. Nonlinear boundaries by default
2. Requires deep networks to perform well
3. Often overfits small datasets
4. **Usually doesn't overfit and works with much less data**
5. Slow classification at inference

<b>**Question ID:**</b>	M07-Q11
<b>**Topic:**</b>	Naïve Bayes assumptions
<b>**Difficulty:**</b>	Easy
<b>**Blooms:**</b>	Understand

**Why is Naïve Bayes called 'naïve' in these slides?**

1. **It assumes features are independent given the class**
2. It assumes k is chosen by cross-validation
3. It assumes infinite training data
4. It assumes a linear decision boundary
5. It assumes labels are uniformly distributed

<b>**Question ID:**</b>	M07-Q12
<b>**Topic:**</b>	Logistic regression
<b>**Difficulty:**</b>	Easy
<b>**Blooms:**</b>	Remember

**Logistic regression models the log-odds (logit) as:**

1. A sum of kernel functions over support vectors
2. **A linear function of input features ( $w^T x$ )**
3. A decision tree depth
4. A constant independent of input features
5. The reciprocal of a Gaussian density

<b>**Question ID:**</b>	M07-Q13
<b>**Topic:**</b>	Sigmoid link
<b>**Difficulty:**</b>	Easy
<b>**Blooms:**</b>	Remember

In the logistic regression derivation provided, which activation function maps  $\mathbf{z} = \mathbf{w}^T \mathbf{x}$  to a probability  $p$ ?

1.  $\sigma(z) = e^z / (1 + e^z)$
2. Hard threshold at  $z=0$
3.  $\tanh(z)$
4.  $\text{Softplus}(z) = \log(1 + e^z)$
5.  $\text{ReLU}(z) = \max(0, z)$

<b>**Question ID:**</b>	M07-Q14
<b>**Topic:**</b>	Linear separators
<b>**Difficulty:**</b>	Easy
<b>**Blooms:**</b>	Understand

Which statement about linear classifiers is consistent with the slides?

1. They never misclassify overlapping classes
2. They can only be used for regression
3. They require kNN distance voting
4. **They find a line/plane/hyperplane that separates classes**
5. They search for polynomial decision boundaries of degree  $\geq 2$

<b>**Question ID:**</b>	M07-Q15
<b>**Topic:**</b>	Baselines to try first
<b>**Difficulty:**</b>	Medium
<b>**Blooms:**</b>	Remember

Which of the following is \*not\* one of the four 'brain-dead' baseline algorithms the slides recommend trying first?

1. Logistic Regression
2. **Support Vector Machines with RBF kernel**
3. Naïve Bayes
4. Linear Discriminant Analysis
5. k-Nearest Neighbors

<b>**Question ID:**</b>	M07-Q16
<b>**Topic:**</b>	Projection quality
<b>**Difficulty:**</b>	Medium
<b>**Blooms:**</b>	Analyze

**Which plot-based intuition is used in the LDA section to contrast poor vs. better projections?**

1. Calibration curves
2. Box plots of residuals
3. ROC curves
4. Precision–recall curves
5. **Histograms of the projected feature**

<b>**Question ID:**</b>	M07-Q17
<b>**Topic:**</b>	Feature engineering motivation
<b>**Difficulty:**</b>	Medium
<b>**Blooms:**</b>	Understand

**The 'Feature Problem' slide motivates which broader idea?**

1. Labels are optional if we have enough features
2. We should avoid features and learn end-to-end only
3. Only image datasets need features
4. **Going from raw inputs to informative features is crucial**
5. Preprocessing is unnecessary with deep learning

<b>**Question ID:**</b>	M07-Q18
<b>**Topic:**</b>	Transparent vs Opaque models
<b>**Difficulty:**</b>	Easy
<b>**Blooms:**</b>	Understand

**Which pair best exemplifies the slides' distinction between transparent and opaque models?**

1. Random forests (transparent) vs. logistic regression (opaque)
2. **Decision trees (transparent) vs. deep neural networks (opaque)**
3. k-means (transparent) vs. linear regression (opaque)
4. Autoencoders (transparent) vs. KNN (opaque)
5. GANs (transparent) vs. PCA (opaque)

<b>**Question ID:**</b>	M07-Q19
<b>**Topic:**</b>	Feature importance in logistic regression
<b>**Difficulty:**</b>	Medium
<b>**Blooms:**</b>	Apply

**According to the M07B slides, what is a \*critical\* step before using logistic regression weights as feature importance?**

1. Apply PCA to all features
2. Normalize or standardize features to comparable scales
3. Remove the intercept term
4. Use L1 regularization
5. Convert all features to binary indicators

**Question ID:**	M07-Q20
**Topic:**	Post-hoc model-agnostic tools
**Difficulty:**	Easy
**Blooms:**	Understand

**Which statement about LIME and SHAP matches the deck?**

1. They require the model to be a decision tree
2. They are post-hoc explanation libraries that can be applied to black boxes
3. They are training algorithms for deep networks
4. They only provide global explanations, not local ones
5. They directly change model weights to be interpretable

**Question ID:**	M07-Q21
**Topic:**	Local vs Global
**Difficulty:**	Easy
**Blooms:**	Understand

**Which pair correctly contrasts local vs. global explanations as discussed?**

1. Local—visualize weights; Global—visualize a single gradient
2. Local—train a smaller model; Global—train a larger model
3. Local—understand model overall; Global—explain one prediction
4. Local—explain a single prediction; Global—summarize model behavior across the dataset
5. Local—optimize hyperparameters; Global—optimize loss

**Question ID:**	M07-Q22
**Topic:**	Surrogate models
**Difficulty:**	Medium
**Blooms:**	Apply

**Which technique is an example of a \*surrogate model\* approach?**

1. Ensembling many neural networks into a single predictor
2. Normalizing inputs to zero mean
3. Using dropout at test time
4. Computing the gradient of the loss w.r.t. inputs
5. **Training an interpretable decision tree to mimic a complex model**

<b>**Question ID:**</b>	M07-Q23
<b>**Topic:**</b>	Feature effects
<b>**Difficulty:**</b>	Medium
<b>**Blooms:**</b>	Understand

**Partial Dependence Plots (PDP) and Individual Conditional Expectation (ICE) primarily help with:**

1. **Understanding feature effects on predictions**
2. Speeding up training
3. Detecting data leakage
4. Generating counterfactuals by optimization
5. Estimating label noise

<b>**Question ID:**</b>	M07-Q24
<b>**Topic:**</b>	Counterfactuals
<b>**Difficulty:**</b>	Medium
<b>**Blooms:**</b>	Analyze

**A counterfactual explanation is best described as:**

1. A visualization of gradients in a neural network
2. An example from the training set closest to a query point
3. A proof that the model is globally optimal
4. **The minimal change to an input that would change the prediction**
5. Random noise added to test robustness

<b>**Question ID:**</b>	M07-Q25
<b>**Topic:**</b>	Motivation for XAI
<b>**Difficulty:**</b>	Easy
<b>**Blooms:**</b>	Understand

**Which of the following best captures the motivation slide 'Why do we need explainable models?'**

1. To increase training speed only
2. To support trust, debugging, fairness, and accountability in decisions (e.g., 'Why didn't I get the loan?')
3. To reduce the need for labeled data
4. Because explainability always improves accuracy
5. So we never need test sets again

**Question ID:**	M07-Q26
**Topic:**	Practical transparency
**Difficulty:**	Medium
**Blooms:**	Analyze

**Which statement aligns with the slides' nuance that 'both can use black-box methods'?**

1. Even simple models can act as black boxes if complexity or access prevents interpretation
2. Opaque models are always interpretable
3. Only neural networks are black boxes by definition
4. Transparent models cannot be used in black-box fashion
5. Black-box methods require no validation

**Question ID:**	M07-Q27
**Topic:**	Feature importance methods
**Difficulty:**	Medium
**Blooms:**	Understand

**Permutation Feature Importance (PFI) belongs to which family of explanation methods?**

1. Fairness constraints at training time
2. Counterfactual generation by gradient descent
3. Exact Shapley value computation for trees only
4. Perturbation-based, model-agnostic post-hoc importance
5. Model training objectives

**Question ID:**	M07-Q28
**Topic:**	XAI libraries
**Difficulty:**	Medium
**Blooms:**	Remember

**Which pairing is accurate per the slides' library list?**

1. SHAP — by scikit-learn core team only
2. LIME — exclusively for image models
3. Captum — official TensorFlow library
4. ELI5 — IBM's enterprise XAI suite
5. InterpretML — Microsoft's explainability toolkit

**Question ID:**	M07-Q29
**Topic:**	TreeSHAP
**Difficulty:**	Medium
**Blooms:**	Understand

**TreeSHAP is a specialized variant of SHAP intended for:**

1. Kernel density estimators
2. Convolutional neural networks only
3. Tree-based ensemble models
4. Hidden Markov Models
5. Model-agnostic linear models

**Question ID:**	M07-Q30
**Topic:**	Fairness notions (context)
**Difficulty:**	Medium
**Blooms:**	Evaluate

**Which statement about fairness and explainability is consistent with the surrounding course materials?**

1. Choosing one fairness metric can create trade-offs with others
2. Fairness is solved by increasing model capacity
3. All statistical definitions of fairness are mutually compatible
4. Only demographic parity matters in practice
5. Explainability eliminates the need for fairness analysis

**Question ID:**	MT-Q1
**Topic:**	N/A
**Difficulty:**	Easy
**Blooms:**	N/A

**Which of the following statements about Python dictionaries is TRUE?**

1. Dictionaries preserve the order of key insertion in all Python versions
2. Keys in a dictionary must be immutable objects
3. Duplicate keys are allowed and the values are stored in a list
4. Dictionary keys and values must be of the same type

**Question ID:**	MT-Q2
**Topic:**	N/A
**Difficulty:**	Easy
**Blooms:**	N/A

**What does the following Pandas code return?  
“python\\newlinedf.groupby  
'count'].reset\_index()\\newline“**

1. A DataFrame of unique categories and the sum of prices
2. A DataFrame summarising the mean and count of prices for each category
3. The original DataFrame unchanged
4. Raises an error because multiple aggregations are not allowed

**Question ID:**	MT-Q3
**Topic:**	N/A
**Difficulty:**	Medium
**Blooms:**	N/A

**Which of the following statements about the Central Limit Theorem is FALSE?**

1. It applies to the distribution of sample means
2. It requires that the population distribution be normal
3. It becomes more accurate as sample size increases
4. It implies that the sampling distribution approaches a normal distribution

**Question ID:**	MT-Q4
**Topic:**	N/A
**Difficulty:**	Medium
**Blooms:**	N/A

Consider the following SQL query:  
`SELECT  
customers.name, SUM(orders.amount) AS total  
FROM customers  
LEFT JOIN orders ON customers.id = orders.customer_id  
GROUP BY customers.name;`  
 What will this query return?

1. A table listing each order with the customer's name
2. The sum of order amounts per customer, including customers with no orders
3. Only customers who have placed at least one order
4. An error because GROUP BY cannot follow a LEFT JOIN

**Question ID:**	MT-Q5
**Topic:**	N/A
**Difficulty:**	Easy
**Blooms:**	N/A

In linear regression, the  $R^2$  (coefficient of determination) represents:

1. The correlation coefficient between the predictor and response
2. The proportion of variance in the response explained by the model
3. The square root of the residual sum of squares
4. The slope of the regression line

**Question ID:**	MT-Q6
**Topic:**	N/A
**Difficulty:**	Medium
**Blooms:**	N/A

Which of the following statements about hypothesis testing is TRUE?

1. A small p-value proves that the null hypothesis is false
2. Failing to reject the null hypothesis means it is true
3. A p-value is the probability, under the null hypothesis, of obtaining a result as extreme as the one observed
4. The null hypothesis always states that there is a difference

**Question ID:**	MT-Q7
**Topic:**	N/A
**Difficulty:**	Easy
**Blooms:**	N/A

**Which of the following statements about cross-validation is TRUE?**

1. In k-fold cross-validation, the value of k must equal the number of observations
2. Cross-validation provides a way to estimate model performance on unseen data
3. Cross-validation is only applicable to classification problems
4. Cross-validation always reduces overfitting

**Question ID:**	MT-Q8
**Topic:**	N/A
**Difficulty:**	Easy
**Blooms:**	N/A

**Write a Python function that returns the count of unique values in a Pandas Series without using ‘nunique()’.**

**Question ID:**	MT-Q9
**Topic:**	N/A
**Difficulty:**	Medium
**Blooms:**	N/A

**Given two arrays ‘a’ and ‘b’ of equal length, write a NumPy expression that computes the cosine similarity between them.**

**Question ID:**	MT-Q10
**Topic:**	N/A
**Difficulty:**	Easy
**Blooms:**	N/A

**Write an SQL query to return the top 3 products with the highest average rating from a table ‘reviews‘ with columns ‘product\_id‘, ‘rating‘.**

**Question ID:**	MT-Q11
**Topic:**	N/A
**Difficulty:**	Medium
**Blooms:**	N/A

**Using Matplotlib, plot the probability density function of a standard normal distribution over the interval [-3, 3].**

**Question ID:**	M04-Q1
**Topic:**	Story vs. Model
**Difficulty:**	Easy
**Blooms:**	Remember

**In this course, a 'model' primarily serves to...**

1. Guarantee that data will fit a normal distribution
2. Tell a compelling narrative that motivates the audience
3. Eliminate randomness from the world entirely
4. Automatically generate high-quality plots without code
5. **Provide a mathematical abstraction that summarizes patterns and supports prediction**

**Question ID:**	M04-Q2
**Topic:**	Deterministic models
**Difficulty:**	Medium
**Blooms:**	Understand

**Which example best illustrates a deterministic model mentioned in lecture?**

1. Choosing chart color palettes
2. Predicting tomorrow's weather with perfect certainty
3. **Using  $F = ma$  to compute motion under known forces**
4. Shuffling labels in a permutation test
5. Estimating coin-flip bias by simulation

**Question ID:**	M04-Q3
**Topic:**	Random Variable
**Difficulty:**	Medium
**Blooms:**	Apply

**A random variable is best described as:**

1. A list of frequencies for categories only
2. A quantity that must be continuous
3. A fixed constant determined once and for all
4. **A function that maps outcomes of an experiment to numbers**

- Any number sampled from a histogram

**Question ID:**	M04-Q4
**Topic:**	Density vs Distribution
**Difficulty:**	Hard
**Blooms:**	Analyze

**For a continuous distribution with density  $p(x)$ , a core property is that:**

- $\sum p(x)$  over the data equals 1
- $p(x)$  must always be less than 1
- The derivative of  $p(x)$  is constant
- $p(x)$  must be symmetric around the mean
- $\int p(x) dx$  over the support equals 1

**Question ID:**	M04-Q5
**Topic:**	Discrete vs Continuous
**Difficulty:**	Easy
**Blooms:**	Evaluate

**In the slides, which notation pairing is used (abuse of notation aside)?**

- Discrete variables must be converted to z-scores
- Only continuous variables can be modeled
- Both must be written as integrals only
- P(X) for discrete distributions,  $p(X)$  for continuous densities
- $p(X)$  for discrete, P(X) for continuous

**Question ID:**	M04-Q6
**Topic:**	Joint vs Conditional
**Difficulty:**	Medium
**Blooms:**	Create

**Which expression corresponds to the conditional probability as defined in the slides?**

- $P(C|X) = P(C) \cdot P(X)$
- $P(X|C) = P(C,X) / P(C)$
- $P(C|X) = P(C,X) / P(X)$
- $P(C|X) = P(C) + P(X)$

$$5. P(C-X) = 1 - P(C,X)$$

**Question ID:**	M04-Q7
**Topic:**	Bayes Theorem
**Difficulty:**	Medium
**Blooms:**	Remember

**In Bayes' theorem, the term 'likelihood' refers to:**

1.  $P(C)$
2.  $P(C-X)$
3.  $P(X-C)$
4.  $P(C \text{ and } X)$  divided by zero
5.  $P(X)$  only when  $X$  is discrete

**Question ID:**	M04-Q8
**Topic:**	MAP Estimate
**Difficulty:**	Hard
**Blooms:**	Understand

**The Maximum A Posteriori (MAP) estimate selects:**

1. Any class with probability above 0.1
2. **The class C that maximizes  $P(C-X)$**
3. The class that maximizes  $P(X-C)$  regardless of priors
4. The class with the smallest label alphabetically
5. The parameter that minimizes variance only

**Question ID:**	M04-Q9
**Topic:**	Maximum Likelihood
**Difficulty:**	Easy
**Blooms:**	Apply

**If priors  $P(C)$  are unknown or assumed equal, maximizing  $P(C-X)$  reduces to maximizing:**

1. The number of features in  $X$
2. The sample median
3.  **$P(X-C)$  (the likelihood)**
4.  $P(X)$  which is constant across  $C$  and thus decisive

## 5. P(C)

**Question ID:**	M04-Q10
**Topic:**	MLE for Normal
**Difficulty:**	Medium
**Blooms:**	Analyze

For a Normal distribution fitted by Maximum Likelihood to i.i.d. data, the parameter estimates are:

1. Any pair that minimizes the max error
2. Both parameters are always zero
3. Sample median for  $\mu$  and IQR for  $\sigma$
4. Sample mean for  $\mu$  and sample standard deviation for  $\sigma$
5. Skewness for  $\mu$  and kurtosis for  $\sigma$

**Question ID:**	M04-Q11
**Topic:**	Frequentist simulation
**Difficulty:**	Medium
**Blooms:**	Evaluate

Observing 27 heads in 40 fair-coin flips is evaluated in lecture by:

1. Replacing missing flips with the mean
2. Assuming the coin is biased without testing
3. Comparing the observed statistic to a simulated null distribution
4. Counting only the last five flips
5. Using a geospatial plot of flips

**Question ID:**	M04-Q12
**Topic:**	Permutation Test
**Difficulty:**	Hard
**Blooms:**	Create

Permutation tests (label shuffling) assess:

1. How to order categories alphabetically
2. Whether variables are exactly deterministic
3. Which distribution family is true a priori
4. Whether the observed group difference could arise by chance under the null

5. The best color map for plots

**Question ID:**	M04-Q13
**Topic:**	Bootstrapping
**Difficulty:**	Easy
**Blooms:**	Remember

**Bootstrapping, as presented, is primarily used to:**

1. Approximate the sampling distribution of a statistic by resampling with replacement
2. Force data to follow a Normal distribution
3. Guarantee smaller variance than the original sample
4. Encrypt sensitive variables
5. Eliminate outliers by deletion

**Question ID:**	M04-Q14
**Topic:**	EDA
**Difficulty:**	Medium
**Blooms:**	Understand

**A central purpose of EDA is to:**

1. Guarantee future market trends
2. Reveal structure, anomalies, and relationships in raw data before modeling
3. Replace statistics with visuals
4. Finalize the causal model and deploy
5. Remove all missing data automatically

**Question ID:**	M04-Q15
**Topic:**	Data Storytelling
**Difficulty:**	Medium
**Blooms:**	Apply

**Why combine narrative with visuals in data storytelling per the lecture?**

1. Emotions prove scientific truth on their own
2. Stories remove uncertainty from estimates
3. Stakeholders only read fiction
4. Narrative replaces the need for any chart
5. Narratives make insights memorable and persuasive beyond clean visuals

**Question ID:**	M04-Q16
**Topic:**	Narrative Components
**Difficulty:**	Hard
**Blooms:**	Analyze

**Which is one of the narrative components highlighted for data stories?**

1. Kernel bandwidth
2. Hash seed
3. ASCII table index
4. **Conflict**
5. Recursion depth

**Question ID:**	M04-Q17
**Topic:**	Story Structure
**Difficulty:**	Easy
**Blooms:**	Evaluate

**A simple story arc used to frame analyses in the slides is:**

1. Read → Write → Sleep
2. Inference → Proof → Axiom
3. **Setup → Conflict → Resolution**
4. North → East → West
5. Encode → Decode → Re-encode

**Question ID:**	M04-Q18
**Topic:**	Save the Cat
**Difficulty:**	Medium
**Blooms:**	Create

**The 'Save the Cat' beat sheet is referenced to emphasize:**

1. **Borrowing proven narrative sequences to structure data stories**
2. Replacing statistics with anecdote
3. Only using 3 slides per talk
4. Using cats to label datasets for CV tasks
5. Avoiding structure in favor of randomness

**Question ID:**	M04-Q19
**Topic:**	Academic Model
**Difficulty:**	Medium
**Blooms:**	Remember

**Which element belongs to the academic paper/presentation model listed?**

1. Always use 3D charts
2. Never disclose limitations
3. Choose a punchline first, skip data
4. GPU cache coherence
5. Why / So what?

**Question ID:**	M04-Q20
**Topic:**	MAP vs ML
**Difficulty:**	Hard
**Blooms:**	Understand

**Compared with ML, MAP will differ when:**

1. Priors over classes/parameters are non-uniform and informative
2. There is no data at all
3. Data are measured in centimeters not inches
4. You shuffle labels many times
5. The likelihood is constant across classes

**Question ID:**	M04-Q21
**Topic:**	Prior
**Difficulty:**	Easy
**Blooms:**	Apply

**In Bayes' theorem, the prior  $P(C)$  represents:**

1. Belief about classes before seeing current data
2. Noise added to stabilize training
3. A probability that must always be 0.5
4. The maximum of the likelihood
5. A resampling weight in bootstrapping

**Question ID:**	M04-Q22
**Topic:**	Evidence
**Difficulty:**	Medium
**Blooms:**	Analyze

In the factorization  $P(C|X) = P(X|C)P(C)/P(X)$ ,  $P(X)$  (the 'evidence') mainly serves to:

1. Set the sampling rate for bootstrapping
2. Choose color bars for heatmaps
3. Encode class labels as integers
4. Maximize the likelihood for each class
5. Normalize the posterior over classes so probabilities sum to 1

**Question ID:**	M04-Q23
**Topic:**	Joint Probability
**Difficulty:**	Medium
**Blooms:**	Evaluate

The joint distribution  $P(C,X)$  can be expressed as:

1.  $P(C|X) P(C)$
2.  $P(X)/P(C)$
3.  $P(X|C) P(C)$
4.  $P(C) + P(X)$
5.  $1 - P(C|X)$

**Question ID:**	M04-Q24
**Topic:**	Posterior normalization
**Difficulty:**	Hard
**Blooms:**	Create

For fixed measurements  $X$ , the slides note that  $\sum^i P(C^i|X)$  equals:

1.  $P(X)$
2. The number of classes
3. 1
4. 0
5.  $E[X]$

**Question ID:**	M04-Q25
**Topic:**	Discrete Distributions
**Difficulty:**	Easy
**Blooms:**	Remember

**Which is an example of a discrete distribution family discussed in the slides?**

1. Fourier
2. Cauchy (continuous heavy-tailed)
3. Binomial
4. Beta-prime (advanced)
5. Gaussian mixture with infinite components

**Question ID:**	M04-Q26
**Topic:**	Continuous Distributions
**Difficulty:**	Medium
**Blooms:**	Understand

**Which is included among continuous distributions referenced in lecture?**

1. Normal
2. Bernoulli
3. Directory
4. Categorical
5. Poisson

**Question ID:**	M04-Q27
**Topic:**	scipy.stats
**Difficulty:**	Medium
**Blooms:**	Apply

**The lecture notes mention practical access to many distributions via:**

1. ‘scipy.stats’
2. ‘sklearn.metrics’
3. ‘os.path’
4. ‘torch.vision’
5. ‘seaborn.misc’

**Question ID:**	M04-Q28
**Topic:**	Viz + Narrative
**Difficulty:**	Hard
**Blooms:**	Analyze

**“Good visualization is not enough” implies that analysts should:**

1. Only present statistically significant results
2. Prefer 3D plots for memorability
3. Pair clear visuals with narrative and purpose tailored to the audience
4. Use animations for every chart
5. Avoid visuals entirely

**Question ID:**	M04-Q29
**Topic:**	Stakeholders
**Difficulty:**	Easy
**Blooms:**	Evaluate

**Which of the following appears on the slides as a stakeholder to consider in the narrative?**

1. Kernel developers
2. Astrologers
3. Compiler vendors
4. Captcha solvers
5. Shareholders

**Question ID:**	M04-Q30
**Topic:**	Analytics Types
**Difficulty:**	Medium
**Blooms:**	Create

**Prescriptive analytics focuses on:**

1. Sorting categories alphabetically
2. Recommending actions based on analysis
3. Replacing domain experts with charts
4. Labelling axes after plotting
5. Describing what happened only

**Question ID:**	M04-Q31
**Topic:**	Analytics Types
**Difficulty:**	Medium
**Blooms:**	Remember

**Diagnostic analytics asks primarily:**

1. Why did this happen?
2. What should we do next?
3. Who is to blame regardless of data?
4. How to export to PDF?
5. What color should bars be?

**Question ID:**	M04-Q32
**Topic:**	Analytics Types
**Difficulty:**	Hard
**Blooms:**	Understand

**Predictive analytics, per the slide taxonomy, addresses:**

1. What happened, historically?
2. What may happen in the future?
3. Where to place legends
4. Which category is funniest?
5. How to set DPI for figures

**Question ID:**	M04-Q33
**Topic:**	MAP Decision
**Difficulty:**	Easy
**Blooms:**	Apply

**Under MAP classification with two classes, a tie in posterior probabilities should be resolved by:**

1. Set both labels at once
2. Always pick the positive class
3. Switch to a different dataset silently
4. Pick randomly without recording
5. An explicit tie-break policy (e.g., choose class with lower cost)

**Question ID:**	M04-Q34
**Topic:**	MLE for Bernoulli/Binomial
**Difficulty:**	Medium
**Blooms:**	Analyze

For coin flips with H heads out of N trials, the ML estimate of  $\theta=P(H)$  is:

1.  $(H-1)/(N-1)$  always
2. 0.5 regardless of data
3.  $(H+1)/(N+2)$  (Laplace smoothing)
4.  $N/H$
5.  $H/N$

**Question ID:**	M04-Q35
**Topic:**	Evidence term
**Difficulty:**	Medium
**Blooms:**	Evaluate

In class-by-class comparisons,  $P(X)$  often cancels because:

1. It is always 1
2. It equals 0 for balanced datasets
3. It is undefined and ignored
4. It is the same for all classes for a fixed X
5. It depends only on priors

**Question ID:**	M04-Q36
**Topic:**	Q-Q Plot
**Difficulty:**	Hard
**Blooms:**	Create

Which tool helps assess whether data follow a theoretical distribution?

1. Stacked area chart
2. Q-Q plot comparing sample vs theoretical quantiles
3. Pie chart
4. Bar chart of categories
5. Chord diagram

**Question ID:**	M04-Q37
**Topic:**	Simulation Precision
**Difficulty:**	Easy
**Blooms:**	Remember

**Repeating simulations many times primarily helps to:**

1. Guarantee significance
2. Avoid computing summaries
3. Quantify variability of estimates (e.g., Monte Carlo error)
4. Reduce the true variance in the population
5. Change the historical data

**Question ID:**	M04-Q38
**Topic:**	Model Evaluation
**Difficulty:**	Medium
**Blooms:**	Understand

**Cross-validation (mentioned for future lectures) is mainly used to:**

1. Sort features alphabetically
2. Compute exact posteriors in closed form
3. Visualize categorical data
4. Estimate generalization performance by repeated train/test splits
5. Fit parameters of a Normal distribution

**Question ID:**	M04-Q39
**Topic:**	Time Series
**Difficulty:**	Medium
**Blooms:**	Apply

**A caution for time series modeling highlighted in the notes is to:**

1. Replace time with row index randomly
2. Always difference until stationary regardless of context
3. Transform units to percentages only
4. Avoid using future observations to train models evaluated on the past
5. Impute time with the mean

**Question ID:**	M04-Q40
**Topic:**	Missing Data
**Difficulty:**	Hard
**Blooms:**	Analyze

**Which Python package was explicitly recommended to visualize missingness patterns?**

1. pytorch.missviz
2. matplotlib.auto\_na
3. **missingno**
4. seaborn.naheat
5. cv2.missing

<b>**Question ID:**</b>	M04-Q41
<b>**Topic:**</b>	Narrative Conflict
<b>**Difficulty:**</b>	Easy
<b>**Blooms:**</b>	Evaluate

**In the narrative example of a social-media environmental crisis, the 'conflict' includes:**

1. Elimination of all randomness in data
2. A change in color theme for charts
3. Switching from PNG to SVG
4. **Sales drop among younger customers following a viral tweet**
5. Sorting bars by height

<b>**Question ID:**</b>	M04-Q42
<b>**Topic:**</b>	Narrative Resolution
<b>**Difficulty:**</b>	Medium
<b>**Blooms:**</b>	Create

**The proposed 'resolution' in that narrative suggests:**

1. Ignoring stakeholders entirely
2. Re-labeling axes to look better
3. Reducing sample size to avoid outliers
4. Suppressing negative comments without analysis
5. **A data-backed pivot to sustainability with cost/payoff visualization**

<b>**Question ID:**</b>	M04-Q43
<b>**Topic:**</b>	Story vs Model
<b>**Difficulty:**</b>	Medium
<b>**Blooms:**</b>	Remember

## Difference between data story and model emphasized in lecture:

1. Story is high-level meaning; model is a mathematical abstraction
2. They are identical in purpose
3. Model always implies causation
4. Story contains equations; model is the talk title
5. Story never uses any charts

**Question ID:**	M04-Q44
**Topic:**	Discrete vs Continuous
**Difficulty:**	Hard
**Blooms:**	Understand

Which statement is TRUE regarding discrete vs. continuous modeling per lecture?

1. Discrete probabilities sum to 1; continuous densities integrate to 1
2. Continuous probabilities sum to 1 over observed bins only
3. Neither requires normalization
4. Discrete distributions cannot have expectations
5. Both integrate to 1 only

**Question ID:**	M04-Q45
**Topic:**	MAP Decision
**Difficulty:**	Easy
**Blooms:**	Apply

The slides define the 'best guess' class given X as:

1.  $\text{argmax}^i P(C^i | X)$
2. Choose the smallest index i
3.  $\text{argmin}^i P(C^i | X)$
4. Random selection
5.  $\text{argmax}^i P(X)$

**Question ID:**	M04-Q46
**Topic:**	Deterministic Models
**Difficulty:**	Medium
**Blooms:**	Analyze

**A finite-state machine (FSM) controlling an elevator is used to illustrate:**

1. A continuous-time Markov chain only
2. A visualization library
3. An image compression algorithm
4. A stochastic bootstrap procedure
5. A deterministic system with defined transitions

**Question ID:**	M04-Q47
**Topic:**	Likelihood
**Difficulty:**	Medium
**Blooms:**	Evaluate

**Why is likelihood often easier to compute than the posterior?**

1. Likelihood ignores the data completely
2. Likelihood equals prior exactly
3. Posterior is independent of X
4. Posterior requires normalization by  $P(X)$ , which involves summing/integrating over classes or parameters
5. Posterior is always undefined

**Question ID:**	M04-Q48
**Topic:**	Simulation
**Difficulty:**	Hard
**Blooms:**	Create

**When no convenient theoretical distribution is available, the notes recommend:**

1. Pick parameters to match your hypothesis
2. Assume Normality without checking
3. Use simulation to approximate uncertainty (e.g., repeated runs)
4. Abandon modeling entirely
5. Report a single run as definitive

**Question ID:**	M04-Q49
**Topic:**	Frequentist Computation
**Difficulty:**	Easy
**Blooms:**	Remember

**The 'frequentist/computational' approach in lecture emphasizes:**

1. Replacing statistics with narrative
2. Tuning figure DPI for accuracy
3. Measuring distance in pixels
4. Selecting priors to encode beliefs
5. **Using resampling/shuffling to approximate sampling distributions**

<b>**Question ID:**</b>	M04-Q50
<b>**Topic:**</b>	EDA Position
<b>**Difficulty:**</b>	Medium
<b>**Blooms:**</b>	Understand

**Why is EDA positioned early in the workflow?**

1. To finalize the deployment pipeline first
2. To increase file sizes for charts
3. To ensure only balanced classes are kept
4. **To surface data quality issues and structure that inform modeling choices**
5. Because modeling is obsolete

<b>**Question ID:**</b>	M04-Q1
<b>**Topic:**</b>	Data Storytelling Components
<b>**Difficulty:**</b>	Easy
<b>**Blooms:**</b>	Remember

**Which trio forms the core 'puzzle pieces' of data storytelling emphasized in the lecture?**

1. **Data, Narrative, Visualizations**
2. Tables, Dashboards, Animation
3. Python, SQL, Excel
4. Fonts, Colors, Icons
5. Neural Nets, GPUs, TPUs

<b>**Question ID:**</b>	M05A-Q1
<b>**Topic:**</b>	Structured vs. Unstructured Data
<b>**Difficulty:**</b>	Easy
<b>**Blooms:**</b>	Remember

**Which statement best captures why data is called 'structured' in this context?**

1. It lives in a single file rather than multiple tables.
2. It is stored as images and audio files.
3. It contains only numeric values with no missing entries.
4. It comes from sensors in time order only.
5. **It has a predefined schema with labeled columns and consistent record shape.**

<b>**Question ID:**</b>	M06-Q1
<b>**Topic:**</b>	ML Process: Problem Framing
<b>**Difficulty:**</b>	Medium
<b>**Blooms:**</b>	Analyze

**Which step most directly ensures that the problem you're solving is actually a machine-learning problem and not a simple rules or query task?**

1. Check whether a labeled target exists and if patterns must generalize to new data
2. Visualize predictions with a confusion matrix
3. Collect more data first
4. Tune hyperparameters with cross-validation
5. Train a baseline linear model

<b>**Question ID:**</b>	MT-Q1
<b>**Topic:**</b>	N/A
<b>**Difficulty:**</b>	Easy
<b>**Blooms:**</b>	N/A

**Which of the following statements about Python dictionaries is TRUE?**

1. Dictionaries preserve the order of key insertion in all Python versions
2. **Keys in a dictionary must be immutable objects**
3. Duplicate keys are allowed and the values are stored in a list
4. Dictionary keys and values must be of the same type