

# Tecnologías de gestión de información no estructurada

## Preguntas generados con ChatGPT a partir del temario

1. What is the main goal of Big Data in textual data management?
  - a) Storing vast amounts of data
  - b) Simplifying data storage structures
  - c) Transforming raw data into actionable knowledge
  - d) Reducing redundancy in databases
2. Which of the following is an example of textual data as Big Data?
  - a) Temperature logs
  - b) Scientific literature
  - c) Financial spreadsheets
  - d) Network configurations
3. What proportion of the world's data generated annually is textual?
  - a) Negligible
  - b) A significant portion
  - c) Entirely textual
  - d) Only a fraction
4. What characterizes structured data?
  - a) Random arrangement of content
  - b) Well-defined schemas
  - c) Exclusively textual representation
  - d) Lack of metadata
5. Why is textual data challenging to process?
  - a) It consumes less storage
  - b) It lacks explicit structure
  - c) It is highly numeric
  - d) It requires no computational interpretation
6. What does Natural Language Processing (NLP) rely on for understanding text?
  - a) Heuristic and statistical approaches
  - b) Manual interpretation
  - c) Only deep semantic models
  - d) Strictly predefined templates
7. Which system allows quick access to relevant textual data?
  - a) NLP models
  - b) Text retrieval systems
  - c) Data clustering tools
  - d) Semantic models
8. What is the focus of text mining in Big Data?
  - a) Organizing data alphabetically
  - b) Generating raw textual data
  - c) Extracting valuable knowledge from text
  - d) Reducing text length
9. In the "pull" mode of accessing information, what role does the user play?
  - a) Passive recipient
  - b) Active seeker
  - c) Content generator
  - d) Recommendation builder
10. What is the goal of a Textual Information System (TIS)?
  - a) Create unstructured data repositories
  - b) Match relevant information to users' needs
  - c) Replace traditional search engines
  - d) Perform manual data classification
11. What process involves grouping similar textual objects?
  - a) Text retrieval
  - b) Clustering
  - c) Categorization
  - d) Visualization
12. Which task generates concise representations of large textual content?
  - a) Extraction
  - b) Summarization
  - c) Navigation
  - d) Filtering
13. What is the main purpose of sentiment analysis in textual data?
  - a) Organizing text lexically
  - b) Extracting and categorizing opinions
  - c) Visualizing content patterns
  - d) Translating textual data
14. How does a recommendation system in a TIS function?
  - a) By extracting patterns in data
  - b) By proactively suggesting relevant information
  - c) By creating unsupervised learning models
  - d) By analyzing user-generated summaries
15. What technique discovers patterns and trends in large textual datasets?
  - a) Filtering
  - b) Thematic analysis
  - c) Snippet generation
  - d) Data reduction
16. Which of the following is a characteristic of unsupervised thematic analysis?
  - a) Predefined themes
  - b) Algorithm-generated topics
  - c) Strict manual labeling
  - d) Reliance on metadata only
17. Which approach is commonly used in most TIS for text representation?
  - a) Semantic deep learning

- b) Bag-of-words model
  - c) Ontology-based parsing
  - d) Strict rule-based systems
18. What is a key advantage of visualizing textual data?
- a) Avoiding analysis complexity
  - b) Enhancing pattern recognition
  - c) Eliminating redundancy
  - d) Simplifying query syntax
19. What does information extraction focus on?
- a) Reducing data size
  - b) Identifying entities and their relationships
  - c) Simplifying textual metadata
  - d) Categorizing user preferences
20. What does a snippet in a search result typically represent?
- a) A dynamic summary based on the query
  - b) An unrelated document
  - c) The metadata of the file
  - d) The full text of the search result
21. What is the primary goal of Natural Language Processing (NLP)?
- a) Programming efficient algorithms
  - b) Enabling machines to understand and generate human language
  - c) Enhancing computer hardware performance
  - d) Managing large datasets
22. Which of the following tasks is associated with lexical analysis?
- a) Determining the meaning of a sentence
  - b) Identifying grammatical relationships
  - c) Breaking down text into meaningful units like words
  - d) Resolving ambiguity in context
23. What is semantic analysis in NLP primarily concerned with?
- a) Understanding sentence meaning
  - b) Resolving syntactic ambiguities
  - c) Analyzing word frequencies
  - d) Determining sentence structure
24. In NLP, what is an example of pragmatic analysis?
- a) Identifying parts of speech
  - b) Determining the intent behind a statement
  - c) Parsing grammatical structures
  - d) Generating language models
25. What does syntactic analysis focus on?
- a) Understanding the relationships between words in a sentence
  - b) Translating words into numerical representations
  - c) Recognizing emotions in text
  - d) Creating visualizations of data
26. Why is NLP considered "complete in AI"?
- a) It requires vast computational resources
  - b) It is as difficult as other complex AI problems
  - c) It focuses on data mining
  - d) It always requires deep learning techniques
27. Which of the following is a common challenge in NLP?
- a) Lack of computational tools
  - b) Limited datasets
  - c) Ambiguities in natural language
  - d) Inefficient programming languages
28. What significant shift occurred in NLP research after the 1980s?
- a) Transition from symbolic to statistical approaches
  - b) Focus on hardware advancements
  - c) Abandonment of speech recognition
  - d) Exclusive reliance on human linguistic knowledge
29. What was a limitation of symbolic approaches to NLP?
- a) Dependency on machine learning techniques
  - b) Inability to scale for real-world applications
  - c) Excessive reliance on unstructured data
  - d) Lack of support from computational linguists
30. Early NLP systems like Eliza relied on:
- a) Deep learning algorithms
  - b) Superficial rule-based techniques
  - c) Statistical analysis of large datasets
  - d) Neural network architectures
31. What does POS tagging in NLP refer to?
- a) Predicting sentence probability
  - b) Assigning parts of speech to words
  - c) Generating text data
  - d) Analyzing grammatical dependencies
32. Which statistical model assumes words in a sequence are generated independently?
- a) Bigram model
  - b) Unigram model
  - c) Neural network model
  - d) Contextual model
33. What is the primary purpose of smoothing in language models?
- a) Enhance training efficiency
  - b) Reduce data dimensionality
  - c) Assign probabilities to unseen events
  - d) Optimize computational speed
34. What is a key advantage of statistical approaches over symbolic ones in NLP?
- a) Higher dependency on handcrafted rules
  - b) Greater robustness to real-world variability
  - c) Requirement for linguistic experts
  - d) Elimination of ambiguity in text
35. Which of the following is an example of a high-level text representation?
- a) Character sequences
  - b) Word embeddings with semantic relationships
  - c) Unstructured text documents
  - d) Paragraph-level sentiment scores
36. In speech recognition, language models help by:

- a) Identifying phonemes
  - b) Predicting probable word sequences
  - c) Filtering noise from audio input
  - d) Translating text to audio
37. How do language models assist in text categorization?
- a) They tokenize input sentences
  - b) They assign probabilities to thematic topics
  - c) They parse grammatical structures
  - d) They generate summaries of text
38. Named Entity Recognition (NER) involves:
- a) Identifying grammatical errors
  - b) Recognizing specific entities like names and locations
  - c) Analyzing syntactic structures
  - d) Detecting sentiment in text
39. What is one limitation of semantic analysis in NLP?
- a) Requires too much memory
  - b) It cannot handle labeled datasets
  - c) Difficulty in processing ambiguous sentences
  - d) Limited applicability in non-English texts
40. Which task would benefit most from discourse analysis?
- a) Translating isolated sentences
  - b) Summarizing multi-paragraph documents
  - c) Identifying word synonyms
  - d) Counting word frequencies
41. What does contextual word embedding refer to?
- a) Representing words with fixed vector sizes
  - b) Capturing meaning based on surrounding words
  - c) Tokenizing text into characters
  - d) Eliminating ambiguity from language
42. A unigram model is most suitable for:
- a) Capturing long-range dependencies
  - b) Representing independent word probabilities
  - c) Parsing syntactic structures
  - d) Analyzing discourse coherence
43. Why is deep analysis of text representation often challenging?
- a) Limited linguistic knowledge
  - b) Dependence on massive labeled datasets
  - c) Over-reliance on character-level models
  - d) Insufficient computational power
44. How does a statistical model infer semantic relationships?
- a) Using predefined rules
  - b) Analyzing word co-occurrences
  - c) Assigning probabilities manually
  - d) Ignoring outliers in data
45. Which feature differentiates pragmatic analysis from semantic analysis?
- a) Focus on sentence grammar
  - b) Dependence on textual context and intent
  - c) Use of statistical models
  - d) Reliance on syntactic parsing
46. What is one of the primary roles of textual data access technologies?
- a) Reducing hardware dependencies
  - b) Enhancing database storage efficiency
  - c) Connecting users with relevant information at the right time
  - d) Automating all data retrieval processes
47. In textual data access, what does the pushmode imply?
- a) Users actively query a system for information
  - b) The system proactively delivers relevant information to users
  - c) Data is manually transferred between users
  - d) Systems focus solely on structured data
48. Which of the following best describes the pullmode in textual data access?
- a) Systems automatically recommend content
  - b) Users initiate the retrieval of relevant data
  - c) Data is filtered through machine learning models
  - d) Retrieval relies solely on pre-defined rules
49. What type of information needs is often associated with the pullmode?
- a) Long-term, stable needs
  - b) Dynamic, ad hoc needs
  - c) Fixed and structured needs
  - d) Automated, unsupervised needs
50. Why is ad hoc retrieval considered challenging?
- a) It requires structured datasets
  - b) There is limited user feedback for optimization
  - c) It only applies to numeric data
  - d) It mandates constant user interaction
51. What is the primary benefit of integrating push and pull modes in a system?
- a) To eliminate user input requirements
  - b) To ensure consistent data analysis
  - c) To maximize flexibility in accessing data
  - d) To automate data structuring processes
52. In multimodal interaction systems, what is one advantage of navigation?
- a) It guarantees accurate query results
  - b) Users can explore topics without a fixed goal
  - c) It eliminates the need for keyword-based searches
  - d) All results are pre-ranked for relevance
53. What does a topic map in a search interface typically support?
- a) Statistical analysis of search queries
  - b) User navigation through interconnected topics
  - c) Automatic data labeling for machine learning
  - d) Compression of large datasets
54. What action allows users to explore documents within a thematic region on a topic map?
- a) Submitting new queries
  - b) Clicking on thematic nodes
  - c) Adjusting system settings
  - d) Switching between datasets

55. What is an example of a long-range jump in a topic map system?
  - a) Zooming in on a child node
  - b) Submitting a new search query
  - c) Navigating between neighboring nodes
  - d) Viewing a local map segment
56. What is the primary goal of a text retrieval (TR) system?
  - a) To store textual data efficiently
  - b) To identify relevant documents for a user query
  - c) To automate database schema creation
  - d) To categorize images alongside text
57. How does a text retrieval system differ from a database system?
  - a) Text retrieval systems manage structured data
  - b) Database systems are not designed for numeric operations
  - c) Text retrieval focuses on unstructured text
  - d) Database systems emphasize vague query handling
58. Why are keyword-based queries often insufficient for text retrieval?
  - a) They only support long-form queries
  - b) They depend on predefined database schemas
  - c) They lack specificity and completeness
  - d) They cannot handle numeric datasets
59. In text retrieval, what is the role of relevance feedback?
  - a) To filter irrelevant datasets
  - b) To improve future query results using user feedback
  - c) To enhance system processing speeds
  - d) To structure unstructured text
60. What is a key challenge in designing text retrieval systems?
  - a) Managing structured data
  - b) Ensuring low computational overhead
  - c) Modeling user information needs
  - d) Standardizing query languages
61. What distinguishes document selection from ranking in retrieval systems?
  - a) Selection involves absolute relevance classification
  - b) Ranking eliminates user decision-making
  - c) Selection depends on user-defined thresholds
  - d) Ranking focuses on binary classification
62. Why is ranking often preferred over selection in document retrieval?
  - a) It simplifies query processing
  - b) It allows users to decide cut-off points based on relevance
  - c) It eliminates query ambiguity
  - d) It is computationally less demanding
63. What is the purpose of scoring functions in ranking?
  - a) To classify documents into structured datasets
  - b) To measure the probability of a document's relevance
  - c) To analyze semantic relationships between terms
  - d) To generate automatic summaries
64. Which strategy is commonly used in ranking-based systems?
  - a) Binary classification
  - b) Probability-based ordering
  - c) Manual document tagging
  - d) Random document sampling
65. What assumption underpins probability ranking principles in retrieval?
  - a) Document relevance is independent of user feedback
  - b) A user's utility for one document is independent of others
  - c) Query complexity is inversely proportional to relevance
  - d) Only exact matches are deemed relevant
66. What is one major challenge in combining push and pull nodes?
  - a) Integrating structured and unstructured data
  - b) Predicting long-term user needs
  - c) Balancing system-driven and user-driven interactions
  - d) Avoiding over-recommendation of irrelevant data
67. How does a multimodal system enhance text retrieval?
  - a) By prioritizing structured over unstructured data
  - b) By seamlessly integrating query and navigation modes
  - c) By automating relevance feedback mechanisms
  - d) By filtering out incomplete datasets
68. Why is document classification more challenging than ranking?
  - a) It requires predefined thresholds for all queries
  - b) It mandates real-time user interaction
  - c) It depends on absolute measures of relevance
  - d) It excludes probabilistic approaches
69. What makes ad hoc information retrieval particularly difficult?
  - a) Reliance on structured data inputs
  - b) Lack of user feedback for query optimization
  - c) Dependence on static datasets
  - d) Overlap with database retrieval techniques
70. Which of the following improves user experience in retrieval systems?
  - a) Using longer keyword queries
  - b) Enhancing thematic navigation tools
  - c) Restricting access to unstructured data
  - d) Avoiding multimodal interactions
71. In similarity-based retrieval models, the relevance of a document is determined by:
  - a) Its metadata attributes compared to other documents.
  - b) The extent to which its content overlaps with the query terms.
  - c) The probability distribution of terms across all documents.
  - d) Its absolute frequency of term occurrences regardless of context.
72. Which of the following best describes the role of the binary random variable  $R$  in probabilistic retrieval models?

- a) It identifies whether a document is correctly indexed.
  - b) It measures the overlap between documents and the query vector.
  - c) It indicates whether a document is relevant to a given query.
  - d) It calculates the similarity score between two documents.
73. What is the primary limitation of the classic probabilistic model when estimating relevance?
- a) It assumes all query terms have equal importance.
  - b) It requires precomputed similarity measures.
  - c) It cannot handle unseen documents or queries.
  - d) It assumes document lengths are always uniform.
74. Which retrieval model explicitly incorporates TF (Term Frequency) and IDF (Inverse Document Frequency) in its scoring?
- a) Vector Space Model with pivoted length normalization.
  - b) Query Likelihood Retrieval Model.
  - c) Boolean Retrieval Model.
  - d) Language Modeling with Dirichlet Smoothing.
75. In the context of vector space models, why is document length normalization crucial?
- a) To ensure that rare terms are penalized in longer documents.
  - b) To equalize the likelihood of term matches across documents of varying lengths.
  - c) To reward longer documents for their higher information density.
  - d) To prioritize documents containing all query terms regardless of length.
76. In the Vector Space Model, a document is represented as:
- a) A sequence of terms and phrases in hierarchical order.
  - b) A graph of term dependencies and their weights.
  - c) A high-dimensional vector where each dimension corresponds to a term.
  - d) A matrix of term frequencies and collection frequencies.
77. Why does the vector space model use a dot product to compute similarity between a document and a query?
- a) It emphasizes the rare terms that occur across all documents.
  - b) It measures the degree of overlap in the term weights of both vectors.
  - c) It normalizes the effect of document length automatically.
  - d) It reduces computational complexity by ignoring term frequency.
78. What challenge arises from representing documents using a bag-of-words approach in vector space models?
- a) The inability to account for term co-occurrence within documents.
  - b) Excessive reliance on document length normalization.
  - c) Overemphasis on stop words due to their high IDF.
  - d) Loss of term frequency information for longer documents.
79. How does TF-IDF weighting address the imbalance caused by frequently occurring terms like "the" or "and"?
- a) It assigns a weight proportional to the absolute term frequency in a document.
  - b) It penalizes common terms by assigning them a lower IDF score.
  - c) It increases the document length to balance the frequency.
  - d) It excludes stop words entirely from the vocabulary.
80. Which of the following represents a key improvement of the pivoted length normalization method over standard TF-IDF?
- a) It uses sublinear transformation to adjust term frequency scores dynamically.
  - b) It applies a binary weighting system for stop words.
  - c) It adjusts the TF-IDF weight based on a document's deviation from the average length.
  - d) It excludes high-frequency terms from the document corpus entirely.
81. If two documents  $d_1$  and  $d_2$  have identical term frequencies but differ in length, which of the following will most likely occur in a well-normalized Vector Space Model?
- a) Both documents will be ranked equally for any query.
  - b)  $d_1$  will be ranked higher if it is shorter.
  - c)  $d_2$  will be ranked higher if it is longer.
  - d) Neither document will be ranked for queries containing rare terms.
82. In TF-IDF weighting, which term would likely have the highest contribution to document relevance?
- a) A common term appearing in almost all documents.
  - b) A moderately rare term with high frequency in a specific document.
  - c) A term occurring equally across all documents.
  - d) A term absent from the query vector.
83. The introduction of sublinear term frequency transformations (e.g., logarithmic scaling) in vector space models primarily aims to:
- a) Prevent excessive penalization of infrequent terms.
  - b) Adjust term weight contributions in long documents.
  - c) Avoid dominance by terms that appear excessively within a single document.
  - d) Reduce the computational complexity of similarity calculations.
84. How does document frequency (DF) influence the weight of a term in the TF-IDF model?
- a) Higher DF values increase the importance of a term.
  - b) Lower DF values lead to higher inverse document frequency (IDF) scores.
  - c) DF directly determines the length of a document vector.
  - d) DF scores are used to normalize term weights across queries.
85. What fundamental assumption underlies the use of cosine similarity in vector space models?
- a) The vectors are normalized to unit length to account for document length differences.
  - b) The angles between vectors are irrelevant in determining similarity.
  - c) Only binary term presence contributes to similarity measures.
  - d) All terms are weighted equally regardless of their importance.

86. What problem does inverse document frequency (IDF) aim to solve in retrieval models?
  - a) Overemphasis on document length in similarity calculations.
  - b) Dominance of common terms that appear in most documents.
  - c) Insufficient weighting of query terms in long documents.
  - d) Redundancy caused by repeated terms in queries.
87. Which of the following best illustrates the issue with using raw term frequency (TF) without transformation in vector space models?
  - a) Rare terms are assigned disproportionately high weights.
  - b) Documents containing repetitive text may dominate rankings.
  - c) Query words that do not occur in any document are assigned a weight of zero.
  - d) Stop words are excluded automatically from the term vector.
88. How does the BM25 model address the diminishing returns of repeated term occurrences in a document?
  - a) By applying a logarithmic scaling to the term frequency
  - b) By normalizing term frequency using document length
  - c) By using an upper-bounded term frequency transformation function
  - d) By excluding terms that occur more than a fixed number of times
89. Why might two documents with the same term frequencies receive different relevance scores in a retrieval model?
  - a) One document contains stop words that are ignored in ranking
  - b) Differences in their lengths affect the normalization factor
  - c) Their similarity is measured with respect to different queries
  - d) Their term distributions across the corpus vary significantly
90. The key conceptual difference between vector space and probabilistic retrieval models is:
  - a) Vector space models rely on cosine similarity, while probabilistic models do not
  - b) Probabilistic models rank documents based on the likelihood of relevance, not similarity
  - c) Vector space models ignore document frequency, while probabilistic models depend on it
  - d) Probabilistic models always require training data, whereas vector space models do not
91. In probabilistic retrieval models, the relevance of a document is expressed as:
  - a) A function of the similarity score between the document and query vectors
  - b) A binary decision based on document length normalization
  - c) A conditional probability given the document and query
  - d) A sum of term weights across all matched terms
92. What assumption underlies the Query Likelihood Model in probabilistic retrieval?
  - a) Documents are equally likely to be relevant
  - b) Queries are generated as samples from the language model of a document
  - c) Users select documents based solely on cosine similarity
  - d) Document length has no influence on term probabilities
93. What issue arises in the Query Likelihood Model if a query term does not appear in a document?
  - a) The document receives an overly high probability score
  - b) The document cannot be ranked due to zero probability
  - c) The term is ignored entirely in the scoring process
  - d) The term contributes a penalty instead of a reward
94. How does smoothing address the zero-probability problem in the Query Likelihood Model?
  - a) By adding a small, non-zero probability to all terms in the document language model
  - b) By excluding rare terms from the calculation entirely
  - c) By normalizing term frequencies across the entire corpus
  - d) By increasing the weight of terms that match query words
95. Which smoothing technique adjusts term probabilities using the collection-wide frequency of words?
  - a) Laplace Smoothing
  - b) Dirichlet Prior Smoothing
  - c) Pivoted Length Normalization
  - d) Term Frequency Smoothing
96. What is the primary role of smoothing in probabilistic retrieval models?
  - a) To increase the weight of frequently occurring terms
  - b) To account for unseen query terms in the document model
  - c) To reduce the length normalization bias
  - d) To penalize documents containing rare terms
97. How does the Dirichlet Prior Smoothing technique differ from Laplace Smoothing?
  - a) It applies a fixed probability to all unseen terms
  - b) It considers both document-specific and collection-wide term frequencies
  - c) It eliminates the need for term frequency normalization
  - d) It assigns equal probabilities to all query terms, regardless of their frequency
98. In the Query Likelihood Model, a document will have a higher score if:
  - a) It contains more query terms with high collection frequency
  - b) Its language model assigns high probabilities to the query terms
  - c) Its length is significantly above the collection average
  - d) It minimizes the number of rare terms included in its content
99. What key limitation does smoothing aim to address in the basic Query Likelihood Model?
  - a) Overemphasis on document length normalization

- b) Underrepresentation of stop words in ranking functions
  - c) Zero probabilities assigned to documents lacking query terms
  - d) Difficulty in ranking documents with high TF-IDF scores
100. How does the concept of a "document language model" contribute to retrieval in probabilistic models?
- a) It represents the document as a distribution over terms for query generation
  - b) It transforms the query vector into a probabilistic form
  - c) It penalizes long documents based on term frequency
  - d) It assigns a static weight to each document in the collection
101. Why is the independence assumption for query terms considered a simplification in probabilistic retrieval models?
- a) It allows term weights to be directly proportional to document length
  - b) It ignores term dependencies but simplifies the computation of probabilities
  - c) It ensures that all terms are treated equally regardless of frequency
  - d) It enables documents to be represented using binary vectors only
102. In language modeling for retrieval, the probability of a query given a document is calculated as:
- a) A sum of term frequencies across all query terms in the document
  - b) A product of probabilities for each query term, conditioned on the document
  - c) A ratio of document length to average document length in the corpus
  - d) A difference between query term frequency and document frequency
103. Which of the following is a key advantage of the Query Likelihood Model?
- a) It directly incorporates cosine similarity for ranking
  - b) It avoids explicit term weighting by relying on term probabilities
  - c) It eliminates the need for document length normalization
  - d) It penalizes documents containing rare terms
104. What factor primarily distinguishes the BM25 model from classical probabilistic models?
- a) The use of a binary relevance assumption
  - b) The incorporation of term frequency saturation and length normalization
  - c) Its reliance on global term distributions rather than document-specific statistics
  - d) The absence of inverse document frequency in scoring
105. How does smoothing with the Dirichlet prior achieve balance between document-specific and collection-wide statistics?
- a) By adding a fixed value to the term frequency for every term in the document
  - b) By adjusting term probabilities based on a prior estimate from the entire corpus
  - c) By normalizing document length against the average document length
  - d) By assigning equal weights to all terms regardless of their frequency
106. Which of the following best explains why BM25 is widely adopted in search engine implementations?
- a) Its complexity allows fine-grained term analysis
  - b) It combines TF, IDF, and length normalization effectively into a flexible scoring function
  - c) It avoids the use of smoothing for unseen terms
  - d) It directly models the semantic meaning of query terms
107. What issue does pivoted length normalization attempt to address that is not handled by basic TF-IDF weighting?
- a) Bias towards documents containing rare terms
  - b) Over-penalization of long documents during ranking
  - c) Underrepresentation of common terms in short documents
  - d) Zero probability of relevance for documents with unmatched query terms
108. Why might BM25+ outperform the original BM25 in certain retrieval scenarios?
- a) It excludes document length normalization entirely
  - b) It includes a small constant to reduce over-penalization of longer documents
  - c) It prioritizes stop words to improve matching for natural language queries
  - d) It replaces IDF with term co-occurrence probabilities
109. How does the concept of "term frequency saturation" improve retrieval functions like BM25?
- a) It ensures a linear increase in relevance scores with term repetition
  - b) It caps the contribution of a term to avoid dominance by excessive occurrences
  - c) It balances the importance of rare and frequent terms uniformly
  - d) It adjusts the weight of terms based on their positional importance
110. In a retrieval scenario using a language modeling approach, why might a query fail to retrieve relevant documents without smoothing?
- a) The query terms might be assigned excessively high probabilities
  - b) The language model does not differentiate between rare and common terms
  - c) The probability of generating unseen query terms is assumed to be zero
  - d) The document length normalization factor is omitted from the calculation
111. What is a common critique of bag-of-words representations in information retrieval?
- a) They overemphasize rare terms by default
  - b) They fail to capture the semantic relationships between terms
  - c) They are computationally expensive compared to probabilistic models
  - d) They require smoothing to handle missing query terms
112. How does the introduction of IDF weighting resolve a fundamental problem in the vector space model?
- a) It ensures that longer documents are ranked higher

- b) It adjusts the weight of terms based on their specificity within the collection
  - c) It prioritizes documents with higher term frequency across the entire corpus
  - d) It guarantees that stop words are completely ignored in retrieval
113. What is the theoretical basis for combining term frequency (TF) and inverse document frequency (IDF) in retrieval models?
- a) To reward documents with balanced lengths and moderate term frequencies
  - b) To balance the importance of query terms based on their relevance and rarity
  - c) To reduce computational complexity in similarity calculations
  - d) To equalize the distribution of term weights across the query and document vectors
114. Which of the following scenarios is most likely to benefit from the use of BM25-F?
- a) Retrieving documents based on structured fields like title and abstract
  - b) Ranking unstructured text documents by cosine similarity
  - c) Identifying semantic relationships in queries with synonyms
  - d) Scoring documents without smoothing for unseen terms
115. Why is the choice of smoothing parameter critical in the Query Likelihood Model?
- a) It controls the weight of rare terms in long documents
  - b) It determines the influence of collection-wide statistics on document scores
  - c) It ensures zero probabilities are assigned to irrelevant documents
  - d) It balances the TF-IDF scores across different document fields
116. What is the primary task of a tokenizer in an IR system?
- a) Compressing data
  - b) Splitting documents into countable features or tokens
  - c) Assigning weights to tokens
  - d) Sorting tokens by frequency
117. Which of the following represents a common tokenization strategy?
- a) Stop word normalization
  - b) Forward indexing
  - c) Whitespace-based tokenization
  - d) Term-at-a-time ranking
118. Why is it efficient to use term IDs instead of string terms in tokenization?
- a) Term IDs are easier to generate
  - b) Term IDs save memory and allow  $O(1)$  lookups
  - c) Term IDs enable stop-word filtering
  - d) Term IDs increase document count accuracy
119. In tokenization, what does the term "feature generation" imply?
- a) Generating unique document identifiers
  - b) Defining the building blocks of document objects
  - c) Calculating term frequency weights
  - d) Filtering irrelevant documents
120. How does a whitespace tokenizer treat the sentence "Data is key"?
- a) Data: 1, is: 1, key: 1
  - b) Data: 1, Key: 1
  - c) data: 1, is: 1, key: 1
  - d) Data: 1, Is: 1, Key: 1
121. Which is NOT a key feature of the inverted index?
- a) Lexicon
  - b) Forward index
  - c) Postings file
  - d) Document frequency
122. What is the primary role of the lexicon in an inverted index?
- a) Mapping documents to terms
  - b) Storing offsets to postings file entries
  - c) Calculating term weights
  - d) Compressing document IDs
123. What is stored in the postings file of an inverted index?
- a) Term frequency scores only
  - b) Positions of terms within documents
  - c) Document metadata
  - d) All documents with zero scores
124. In indexing, what is the significance of merging runs?
- a) To improve tokenization efficiency
  - b) To create a single sorted postings file
  - c) To calculate document weights
  - d) To map term IDs to lexicon entries
125. What happens during the creation of a forward index?
- a) Terms are mapped to documents
  - b) Documents are mapped to term lists
  - c) Lexicons are compressed for fast retrieval
  - d) Token frequencies are smoothed
126. Why is term-at-a-time ranking preferred in many IR systems?
- a) It scores only documents with non-zero scores
  - b) It requires no accumulators
  - c) It scores all documents regardless of relevance
  - d) It avoids stop-word removal
127. Which scoring algorithm uses a priority queue to maintain top-k documents?
- a) Term-at-a-time
  - b) Document-at-a-time
  - c) BM25 scoring
  - d) Proximity ranking
128. How does a scorer utilize the inverted index?
- a) By retrieving document metadata
  - b) By scanning the forward index
  - c) By scoring only documents with query terms
  - d) By reducing storage requirements
129. What is the purpose of filtering in document ranking?



- a) To update score accumulators
  - b) To skip irrelevant documents based on metadata
  - c) To reduce the lexicon size
  - d) To prioritize proximity-based matches
130. In scoring algorithms, what is typically used to measure term relevance?
- a) Term frequency
  - b) Term position
  - c) Document length
  - d) IDF weighting
131. Why is compression critical in search engine implementations?
- a) To improve tokenization accuracy
  - b) To reduce memory usage and increase read efficiency
  - c) To simplify query processing
  - d) To support index sharding
132. How does caching improve query performance?
- a) By precomputing scores for all documents
  - b) By storing frequently accessed postings in memory
  - c) By reducing the size of the lexicon
  - d) By removing irrelevant documents
133. Which caching strategy follows Zipf's law?
- a) Term-at-a-time caching
  - b) Document metadata caching
  - c) Least Recently Used (LRU)
  - d) Priority-based caching
134. What is index sharding?
- a) Dividing the index into smaller sections
  - b) Merging multiple postings files
  - c) Storing lexicons in memory
  - d) Compressing document IDs
135. How is a shard typically assigned in a distributed search engine?
- a) By document frequency
  - b) By term positions
  - c) By the number of nodes or threads
  - d) By lexicon size
136. What is the main advantage of using a postings file in an inverted index?
- a) It improves the storage format for metadata
  - b) It enables quick lookup of documents containing specific terms
  - c) It reduces the need for document clustering
  - d) It simplifies query tokenization
137. Why is it essential for the postings file to store term positions within documents?
- a) To track the term frequency across the corpus
  - b) To implement proximity-based heuristics for phrase matching
  - c) To remove irrelevant terms during indexing
  - d) To simplify lexicon creation
138. Which of the following is a critical design consideration for the forward index?
- a) Mapping terms to their positions in the postings file
  - b) Allowing efficient access to all terms in a specific document
  - c) Reducing the size of the lexicon stored in memory
  - d) Supporting stop-word removal during scoring
139. How does index sharding optimize query performance?
- a) By ensuring all terms are stored in a single index
  - b) By parallelizing search tasks across multiple shards
  - c) By merging postings files during query processing
  - d) By caching frequently accessed terms
140. Which step in the indexing process ensures efficient term lookup for queries?
- a) Assigning unique term IDs during tokenization
  - b) Sorting terms by frequency in the lexicon
  - c) Merging intermediate runs into a single postings file
  - d) Removing duplicates from the forward index
141. Which scoring model is most associated with term frequency-inverse document frequency (TF-IDF)?
- a) Okapi BM25
  - b) Vector Space Model
  - c) Boolean Retrieval Model
  - d) Proximity-Based Ranking
142. What is the key difference between term-at-a-time and document-at-a-time ranking?
- a) The use of lexicons in scoring
  - b) The method of updating score accumulators
  - c) The requirement for compression in lexicon files
  - d) The inclusion of metadata in ranking scores
143. Why is it unnecessary to score documents with zero-term matches during ranking?
- a) Their scores will always be zero
  - b) They contain no relevant metadata
  - c) They cannot be cached effectively
  - d) They are not stored in the forward index
144. In document-at-a-time ranking, what data structure is used to maintain the top-k documents?
- a) Linked list
  - b) Priority queue
  - c) Hash map
  - d) Bloom filter
145. Which ranking strategy minimizes memory usage for score accumulators?
- a) Term-at-a-time ranking
  - b) Document-at-a-time ranking
  - c) Filtering-based ranking
  - d) BM25 with priority queue
146. Why is the lexicon typically stored in memory?
- a) To allow fast lookup of term-to-document mappings
  - b) To ensure all query terms are scored
  - c) To avoid compression overhead during scoring
  - d) To cache frequently used documents
147. How does compression of the postings file improve query performance?

- a) By reducing the size of the lexicon
  - b) By decreasing disk seek times for term lookups
  - c) By optimizing document scoring algorithms
  - d) By caching term IDs in memory
148. Which strategy ensures that the cache stores only the most useful terms?
- a) Randomized eviction
  - b) Least Recently Used (LRU) policy
  - c) Term frequency normalization
  - d) Proximity-based filtering
149. What is the role of filters in the scoring process?
- a) To prioritize documents with higher term frequencies
  - b) To exclude documents that do not meet specific criteria
  - c) To compress the postings file for faster lookups
  - d) To limit the number of terms in the query
150. What is the purpose of a priority queue in query processing?
- a) To store all document scores temporarily
  - b) To maintain a list of the k most relevant documents
  - c) To ensure efficient term lookups in the lexicon
  - d) To compress term positions within documents
151. In the context of search engines, what does "relevance feedback" typically involve?
- a) Adjusting term weights based on user interaction data
  - b) Compressing the lexicon for faster retrieval
  - c) Removing irrelevant documents from the index
  - d) Updating stop-word lists for tokenization
152. Why is document filtering commonly applied before scoring in IR systems?
- a) To reduce the size of the lexicon
  - b) To exclude documents that do not satisfy user-defined constraints
  - c) To prioritize terms with high frequency
  - d) To enable caching of frequently accessed documents
153. Which method is commonly used to merge multiple partial inverted indexes during indexing?
- a) Binary search
  - b) Merge sort
  - c) Quick sort
  - d) Hash mapping
154. What is the primary advantage of using document meta-data in ranking?
- a) It allows faster compression of postings files
  - b) It enables advanced filters such as date or document type constraints
  - c) It reduces the need for proximity heuristics
  - d) It simplifies the process of assigning term IDs
155. Which ranking method optimizes query time by processing documents in their shard locations?
- a) Document-at-a-time ranking
  - b) Term-at-a-time ranking
  - c) Distributed ranking with index sharding
  - d) BM25 ranking with LRU caching
156. What is the primary purpose of search engine evaluation?
- a) To design new algorithms
  - b) To compare retrieval methods and assess their utility
  - c) To increase computational efficiency
  - d) To refine user query formulations
157. What are the three main dimensions of search engine evaluation?
- a) Precision, Recall, Relevance
  - b) Effectiveness, Usability, Efficiency
  - c) Accuracy, Reliability, Scalability
  - d) Ranking, Indexing, Filtering
158. Which evaluation methodology is central to modern search engine testing?
- a) A-B testing
  - b) Pooling
  - c) Cranfield evaluation methodology
  - d) Mean Reciprocal Rank testing
159. In the Cranfield methodology, what is a critical component for fair comparison?
- a) A unique query set for each method
  - b) Relevance judgments consistent across systems
  - c) Randomized query sampling
  - d) The inclusion of user feedback during testing
160. How is relevance defined in binary judgments?
- a) By the frequency of terms in a document
  - b) As either relevant or non-relevant for a specific query
  - c) By the user's time spent on a document
  - d) As marginally or highly relevant
161. What is the main limitation of using binary relevance judgments?
- a) They are computationally expensive
  - b) They oversimplify the relevance spectrum
  - c) They do not account for system efficiency
  - d) They cannot be reused across queries
162. Which of the following is NOT a benefit of Cranfield methodology?
- a) Reusability of test collections
  - b) Objectivity in algorithm comparisons
  - c) Real-time user interaction during testing
  - d) Compatibility with various retrieval systems
163. In the context of evaluation, what are relevance judgments?
- a) Binary ratings assigned by system administrators
  - b) User-assigned labels indicating document utility
  - c) Automatically generated labels from test collections
  - d) Annotations based on computational heuristics
164. Why is user involvement critical in search engine evaluation?
- a) To determine query efficiency
  - b) To establish relevance judgments
  - c) To perform system ranking
  - d) To compute MAP scores
165. What distinguishes search engine evaluation from database evaluation?

- a) The need for computational efficiency
  - b) The emphasis on user feedback and ranking
  - c) The reliance on binary classification
  - d) The focus on query optimization
166. Precision measures...
- a) The proportion of retrieved documents that are relevant
  - b) The number of relevant documents retrieved out of the total
  - c) The time taken to retrieve documents
  - d) The average relevance score of retrieved documents
167. Recall measures...
- a) The total number of documents retrieved
  - b) The proportion of relevant documents retrieved
  - c) The completeness of the query processing
  - d) The ranking accuracy of the system
168. What is the ideal precision and recall value for a system?
- a) 1.0 for precision, 0.5 for recall
  - b) Both should be close to 0.8
  - c) Both should be 1.0
  - d) 0.5 for precision, 1.0 for recall
169. The F1 score combines...
- a) Precision and ranking
  - b) Recall and ranking
  - c) Precision and recall
  - d) Precision and user interaction
170. Why is F1 score preferred over arithmetic mean of precision and recall?
- a) It emphasizes high recall values
  - b) It balances precision and recall effectively
  - c) It simplifies evaluation metrics
  - d) It penalizes algorithms with low recall
171. What is the formula for F1 score?
- a)  $F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
  - b)  $F1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
  - c)  $F1 = \frac{\text{Precision} \cdot \text{Recall}}{2}$
  - d)  $F1 = \frac{\text{Recall}}{\text{Precision}}$
172. Why is high recall often associated with low precision?
- a) Systems retrieve too many irrelevant documents
  - b) User interactions reduce the ranking accuracy
  - c) Systems cannot filter non-relevant documents efficiently
  - d) Query terms increase the lexicon size
173. What does precision@k represent?
- a) The percentage of queries answered correctly within k seconds
  - b) The proportion of relevant documents in the top-k results
  - c) The average precision for k queries
  - d) The relevance score weighted by k factors
174. Why is F1 score considered a harmonic mean?
- a) It averages precision and recall without bias
  - b) It smoothens variance in precision-recall tradeoff
  - c) It penalizes extreme values of precision or recall
  - d) It prioritizes efficiency over accuracy
175. A system with perfect recall but low precision would likely...
- a) Return only a subset of relevant documents
  - b) Retrieve all relevant and many irrelevant documents
  - c) Exclude some relevant documents
  - d) Return documents in random order
176. Why is it important to evaluate ranked lists in search engine evaluation?
- a) To ensure all documents are retrieved
  - b) To consider the position of relevant documents
  - c) To remove irrelevant documents from the corpus
  - d) To optimize query speed
177. How is precision-recall (PR) curve generated?
- a) By plotting precision at every rank of the retrieved list
  - b) By averaging the number of retrieved documents
  - c) By calculating the F1 score across all queries
  - d) By plotting the ranking of irrelevant documents
178. What does MAP (Mean Average Precision) measure?
- a) Average precision across multiple queries
  - b) Average recall across all ranked lists
  - c) Precision for the top k documents
  - d) Total number of relevant documents in the corpus
179. Why is MAP preferred over simple precision for comparing algorithms?
- a) It considers the rank of every relevant document
  - b) It penalizes irrelevant documents in the corpus
  - c) It ignores user preferences
  - d) It is computationally easier to calculate
180. What is the significance of a high MAP score?
- a) It indicates the system returns fewer documents
  - b) It means relevant documents are consistently ranked high
  - c) It ensures that recall is always greater than precision
  - d) It represents fewer false positives in ranking
181. In MAP calculation, what does the denominator in the precision formula represent?
- a) Total number of retrieved documents
  - b) Total number of relevant documents in the collection
  - c) Number of relevant documents retrieved so far
  - d) Total number of documents in the query
182. How does gMAP differ from MAP?
- a) gMAP emphasizes easy queries over difficult ones
  - b) gMAP penalizes systems with low performance on certain queries
  - c) gMAP averages precision over fewer queries
  - d) gMAP does not consider ranking positions
183. What is the reciprocal rank in known-item search?
- a) The inverse of the rank of the first relevant document
  - b) The average rank of all retrieved documents

- c) The position of the least relevant document
  - d) The harmonic mean of all ranks
184. What is the Mean Reciprocal Rank (MRR) used for?
- a) Evaluating binary relevance judgments
  - b) Comparing systems with high MAP scores
  - c) Measuring efficiency in multi-level judgments
  - d) Assessing performance in single-relevant-item tasks
185. How does average precision differ from precision@k?
- a) Average precision considers all retrieved documents, while precision@k focuses on the top k results
  - b) Precision@k evaluates ranking position, but average precision does not
  - c) Average precision is used for multi-level judgments, while precision@k is for binary judgments
  - d) They are calculated identically, but for different queries
186. What is the purpose of Normalized Discounted Cumulative Gain (NDCG)?
- a) To measure system efficiency in user studies
  - b) To evaluate ranked lists with multiple levels of relevance
  - c) To optimize retrieval for binary relevance judgments
  - d) To assess recall in large datasets
187. In NDCG, why are gains discounted by their rank?
- a) To emphasize higher-ranked documents more than lower-ranked ones
  - b) To reduce computational complexity
  - c) To adjust for binary relevance scores
  - d) To improve retrieval speed
188. How is the ideal DCG (IDCG) defined for a query?
- a) The DCG of a randomly ordered ranked list
  - b) The DCG of a perfectly ranked list of relevant documents
  - c) The total gain of the top k documents in any order
  - d) The harmonic mean of precision and recall scores
189. What does a higher NDCG score indicate?
- a) The system retrieves fewer irrelevant documents
  - b) Relevant documents are ranked closer to the top
  - c) The system uses fewer computational resources
  - d) Precision increases at the expense of recall
190. What is the primary benefit of NDCG over MAP?
- a) It handles non-binary relevance judgments
  - b) It simplifies ranking calculations
  - c) It focuses only on precision@10
  - d) It measures efficiency instead of accuracy
191. How are relevance levels used in calculating DCG?
- a) They are multiplied by their rank
  - b) They are discounted logarithmically by their position
  - c) They are converted to binary values before summation
  - d) They are averaged over all queries
192. Why is normalization important in NDCG?
- a) To scale scores between 0 and 1 for comparison across queries
  - b) To remove bias caused by query length
  - c) To adjust for variations in corpus size
  - d) To simplify multi-level judgments
193. What is a potential drawback of NDCG?
- a) It is insensitive to document order
  - b) It cannot be applied to binary relevance judgments
  - c) It relies heavily on manual relevance judgments
  - d) It ignores the top k documents in the ranking
194. What is the role of pooling in test collection evaluation?
- a) To evaluate multi-level judgments more efficiently
  - b) To combine top-k results from multiple systems for judgment
  - c) To automate the labeling of non-relevant documents
  - d) To speed up relevance judgment calculations
195. Why are unjudged documents in the pool assumed to be non-relevant?
- a) To reduce human effort in relevance judgment
  - b) To simplify ranking algorithms
  - c) To minimize computational overhead
  - d) To improve the F1 score
196. What is the primary advantage of dense word embeddings over one-hot encoding?
- a) They improve the efficiency of hardware computation
  - b) They capture semantic relationships between words
  - c) They require fewer parameters to train neural networks
  - d) They eliminate the need for labeled datasets
197. What hypothesis underpins the development of word embeddings?
- a) Latent Semantic Hypothesis
  - b) Distributional Hypothesis
  - c) Frequency Hypothesis
  - d) Neural Transformation Hypothesis
198. How do classical vectorial representations represent words in text?
- a) As dense embeddings trained using neural networks
  - b) As binary or weighted vectors in a term-document matrix
  - c) As sequences of characters with frequency counts
  - d) As graph-based nodes connected by edges
199. What is a major limitation of one-hot encoding?
- a) It requires pre-trained embeddings
  - b) It does not encode semantic similarity between words
  - c) It cannot represent sparse data efficiently
  - d) It cannot handle high-dimensional vectors
200. What is the "curse of dimensionality" in one-hot encoding?
- a) Data sparsity caused by high-dimensional vectors
  - b) The inability to train deep learning models efficiently
  - c) Overfitting caused by high variance in data
  - d) Lack of representation for unknown words
201. What mathematical operation is often used to measure similarity in classical vectorial models?
- a) Dot product

- b) Euclidean distance
  - c) Cosine similarity
  - d) Matrix multiplication
202. How do word embeddings improve upon classical models?
- a) By using supervised learning on labeled datasets
  - b) By leveraging co-occurrence statistics in dense spaces
  - c) By optimizing sparse matrix multiplications
  - d) By replacing one-hot encoding with binary representations
203. Which term describes embeddings that capture both syntactic and semantic regularities?
- a) Feature-based embeddings
  - b) Fixed-length distributed representations
  - c) Probabilistic latent representations
  - d) Sparse vector embeddings
204. In word embeddings, what does the equation "king + woman - man = queen" illustrate?
- a) One-hot encoding
  - b) Distributional semantics
  - c) Arithmetic on embedding vectors
  - d) Semantic anomaly detection
205. What problem does FastText address that Word2Vec struggles with?
- a) Handling large vocabularies
  - b) Understanding subword-level semantics
  - c) Reducing training time for embeddings
  - d) Capturing context-sensitive representations
206. Which model introduced neural network-based embeddings?
- a) Skip-Gram
  - b) CBOW
  - c) NNLM (Bengio et al. 2003)
  - d) Glove
207. What is the main difference between Skip-Gram and CBOW in Word2Vec?
- a) Skip-Gram predicts the center word from context words, while CBOW does the reverse
  - b) CBOW uses large corpora, while Skip-Gram uses small datasets
  - c) Skip-Gram is faster to train than CBOW
  - d) CBOW captures word positions, while Skip-Gram ignores them
208. What optimization method is commonly used in Word2Vec to improve efficiency?
- a) Hierarchical Softmax
  - b) Negative Sampling
  - c) Both a and b
  - d) Backpropagation
209. How does FastText improve over Word2Vec for rare words?
- a) By encoding characters instead of words
  - b) By using n-gram embeddings
  - c) By normalizing co-occurrence counts
  - d) By extending the context window
210. What distinguishes ELMo embeddings from Word2Vec?
- a) ELMo produces context-dependent embeddings
  - b) ELMo uses only CBOW architecture
  - c) ELMo ignores sentence-level context
  - d) ELMo embeddings are static
211. What is the defining feature of Transformer-based models like BERT?
- a) Attention mechanisms
  - b) N-gram modeling
  - c) Sparse embeddings
  - d) Bi-directional RNNs
212. In Transformers, what does "self-attention" achieve?
- a) It computes relationships between words in the context
  - b) It reduces training complexity in embeddings
  - c) It eliminates the need for tokenization
  - d) It focuses on sentence-level classification
213. What type of tasks can BERT embeddings handle well?
- a) Sequence generation tasks
  - b) Semantic search and information retrieval
  - c) Machine translation
  - d) Graph-based learning
214. How does GPT differ fundamentally from BERT?
- a) GPT uses only encoders, while BERT uses decoders
  - b) GPT is unidirectional, while BERT is bidirectional
  - c) GPT uses fine-tuning, while BERT does not
  - d) GPT models context through n-grams
215. Which model uses both encoder and decoder structures?
- a) GPT
  - b) BERT
  - c) T5
  - d) CBOW
216. What innovation is central to the GloVe model?
- a) Using word-context matrices for embeddings
  - b) Ratios of word co-occurrence probabilities
  - c) Replacing Softmax with hierarchical layers
  - d) Applying subword embeddings for rare words
217. What is the purpose of masked language modeling in BERT?
- a) To predict missing words in a sentence
  - b) To calculate co-occurrence probabilities
  - c) To encode positional relationships between words
  - d) To reduce training time
218. What challenge is addressed by subword-level embeddings in FastText?
- a) Handling out-of-vocabulary (OOV) words
  - b) Increasing context window size
  - c) Learning embeddings faster
  - d) Improving batch training
219. What is a key feature of contextual embeddings like ELMo?
- a) They assign the same embedding to a word regardless of its context

- b) They adapt embeddings based on the word's context in a sentence
  - c) They use binary vectors for efficiency
  - d) They rely on fixed-length dense vectors
220. Why are Transformers considered cutting-edge in NLP?
- a) They rely on n-gram models
  - b) They eliminate sequential processing with attention mechanisms
  - c) They focus only on text classification tasks
  - d) They require no pre-training
221. How does T5 (Text-to-Text Transfer Transformer) differ from BERT and GPT?
- a) It uses an encoder-only architecture
  - b) It treats every NLP task as a sequence-to-sequence problem
  - c) It relies on character-level embeddings for better accuracy
  - d) It generates embeddings without training
222. What is the main limitation of Word2Vec with respect to out-of-vocabulary (OOV) words?
- a) It cannot generate embeddings for words not in the training corpus
  - b) It requires significant computational resources
  - c) It produces context-insensitive embeddings
  - d) It fails to capture semantic relationships
223. In FastText, how are embeddings generated for OOV words?
- a) By averaging the embeddings of similar words
  - b) By constructing embeddings from n-grams of the word
  - c) By using a random initialization followed by fine-tuning
  - d) By ignoring OOV words entirely
224. What distinguishes context-specific embeddings from static embeddings?
- a) Context-specific embeddings adjust based on surrounding words
  - b) Static embeddings change with fine-tuning
  - c) Context-specific embeddings use binary representations
  - d) Static embeddings are computationally faster to train
225. What is a practical application of BERT in information retrieval?
- a) Generating summaries for long documents
  - b) Matching queries to relevant documents using semantic meaning
  - c) Predicting the next sentence in a paragraph
  - d) Translating documents between languages
226. What is a "masked token" in the context of BERT pre-training?
- a) A placeholder used to predict missing words during training
  - b) A token that represents rare words in the dataset
  - c) A token excluded from the training corpus
  - d) A token that indicates sentence boundaries
227. How do attention mechanisms improve NLP models?
- a) By reducing the computational cost of matrix operations
  - b) By focusing on relevant parts of the input sequence
  - c) By replacing embeddings with binary representations
  - d) By optimizing word similarity measures
228. What is the role of "self-attention" in Transformers?
- a) It eliminates dependencies on labeled data
  - b) It encodes relationships between words in the same input sequence
  - c) It maps embeddings to a hierarchical structure
  - d) It computes probabilities for masked tokens
229. Which property of embeddings enables analogical reasoning (e.g., "king + woman - man = queen")?
- a) High-dimensionality
  - b) Semantic regularities in vector space
  - c) Syntactic dependencies
  - d) Contextual alignment of embeddings
230. What is the main advantage of using GloVe over Word2Vec?
- a) GloVe combines global co-occurrence statistics with local context
  - b) GloVe requires less training data than Word2Vec
  - c) GloVe embeddings are better for subword-level tasks
  - d) GloVe supports unsupervised learning natively
231. How does ELMo produce context-dependent embeddings?
- a) By using bidirectional LSTM-based language models
  - b) By training on fixed-length n-grams
  - c) By employing hierarchical Softmax
  - d) By reducing dimensionality of one-hot vectors
232. Why is the Transformer architecture considered scalable?
- a) It uses fixed-size embeddings for all datasets
  - b) It processes input sequences in parallel using attention
  - c) It eliminates the need for recurrent connections
  - d) It is optimized for small-scale tasks
233. What is the purpose of positional encoding in Transformers?
- a) To assign unique IDs to tokens
  - b) To encode the order of words in a sequence
  - c) To train embeddings faster
  - d) To replace attention mechanisms
234. Which Transformer model is used primarily for text generation tasks?
- a) BERT
  - b) GPT
  - c) T5
  - d) FastText
235. What improvement does the T5 model offer over BERT and GPT?
- a) It uses both encoder and decoder structures for sequence-to-sequence tasks
  - b) It eliminates the need for fine-tuning
  - c) It reduces the training data required for embeddings
  - d) It focuses only on text classification tasks

236. What is a key challenge with embeddings when dealing with OOV words?
- They require labeled training data for embeddings
  - They cannot represent words not seen during training
  - They increase the dimensionality of embeddings
  - They reduce the semantic accuracy of embeddings
237. How does FastText address the OOV problem?
- By breaking words into subword units
  - By leveraging co-occurrence counts
  - By pre-training on larger corpora
  - By interpolating missing embeddings
238. What is a significant drawback of count-based models like LSA?
- They cannot capture semantic regularities
  - They lack scalability for large vocabularies
  - They ignore global co-occurrence statistics
  - They are computationally less efficient than prediction-based models
239. Why is fine-tuning necessary after pretraining models like BERT?
- To adapt the pre-trained model to specific downstream tasks
  - To reduce the overfitting of embeddings
  - To remove irrelevant data from the training corpus
  - To improve the unsupervised learning phase
240. Which of the following is an example of a downstream task for embeddings?
- Text classification
  - Language model pretraining
  - Vector space optimization
  - Corpus augmentation
241. What differentiates subword embeddings from word embeddings?
- Subword embeddings use smaller context windows
  - Subword embeddings can represent parts of words, such as prefixes or suffixes
  - Subword embeddings require labeled datasets
  - Subword embeddings ignore rare words in the corpus
242. What aspect of embedding vectors makes them suitable for clustering tasks?
- High sparsity
  - Semantic proximity in vector space
  - Uniform dimensionality across all tokens
  - Static nature of vectors
243. Why are attention mechanisms critical for machine translation?
- They provide context for each word in the input and output sequences
  - They replace embeddings with binary representations
  - They align the embeddings to word frequencies
  - They reduce computational overhead
244. How does the Glove model derive semantic information?
- By focusing on ratios of co-occurrences between word pairs
  - By using hierarchical Softmax for faster training
  - By splitting words into subword units
  - By reducing embeddings to binary vectors
245. What is a key feature of embeddings produced by neural network language models?
- They encode both syntactic and semantic properties of words
  - They eliminate the need for labeled datasets entirely
  - They require fewer dimensions than count-based models
  - They are static and context-independent
246. What is the primary challenge for web crawlers in handling dynamic content?
- Respecting the robots.txt file
  - Recognizing duplicate pages
  - Managing traps with infinite dynamically generated content
  - Handling different file types
247. What does the robots.txt file specify for a crawler?
- Pages that should be indexed
  - Pages that should not be crawled
  - Pages that need higher ranking
  - Pages that require frequent updates
248. Focused crawling is characterized by:
- Randomly selecting URLs from a seed set
  - Prioritizing pages related to a specific topic
  - Downloading all forum posts on the web
  - Crawling only static pages
249. What is the most common crawling strategy to balance server load?
- Depth-first search
  - Random crawling
  - Breadth-first search
  - Incremental crawling
250. How does a crawler detect hidden URLs?
- By parsing robots.txt
  - Through manual discovery
  - By analyzing unlinked content sources
  - By following user interactions
251. What innovation allows Google to manage files across large clusters?
- Distributed MapReduce
  - Google File System (GFS)
  - Hadoop Index Manager
  - Parallel Computing Framework
252. What is the purpose of the shuffle and sort step in MapReduce?
- To compress the input data
  - To distribute keys among reducers
  - To delete duplicate records
  - To reduce memory overhead
253. Which challenge is unique to web-scale indexing compared to traditional indexing?
- Creating inverted indexes

- b) Processing data in parallel across machines
  - c) Ranking documents by relevance
  - d) Storing indexes on single machines
254. What does the Map phase in MapReduce typically output in an inverted indexing task?
- a) Word-frequency pairs
  - b) Document-summary pairs
  - c) Tokenized sentences
  - d) Anchor-text descriptions
255. Why is it beneficial to use a distributed file system like GFS in indexing?
- a) Ensures data is stored in a single location
  - b) Provides robust fault tolerance
  - c) Accelerates file compression
  - d) Optimizes single-threaded processing
256. Which statement about PageRank is true?
- a) It only counts direct inlinks
  - b) It assumes every page has at least one pseudo-link
  - c) It ignores anchor text
  - d) It cannot be combined with other ranking algorithms
257. In the random surfer model, what is the significance of parameter  $\alpha$ ?
- a) It determines the likelihood of jumping to any random page
  - b) It controls the rate of damping in PageRank
  - c) It represents the average time spent on a page
  - d) It balances inlink and outlink weightage
258. HITS differs from PageRank because it:
- a) Focuses only on authority scores
  - b) Computes both hub and authority scores
  - c) Requires random jump probabilities
  - d) Ignores adjacency matrix multiplication
259. What role does anchor text play in link-based ranking?
- a) Improves ranking robustness
  - b) Provides additional descriptive context for target pages
  - c) Prevents spam links from influencing scores
  - d) Simplifies inverted index creation
260. How does PageRank handle dead-end pages (zero out-links)?
- a) Ignores them in computations
  - b) Smooths the matrix with a random jump probability
  - c) Deletes such pages from the graph
  - d) Relies solely on anchor text for ranking
261. What is the primary goal of learning to rank in search engines?
- a) Optimizing storage efficiency
  - b) Combining multiple features into a single ranking function
  - c) Maximizing crawl frequency
  - d) Creating anchor-text-based relevance scores
262. Logistic regression in ranking assumes:
- a) Relevance is a nonlinear combination of features
  - b) Feature weights are static
  - c) Features can be combined linearly
  - d) Ranking scores must be binary
263. Which feature might NOT be included in a learning-to-rank model?
- a) PageRank score
  - b) URL format
  - c) User query
  - d) Crawler delay
264. Why do advanced ranking models optimize metrics like MAP or NDCG?
- a) To handle noisy data
  - b) To align ranking with user satisfaction
  - c) To predict relevance without human feedback
  - d) To simplify training data preparation
265. How can personalization improve ranking in learning to rank?
- a) By weighting anchor text more heavily
  - b) By creating user-specific PageRank scores
  - c) By limiting the scope of crawled pages
  - d) By prioritizing static over dynamic pages
266. Vertical search engines are more effective because they:
- a) Use general scoring functions
  - b) Focus on a specialized domain
  - c) Avoid user feedback
  - d) Only index static pages
267. Lifelong learning in search engines enables:
- a) Static indexing and scoring
  - b) Continuous improvement based on user behavior
  - c) Prioritization of historical relevance
  - d) Manual tuning of ranking functions
268. Integrating search and recommendation systems can provide:
- a) Randomized ranking
  - b) Enhanced decision support
  - c) Spam-resistant indexing
  - d) Faster crawling
269. Intelligent systems aim to connect which three elements?
- a) User-Agent-Query
  - b) Data-User-Service
  - c) Link-Anchor-Hub
  - d) Query-Document-Feedback
270. What does Google's Knowledge Graph exemplify?
- a) Bag-of-words implementation
  - b) Large-scale semantic analysis
  - c) Link-based ranking
  - d) URL optimization
271. Why is incremental crawling beneficial?
- a) It avoids re-crawling unchanged pages
  - b) It increases the speed of initial indexing
  - c) It reduces the number of servers required
  - d) It focuses only on high-utility pages



272. A focused crawler is more efficient than a general crawler because it:
- Operates without a seed set
  - Avoids parsing robots.txt files
  - Targets a pre-specified topic
  - Ignores page duplication
273. Which of the following best describes parallel crawling?
- Analyzing one server at a time
  - Simultaneously fetching pages across multiple threads or machines
  - Restricting crawling to local networks
  - Alternating between depth-first and breadth-first crawling
274. What is a critical consideration for crawlers when interacting with a web server?
- Matching all metadata
  - Avoiding excessive requests to prevent overload
  - Parsing all HTML elements
  - Ignores embedded media files
275. What kind of information might a crawler prioritize when respecting the robots.txt file?
- URL patterns disallowed for crawling
  - Priority levels of indexed pages
  - External link frequency
  - Target page's PageRank
276. What is the purpose of MapReduce in indexing?
- To normalize input text data
  - To facilitate parallel processing of data across large clusters
  - To extract semantic information from pages
  - To minimize query time
277. What component ensures fault tolerance in the MapReduce framework?
- Load balancer
  - Reducer fallback mechanisms
  - Task re-execution on other servers
  - Enhanced scheduler prioritization
278. How does the Google File System (GFS) ensure reliability of stored data?
- Dynamic re-indexing
  - Replicating data chunks across multiple servers
  - Incremental page fetching
  - Regular checksum verifications
279. In MapReduce, what happens to the outputs of the Map function?
- They are directly added to the index
  - They are grouped and sorted by key for the Reduce function
  - They are analyzed for spam detection
  - They are used to update the PageRank
280. Why is it necessary to distribute an inverted index across machines?
- To handle the scale of web data efficiently
  - To ensure content relevance
  - To prioritize link-based ranking
  - To minimize downtime during retrieval
281. Which algorithm introduced the concept of authority and hub scores?
- PageRank
  - HITS
  - BM25
  - Personalized PageRank
282. What differentiates Personalized PageRank from standard PageRank?
- Incorporation of random jumps
  - Use of query-specific relevance
  - Focus on external links
  - Elimination of anchor text
283. How does the HITS algorithm define a good hub page?
- A page with many inlinks
  - A page that links to high-authority pages
  - A page that is frequently updated
  - A page that balances inlinks and outlinks
284. PageRank captures indirect citations by:
- Counting secondary and tertiary links recursively
  - Penalizing pages with too many outlinks
  - Assigning equal weight to all links
  - Ignores low-quality pages
285. What is the primary limitation of HITS compared to PageRank?
- HITS does not account for link quality
  - HITS is query-dependent
  - HITS cannot handle spam effectively
  - HITS disregards link structure in ranking
286. Why are multiple features combined in learning-to-rank models?
- To remove redundant data
  - To improve robustness and accuracy of ranking
  - To minimize query latency
  - To prioritize feedback from users
287. What is a potential drawback of regression-based learning-to-rank methods?
- Lack of scalability
  - Inability to optimize ranking metrics like NDCG directly
  - Dependence on PageRank
  - Ignores link-based features
288. Which metric is often used to evaluate ranking accuracy in training data?
- Precision at N
  - Mean Reciprocal Rank (MRR)
  - Normalized Discounted Cumulative Gain (NDCG)
  - ROC-AUC
289. What distinguishes logistic regression from linear regression in ranking tasks?
- Logistic regression maps scores to probabilities between 0 and 1

- b) Logistic regression ignores feature weights
  - c) Logistic regression requires no training data
  - d) Logistic regression is unsuitable for ranking
290. Advanced machine learning ranking algorithms aim to:
- a) Prioritize simplicity over accuracy
  - b) Directly optimize ranking metrics
  - c) Eliminate the need for feature extraction
  - d) Generate new queries for users
291. Vertical search engines differ from general search engines by:
- a) Using broader document corpora
  - b) Catering to specialized user groups or domains
  - c) Relying only on anchor text analysis
  - d) Avoiding machine learning-based ranking
292. Personalized search engines benefit users by:
- a) Ignoring search history for privacy
  - b) Adapting results based on individual preferences
  - c) Using static ranking algorithms
  - d) Favoring general queries over specific ones
293. What is the goal of integrating recommendation systems into search engines?
- a) To improve crawling efficiency
  - b) To provide users with task-specific support
  - c) To reduce server load
  - d) To minimize reliance on link-based ranking
294. The "Data-User-Service Triangle" emphasizes:
- a) Focus on vertical search engines
  - b) Linking data, users, and services for intelligent systems
  - c) Replacing crawling with predictive analytics
  - d) Standardizing all search algorithms
295. How does Google's Knowledge Graph improve search?
- a) By enhancing URL ranking
  - b) By introducing semantic representations of entities and relationships
  - c) By using traditional bag-of-words models
  - d) By simplifying web crawling
296. Lifelong learning in search engines relies on:
- a) Static indexing of historical data
  - b) Continuous improvement based on user interaction data
  - c) User-provided manual rankings
  - d) Eliminating relevance feedback
297. Intelligent systems aim to improve productivity by:
- a) Automating every step of a user's workflow
  - b) Combining user effort with machine efficiency
  - c) Replacing human input in decision-making
  - d) Ignoring contextual user behavior
298. Which future trend involves combining search, navigation, and recommendation?
- a) Dynamic indexing
  - b) Integrated information management
  - c) Anchor text analysis
  - d) Static PageRank optimization
299. How does interactive task support differ from traditional search?
- a) It focuses on crawling static pages
  - b) It combines human and machine collaboration for task completion
  - c) It eliminates the need for relevance feedback
  - d) It prioritizes speed over accuracy
300. Which feature best characterizes modern intelligent systems?
- a) Focus on broad generalization
  - b) Emphasis on semantic analysis and personalized interaction
  - c) Reliance on manual query processing
  - d) Avoidance of machine learning
301. What is the primary goal of text data access in comparison to text data analysis?
- a) To summarize large datasets
  - b) To focus on relevant data for further processing
  - c) To create predictive models from data
  - d) To emphasize sentiment analysis
302. Which scenario highlights the importance of advanced text analysis tools?
- a) When manual text processing is sufficient
  - b) When the dataset is small and static
  - c) For time-critical applications requiring fast decision-making
  - d) When search engines alone suffice for text retrieval
303. What is a key benefit of text analysis in scientific research?
- a) Automating document indexing
  - b) Integrating terminology across disciplines
  - c) Simplifying the extraction of word associations
  - d) Reducing the need for search engines
304. A disaster management system can benefit from real-time text analysis by:
- a) Creating predictive models for stock markets
  - b) Extracting topic clusters from social media discussions
  - c) Monitoring warning signs in tweets about potential natural disasters
  - d) Organizing data into hierarchical categories
305. Text mining is often referred to as a "datascope" because:
- a) It focuses solely on text summarization
  - b) It reveals hidden patterns and knowledge in large datasets
  - c) It exclusively analyzes opinions and sentiments
  - d) It automates the manual processes of text translation
306. Humans are described as "subjective sensors" because they:
- a) Collect numerical and relational data
  - b) Provide unbiased data about the real world
  - c) Express observations and opinions through text
  - d) Record multimedia data like video or audio

307. What advantage does treating text data as human sensor data provide?
- Enables integration with physical sensor data
  - Allows automated sentiment tagging
  - Simplifies clustering algorithms
  - Reduces ambiguity in natural language processing
308. In the data mining framework, text data is particularly valuable because it:
- Is easier to process than numerical data
  - Provides insight into user preferences and opinions
  - Requires fewer algorithms to analyze
  - Contains fewer dimensions than non-text data
309. A general text mining problem aims to:
- Create standardized datasets from text
  - Transform text into actionable knowledge for decision-making
  - Analyze text without using non-textual metadata
  - Predict trends based solely on numerical data
310. How can non-text data enhance text analysis?
- By replacing text-based clustering techniques
  - By providing context, such as time or location, for text interpretation
  - By reducing the need for feature extraction in text
  - By enabling real-time text translation
311. What is a key challenge when mining knowledge from text data?
- Extracting objective statements separately from subjective opinions
  - Representing data in relational databases
  - Identifying errors in physical sensor readings
  - Avoiding repetitive clustering tasks
312. Mining knowledge about the observer (text producer) involves:
- Extracting the main topic of the text
  - Predicting the author's sentiment or mood
  - Identifying the publication date of the text
  - Summarizing key numerical data points
313. What is a typical goal of text-based predictive analytics?
- Visualizing the structure of text data
  - Forecasting non-textual variables based on text correlations
  - Summarizing documents into key topics
  - Identifying synonyms and collocations
314. What is one way to generate effective features for text-based predictive models?
- Use high-level semantic features like topics instead of individual words
  - Ignore variations in word usage across contexts
  - Rely exclusively on raw word counts
  - Minimize the use of text mining algorithms
315. Associating non-text data with text analysis allows:
- Static representations of knowledge
  - Context-sensitive insights and trend detection
  - Reduction of computational complexity
  - Easier extraction of synonyms from documents
316. Which task falls under the category of mining knowledge about the observed world?
- Analyzing public sentiment towards policies
  - Identifying linguistic patterns in English texts
  - Extracting factual statements about entities or events
  - Predicting future events based on text trends
317. What is one major distinction between mining knowledge about text producers and observed worlds?
- Producers provide objective facts, while the observed world includes opinions
  - Text producers often include subjective statements, while observed worlds provide facts
  - Observed worlds rely on metadata, while producers focus on clustering
  - Producers focus on syntactic rules, while observed worlds involve lexicons
318. Predictive models that leverage text and non-text data often benefit from:
- Ignoring correlations between data types
  - Combining features from text-based topics and historical numerical data
  - Reducing all metadata to a single category
  - Prioritizing physical sensors over subjective text
319. Sentiment analysis primarily aims to:
- Summarize factual information in text
  - Understand subjective opinions in text data
  - Predict numerical values from non-text data
  - Enhance text-based clustering techniques
320. How does clustering contribute to text mining?
- It predicts future variables based on word usage
  - It groups similar text objects for exploratory analysis
  - It focuses exclusively on summarizing text data
  - It automates metadata extraction
321. Joint analysis of text and non-text data enables:
- Standalone text classification models
  - A deeper understanding of context-sensitive knowledge
  - Simplified natural language processing tasks
  - Elimination of metadata dependencies
322. Which metadata is commonly associated with text data?
- User browsing history
  - Sensor calibration details
  - Time and location of text creation
  - Physical sensor measurements
323. Why are high-level features like topics preferred for text mining models?
- They simplify sentiment analysis
  - They reduce computational cost significantly
  - They address ambiguity and variations in word usage
  - They remove all subjective elements from text
324. Text categorization differs from text clustering because it:
- Groups similar objects for exploration

- b) Assigns text to predefined categories
  - c) Ignores linguistic patterns in text
  - d) Focuses only on opinion mining
325. What is a common use case for text summarization?
- a) Discovering new associations between words
  - b) Generating concise overviews of document content
  - c) Identifying author sentiment across a corpus
  - d) Predicting stock market trends
326. What is the purpose of combining text data analysis with data mining techniques?
- a) To streamline document clustering
  - b) To extract actionable knowledge for decision-making
  - c) To eliminate subjective elements from datasets
  - d) To prioritize multimedia data over textual data
327. Why is the division between text data access and text data analysis stages described as "artificial"?
- a) Because both stages rely on identical algorithms
  - b) Because sophisticated applications interleave these stages iteratively
  - c) Because text access inherently includes text mining features
  - d) Because both stages ignore non-textual metadata
328. How can text data analysis enhance business intelligence?
- a) By identifying real-time trends in sensor-generated data
  - b) By revealing customer opinions and comparing competitor products
  - c) By optimizing the storage of customer reviews
  - d) By clustering numerical data into actionable categories
329. How do "subjective sensors" differ from traditional physical sensors?
- a) Subjective sensors provide unbiased numerical data
  - b) Subjective sensors interpret and express personal observations of the real world
  - c) Subjective sensors rely solely on relational database formats
  - d) Subjective sensors are programmed to analyze only text data
330. Joint mining of text and non-text data enables:
- a) Separate analysis of text and numerical data
  - b) Integration of context-sensitive insights with metadata
  - c) Complete automation of all text summarization tasks
  - d) Reduction of storage space for multimedia data
331. When combining text and non-text data, what can temporal metadata (e.g., timestamps) enable?
- a) Real-time clustering of unrelated text data
  - b) Discovery of trends and comparisons across time periods
  - c) Replacement of text data with numerical data
  - d) Isolation of subjective opinions from metadata
332. What is a significant benefit of combining predictive text models with historical non-text data?
- a) Simplified storage of text summaries
  - b) Improved accuracy of predictions through complementary features
  - c) Elimination of data preprocessing requirements
  - d) Automation of all text mining tasks
333. Which text mining task focuses on identifying and interpreting patterns in lexical data?
- a) Opinion mining
  - b) Word association mining
  - c) Predictive analytics
  - d) Topic modeling
334. In predictive analytics, why might semantic topics be more effective than word-level features?
- a) They reduce the need for clustering
  - b) They provide higher-level insights, addressing ambiguity and variations in word usage
  - c) They eliminate all subjective elements from text data
  - d) They rely solely on metadata for contextual analysis
335. What is a primary objective of topic modeling in text mining?
- a) Predicting numerical trends based on text
  - b) Extracting high-level themes from large corpora of text
  - c) Associating metadata with clustered text
  - d) Reducing redundancy in text summarization
336. Opinion mining and sentiment analysis focus on:
- a) Identifying factual statements about real-world entities
  - b) Extracting subjective attitudes and emotions from text
  - c) Combining text data with geographical metadata
  - d) Clustering text data into exploratory groups
337. What distinguishes text categorization from clustering?
- a) Categorization assigns text to predefined classes, while clustering groups similar texts without predefined labels
  - b) Categorization uses exploratory analysis, while clustering focuses on summaries
  - c) Categorization relies solely on numerical metadata, while clustering uses lexical data
  - d) Categorization is less structured than clustering
338. How can metadata, such as location, enhance the analysis of text data?
- a) By creating standalone summaries for documents
  - b) By providing new perspectives for comparative analysis
  - c) By clustering documents based solely on geographical factors
  - d) By ignoring temporal trends in text data
339. Which text mining task might benefit the most from integrating multimedia non-text data?
- a) Sentiment analysis
  - b) Predictive analytics
  - c) Word association mining
  - d) Clustering
340. Associating topics with time can generate:
- a) Trend analysis of textual themes over specific periods

- b) Sentiment predictions for text producers
  - c) Standalone summaries of time-specific clusters
  - d) Metadata-free clustering models
341. What is a potential use case for text mining in disaster response?
- a) Automating document summarization for past events
  - b) Monitoring social media for early warning signs of natural disasters
  - c) Identifying customer sentiment in product reviews
  - d) Summarizing technical documentation for first responders
342. Why might a text mining model integrate data from social media platforms?
- a) To simplify text preprocessing requirements
  - b) To analyze real-time opinions or trends on specific topics
  - c) To reduce dependency on physical sensors
  - d) To increase clustering efficiency in text datasets
343. Predictive analytics in text mining often leverages correlations between:
- a) Metadata and numerical data
  - b) Text content and external real-world variables
  - c) Word frequency and document size
  - d) Clustering algorithms and lexical features
344. How does joint analysis of text and metadata enhance clustering?
- a) It reduces computational complexity
  - b) It enables clustering based on time, location, or other contextual factors
  - c) It eliminates subjective data elements
  - d) It avoids the need for word-level feature extraction
345. What is a common output of text summarization tasks?
- a) Extracted real-world variables for forecasting
  - b) Concise representations of key information in a text
  - c) Sentiment analysis of subjective opinions
  - d) Trends in geographical metadata
346. What is the primary goal of clustering in text analysis?
- a) Summarizing document content
  - b) Grouping similar objects to reveal inherent structure
  - c) Predicting future trends in text datasets
  - d) Extracting semantic relationships between terms
347. Which of the following is a unique advantage of clustering algorithms?
- a) They require labeled training data
  - b) They can only process small text datasets effectively
  - c) They are unsupervised and applicable to any text dataset
  - d) They rely on predefined similarity metrics for optimal performance
348. What does term clustering typically enable in text mining applications?
- a) Improved document summarization
  - b) Creation of a thesaurus or finding semantic concepts
  - c) Real-time sentiment analysis
  - d) Redundant removal of query results
349. Clustering results are often used to:
- a) Identify predefined categories in a dataset
  - b) Provide an overview of data for exploratory analysis
  - c) Replace document similarity measures
  - d) Automate topic summarization tasks
350. What is the concept of "clustering bias"?
- a) The tendency of clusters to include noise from unrelated documents
  - b) The specific perspective used to define similarity for clustering
  - c) The inability of clustering algorithms to process high-dimensional data
  - d) The use of probabilistic models in clustering methods
351. Which type of clustering method merges smaller groups incrementally to form larger clusters?
- a) Divisive clustering
  - b) Agglomerative clustering
  - c) Model-based clustering
  - d) Hierarchical divisive clustering
352. What is a key characteristic of model-based clustering?
- a) Hard cluster assignment
  - b) Creation of a dendrogram
  - c) Assignment of objects to clusters with probabilistic distributions
  - d) Use of fixed similarity metrics for partitioning
353. What is an advantage of similarity-based clustering over model-based clustering?
- a) Ability to perform soft cluster assignment
  - b) Direct reliance on predefined similarity functions
  - c) Flexibility to encode complex probabilistic constraints
  - d) Application to both text and multimedia data
354. In the K-means algorithm, what role does the centroid play?
- a) Acts as a measure of cluster similarity
  - b) Represents the central point of a cluster
  - c) Identifies the largest outlier in a cluster
  - d) Maximizes inter-cluster separation
355. Which clustering approach generates clusters by repeatedly dividing a dataset?
- a) Agglomerative clustering
  - b) K-means clustering
  - c) Complete-link clustering
  - d) Average-link clustering
356. What is a key criterion for evaluating the coherence of a cluster?
- a) The separation of objects in different clusters
  - b) The similarity of objects within the same cluster
  - c) The utility of clusters in specific applications
  - d) The hierarchical structure of clusters
357. Why is choosing the right feature representation important in clustering?
- a) To increase the number of clusters generated
  - b) To minimize computational complexity

- c) To capture crucial concepts that differentiate clusters  
d) To reduce the reliance on similarity functions
358. What is one challenge in evaluating term clustering?
- Lack of predefined document categories
  - Difficulty in defining word-to-word similarity algorithms
  - Over-reliance on cluster labels for scoring
  - Inability to handle high-dimensional data
359. Which measure ensures that clusters are far from each other?
- Coherence
  - Separation
  - Utility
  - Soft clustering
360. When evaluating clusters for search applications, which metric can be used to assess improvement?
- Precision at N
  - Mean Average Precision (MAP)
  - Normalized Cumulative Discounted Gain (NDCG)
  - All of the above
361. What does the single-link clustering algorithm focus on?
- Maximum distance between elements in clusters
  - Minimum distance between elements in clusters
  - Average distance between cluster centroids
  - Random assignment of clusters
362. How does complete-link clustering differ from single-link clustering?
- It merges clusters with the largest maximum distance between elements
  - It produces looser clusters with larger diameters
  - It prioritizes clusters with overlapping elements
  - It keeps clusters compact by minimizing intra-cluster distances
363. What is the significance of the Expectation-Maximization (EM) algorithm in K-means?
- It automates the selection of similarity metrics
  - It alternates between assigning points to clusters and updating centroids
  - It replaces similarity-based clustering methods entirely
  - It focuses on hierarchical clustering alone
364. How does hierarchical clustering represent the merging of clusters?
- As a weighted similarity matrix
  - As a term-document frequency chart
  - As a dendrogram
  - As a centroid-based assignment graph
365. When combining clustering with search engines, what feature can clustering provide?
- Faster indexing of documents
  - Summarization of query results
  - Identification of redundant metadata
  - Automatic query reformulation
366. What is a unique feature of supervised clustering compared to unsupervised clustering?
- It automatically generates predefined categories.
  - It incorporates user-provided constraints on cluster membership.
  - It requires hierarchical structures for clusters.
  - It eliminates the need for similarity measures.
367. Which aspect of clustering is typically not automated and often requires human input?
- Determining the number of clusters
  - Calculating intra-cluster coherence
  - Generating cluster labels
  - Computing inter-cluster separation
368. In text clustering, why might a user want to set the number of clusters beforehand?
- To ensure faster evaluation metrics computation
  - To tailor the clustering output to specific application needs
  - To automatically label the clusters
  - To minimize redundancy in document parsing
369. Which of the following is a practical application of term clustering?
- Generating more diverse search results
  - Reducing memory requirements for clustering algorithms
  - Enabling query expansion by identifying similar terms
  - Automating the generation of training data
370. How does Brown clustering differ from basic agglomerative clustering?
- It incorporates a probabilistic approach to clustering terms hierarchically.
  - It relies on cosine similarity exclusively.
  - It is designed specifically for document-level clustering.
  - It avoids the use of similarity metrics entirely.
371. What makes evaluating utility in clustering more challenging than coherence or separation?
- It requires computing cluster purity scores.
  - It depends on the specific application and the overall effectiveness of the system.
  - It focuses entirely on statistical measures like F1 or MAP.
  - It eliminates the need for human evaluation.
372. Why are predefined categories useful in evaluating clustering algorithms?
- They simplify cluster labeling tasks.
  - They provide a benchmark to compare the clustering output against ground truth.
  - They eliminate the need for similarity metrics.
  - They reduce computation time for hierarchical clustering.
373. What is the main limitation of K-means regarding evaluation of clusters?
- Sensitivity to outliers
  - Difficulty in interpreting centroid values
  - Dependence on initial random centroids
  - Over-reliance on cosine similarity

374. When using clustering in exploratory analysis, how can iterative refinement help?
- a) By allowing users to adjust similarity thresholds dynamically
  - b) By automating the labeling of all clusters
  - c) By generating additional metadata for documents
  - d) By preventing over-segmentation of clusters
375. Which clustering measure is most useful when comparing text documents with varying term frequencies?
- a) Jaccard similarity
  - b) Euclidean distance
  - c) Cosine similarity
  - d) Mutual information
376. In search engines, clustering is particularly useful for:
- a) Identifying outliers in query logs
  - b) Organizing retrieval results into groups for easier navigation
  - c) Reducing latency during document indexing
  - d) Automating user intent prediction
377. How can hierarchical clustering improve the navigation of a text corpus?
- a) By providing a tree structure that allows users to drill down into specific clusters
  - b) By generating a single large cluster for faster processing
  - c) By summarizing documents based on cosine similarity
  - d) By automating document labeling during analysis
378. Which task can benefit from clustering both terms and documents simultaneously?
- a) Topic modeling
  - b) Sentiment analysis
  - c) Metadata extraction
  - d) Query reformulation
379. What is the primary advantage of using soft clustering methods?
- a) They reduce computation time by avoiding similarity calculations.
  - b) They allow objects to belong to multiple clusters with associated probabilities.
  - c) They provide deterministic output for hierarchical clustering.
  - d) They eliminate user-defined parameters like the number of clusters.
380. What is a typical way to handle outliers in clustering algorithms like K-means?
- a) Reassign them to the nearest cluster during centroid computation
  - b) Ignore them entirely during the clustering process
  - c) Create dedicated clusters exclusively for outliers
  - d) Normalize similarity metrics to account for outliers
381. What is the primary difference between clustering and categorization techniques for text data?
- a) Clustering uses supervised learning, while categorization is unsupervised.
  - b) Clustering generates predefined labels, whereas categorization creates clusters.
  - c) Categorization uses predefined categories, while clustering groups based on similarity.
  - d) Categorization is always hierarchical, while clustering is flat.
382. In the context of text categorization, what is the purpose of "text annotation"?
- a) To train a clustering model using labeled data.
  - b) To infer properties about entities from text.
  - c) To represent text with multiple levels, such as keywords and categories.
  - d) To eliminate redundancy in text corpora.
383. Which approach is most suitable for categorization when no labeled training data is available?
- a) Rule-based manual categorization.
  - b) Supervised machine learning.
  - c) Generative classifiers.
  - d) Lazy learning algorithms.
384. What is a primary limitation of the manual rule-based approach in text categorization?
- a) It is overly reliant on training data.
  - b) It cannot handle categorical hierarchies.
  - c) It is labor-intensive and does not scale well.
  - d) It fails in the absence of domain knowledge.
385. How do lazy learners like k-NN differ from discriminative classifiers?
- a) Lazy learners use rule-based systems, while discriminative classifiers use probabilistic models.
  - b) Lazy learners require no explicit training, while discriminative classifiers rely on training data.
  - c) Lazy learners are computationally faster than discriminative classifiers during testing.
  - d) Lazy learners are inherently multiclass, while discriminative classifiers are binary.
386. In Naive Bayes classifiers, why is the independence assumption considered "naive"?
- a) It assumes a linear relationship among features.
  - b) It ignores dependencies between features.
  - c) It only applies to binary classification.
  - d) It overestimates the significance of rare features.
387. What is the primary goal of feature selection in text categorization?
- a) To increase the dimensionality of the feature space.
  - b) To improve classification performance by identifying discriminative features.
  - c) To eliminate the need for labeled training data.
  - d) To generate novel features for training data.
388. Which feature representation is likely most effective for capturing sentiment in text?
- a) Average sentence length.
  - b) Unigram and bigram word tokens.
  - c) Part-of-speech tags.
  - d) Sentence parse trees.
389. What is a significant computational challenge of using k-NN for text categorization?
- a) Building the training dataset.
  - b) Calculating centroids for each category.

- c) Performing a search query for each test instance.  
d) Storing weights for feature vectors.
390. Which of the following is an example of a discriminative classifier?
- a) Naive Bayes.  
b) k-Nearest Neighbors.  
c) Support Vector Machines.  
d) Centroid Classifier.
391. How does SVM ensure robustness in text categorization?
- a) By using probabilistic feature distributions.  
b) By maximizing the margin between decision boundaries and classes.  
c) By relying on local structure in the feature space.  
d) By combining rule-based heuristics with training data.
392. Which evaluation metric is most commonly used to assess text categorization performance?
- a) Cosine similarity.  
b) F1 score.  
c) Term frequency-inverse document frequency (TF-IDF).  
d) Perplexity.
393. What is the main advantage of using cross-validation in text categorization evaluation?
- a) It reduces the size of the training set.  
b) It provides a higher accuracy score.  
c) It helps detect overfitting by testing on multiple splits.  
d) It simplifies feature selection.
394. What does a confusion matrix primarily show in text categorization evaluation?
- a) The overall accuracy of the classifier.  
b) The relationship between features and labels.  
c) True positives, false positives, and class-level predictions.  
d) The effect of training data on testing accuracy.
395. Why is unigram word representation often insufficient for sentiment classification?
- a) It fails to capture negations like "not good."  
b) It increases the computational complexity.  
c) It cannot differentiate synonyms.  
d) It does not handle rare words effectively.
396. What does the perceptron classifier aim to optimize during training?
- a) The centroid distance for each feature vector.  
b) The classification margin for multiclass problems.  
c) The weight vector that separates the classes linearly.  
d) The probabilities of features for each class label.
397. What is the key disadvantage of bigram word representations?
- a) They fail to capture negations.  
b) They miss out on rare informative single words.  
c) They are computationally less efficient than unigram representations.  
d) They introduce ambiguity in sentiment classification.
398. Which approach combines multiple binary classifiers to handle multiclass problems?
- a) One-vs-All and All-vs-All.  
b) Lazy learning.  
c) Naive Bayes smoothing.  
d) Structural feature extraction.
399. What is the role of smoothing in Naive Bayes classification?
- a) To enhance feature selection.  
b) To handle unseen or rare feature occurrences.  
c) To normalize feature weights.  
d) To convert probabilities into confidence scores.
400. Why are structural features combined with lexical features effective in text categorization?
- a) They reduce the dimensionality of the feature space.  
b) They capture orthogonal perspectives of the text.  
c) They improve the speed of classification algorithms.  
d) They eliminate the need for additional training data.
401. What is the primary reason for combining multiple feature sets (e.g., unigram and bigram words) in text categorization?
- a) To reduce computational complexity.  
b) To handle sparsity in the training data.  
c) To capture diverse linguistic contexts and enrich the feature space.  
d) To simplify the training process.
402. What does a centroid represent in the nearest-centroid classifier?
- a) A single document vector closest to the class label.  
b) The average vector of all document vectors for a given class.  
c) The document vector with the maximum similarity to all other documents in the class.  
d) The highest-weighted features of the class.
403. Which of the following is a major limitation of k-NN in text categorization?
- a) It cannot handle multiclass problems.  
b) It requires labeled training data.  
c) It is computationally expensive during testing.  
d) It assumes a linear relationship between features.
404. What distinguishes a generative classifier like Naive Bayes from a discriminative classifier like SVM?
- a) Generative classifiers model data distribution, while discriminative classifiers optimize decision boundaries.  
b) Generative classifiers require labeled data, while discriminative classifiers do not.  
c) Generative classifiers operate in low-dimensional spaces, while discriminative classifiers operate in high-dimensional spaces.  
d) Generative classifiers predict probabilities, while discriminative classifiers predict labels directly.
405. Why might a classifier default to predicting the majority class label in an imbalanced dataset?
- a) To maximize precision.  
b) To simplify feature extraction.  
c) To minimize the loss function during training.



- d) To avoid the computational cost of handling minority classes.
406. What is the main advantage of using the perceptron algorithm for linear classification?
- It works well with non-linear datasets.
  - It is robust against noisy training data.
  - It provides an interpretable weight vector for features.
  - It can calculate exact probabilities for classification.
407. How can edit features, like those in SYNTACTICDIFF, improve text classification tasks?
- By reducing feature dimensionality through tokenization.
  - By capturing structural differences between source and reference texts.
  - By enabling faster parsing of grammatical errors.
  - By eliminating noise in low-frequency words.
408. In n-fold cross-validation, what does a high variance in evaluation metrics between folds indicate?
- The algorithm is robust to unseen data.
  - The algorithm has overfitted certain subsets of data.
  - The dataset lacks sufficient features for classification.
  - The dataset is too balanced for effective evaluation.
409. What is a significant computational tradeoff between k-NN and Naive Bayes classifiers?
- k-NN is slower during training, while Naive Bayes is slower during testing.
  - k-NN is faster during training, while Naive Bayes is faster during testing.
  - Both are equally computationally expensive.
  - k-NN requires more memory, while Naive Bayes requires more labeled data.
410. How does SVM create a robust decision boundary for linear classification?
- By assigning weights to individual features.
  - By maximizing the margin between classes in the feature space.
  - By iteratively refining the centroid of each class.
  - By predicting probabilities for class membership.
411. What is a typical use case for a hierarchical categorization system?
- Clustering documents based on similarity.
  - Classifying documents with nested subcategories.
  - Predicting the sentiment of a document.
  - Tokenizing documents into sentence-level features.
412. Why is smoothing necessary for Naive Bayes in sparse datasets?
- To reduce the dimensionality of the feature vectors.
  - To prevent zero probabilities for unseen features.
  - To increase the independence of features.
  - To enable multiclass classification.
413. How does the decision boundary of a perceptron differ from that of k-NN?
- Perceptron considers global data distribution, while k-NN focuses on local neighbors.
  - Perceptron creates non-linear boundaries, while k-NN creates linear ones.
  - Perceptron requires labeled data, while k-NN does not.
  - Perceptron uses centroids, while k-NN uses prototypes.
414. What is the role of grammatical parse trees in feature extraction?
- They eliminate the need for unigram and bigram tokens.
  - They capture high-level syntactic structures for classification.
  - They simplify the classification of short text documents.
  - They ensure semantic coherence in category labels.
415. Which classifier would you choose for a text categorization task with extremely imbalanced class labels?
- k-Nearest Neighbors.
  - Naive Bayes with smoothing.
  - Support Vector Machines.
  - Rule-based classifier.
416. In the context of feature engineering, why is the choice of tokenizer critical for classification accuracy?
- It determines the computational complexity of the classifier.
  - It directly affects the feature representation quality.
  - It eliminates redundancy in feature extraction.
  - It normalizes the input data distribution.
417. How can weighting schemes improve the performance of k-NN classifiers?
- By balancing class distribution in training data.
  - By assigning higher influence to closer neighbors.
  - By eliminating the need for centroids in classification.
  - By simplifying the tokenization process.
418. What is the significance of the diagonal in a confusion matrix for text categorization?
- It represents the overall accuracy of the classifier.
  - It shows the misclassification rates for each label.
  - It contains the true positive rates for each category.
  - It indicates class imbalance in the training data.
419. Which recent advances in machine learning have significantly improved abstractive summarization?
- Support Vector Machines
  - Transformer-based architectures
  - Decision Trees
  - Rule-based systems
420. Why might extractive summarization methods perform better in evaluations using ROUGE?
- They are less computationally expensive.
  - They contain exact text fragments from the original document.
  - They can generate summaries without training.
  - They require less preprocessing than abstractive methods.
421. What is a limitation of using cosine similarity in MMR-based summarization?
- It cannot measure redundancy.

- b) It fails to capture semantic similarity between sentences.
  - c) It increases computational costs exponentially.
  - d) It prioritizes irrelevant sentences.
422. Which summarization technique is best suited for multi-modal content (e.g., text and images)?
- a) Extractive summarization.
  - b) Multimodal summarization with transformers.
  - c) MMR-based extractive summarization.
  - d) Rule-based sentence extraction.
423. What is a potential ethical concern associated with automated summarization?
- a) High computational costs.
  - b) Misrepresentation due to biased summaries.
  - c) Excessive redundancy in outputs.
  - d) The inability to process large documents.
424. What is a primary challenge in multilingual summarization?
- a) High training data requirements for each language.
  - b) The inability to use ROUGE for evaluation.
  - c) Over-reliance on extractive summarization techniques.
  - d) Lack of similarity metrics across languages.
425. Which summarization technique is most suited for real-time applications?
- a) Extractive summarization using precomputed vectors.
  - b) Abstractive summarization with recurrent neural networks.
  - c) Online summarization algorithms.
  - d) Summarization using pre-trained large language models.
426. What distinguishes neural approaches to summarization from traditional techniques?
- a) Neural approaches require no evaluation metrics.
  - b) Neural methods can learn contextual and semantic relationships.
  - c) Neural methods are only suitable for extractive tasks.
  - d) Neural approaches do not rely on training data.
427. Why is summarization important in fields like financial analysis?
- a) It automates trading decisions.
  - b) It consolidates vast text and data into human-readable formats.
  - c) It replaces traditional reporting systems.
  - d) It eliminates the need for manual data labeling.
428. Which metric is commonly used for abstractive summarization evaluation?
- a) Mean Reciprocal Rank (MRR)
  - b) ROUGE
  - c) Jaccard Index
  - d) KL-Divergence

## Answer Key

1. c)	55. b)	109. b)	163. b)	217. a)	271. a)	325. b)	379. b)
2. b)	56. b)	110. c)	164. b)	218. a)	272. c)	326. b)	380. a)
3. b)	57. c)	111. b)	165. b)	219. b)	273. b)	327. b)	381. c)
4. b)	58. c)	112. b)	166. a)	220. b)	274. b)	328. b)	382. c)
5. b)	59. b)	113. b)	167. b)	221. b)	275. a)	329. b)	383. a)
6. a)	60. c)	114. a)	168. c)	222. a)	276. b)	330. b)	384. c)
7. b)	61. a)	115. b)	169. c)	223. b)	277. c)	331. b)	385. b)
8. c)	62. b)	116. b)	170. b)	224. a)	278. b)	332. b)	386. b)
9. b)	63. b)	117. c)	171. b)	225. b)	279. b)	333. b)	387. b)
10. b)	64. b)	118. b)	172. a)	226. a)	280. a)	334. b)	388. b)
11. b)	65. b)	119. b)	173. b)	227. b)	281. b)	335. b)	389. c)
12. b)	66. c)	120. a)	174. c)	228. b)	282. b)	336. b)	390. c)
13. b)	67. b)	121. b)	175. b)	229. b)	283. b)	337. a)	391. b)
14. b)	68. c)	122. b)	176. b)	230. a)	284. a)	338. b)	392. b)
15. b)	69. b)	123. b)	177. a)	231. a)	285. b)	339. b)	393. c)
16. b)	70. b)	124. b)	178. a)	232. b)	286. b)	340. a)	394. c)
17. b)	71. b)	125. b)	179. a)	233. b)	287. b)	341. b)	395. a)
18. b)	72. c)	126. a)	180. b)	234. b)	288. c)	342. b)	396. c)
19. b)	73. c)	127. b)	181. b)	235. a)	289. a)	343. b)	397. b)
20. a)	74. a)	128. c)	182. b)	236. b)	290. b)	344. b)	398. a)
21. b)	75. b)	129. b)	183. a)	237. a)	291. b)	345. b)	399. b)
22. c)	76. c)	130. d)	184. d)	238. d)	292. b)	346. b)	400. b)
23. a)	77. b)	131. b)	185. a)	239. a)	293. b)	347. c)	401. c)
24. b)	78. a)	132. b)	186. b)	240. a)	294. b)	348. b)	402. b)
25. a)	79. b)	133. c)	187. a)	241. b)	295. b)	349. b)	403. c)
26. b)	80. c)	134. a)	188. b)	242. b)	296. b)	350. b)	404. a)
27. c)	81. b)	135. c)	189. b)	243. a)	297. b)	351. b)	405. c)
28. a)	82. b)	136. b)	190. a)	244. a)	298. b)	352. c)	406. c)
29. b)	83. c)	137. b)	191. b)	245. a)	299. b)	353. b)	407. b)
30. b)	84. b)	138. b)	192. a)	246. c)	300. b)	354. b)	408. b)
31. b)	85. a)	139. b)	193. c)	247. b)	301. b)	355. b)	409. b)
32. b)	86. b)	140. c)	194. b)	248. b)	302. c)	356. b)	410. b)
33. c)	87. b)	141. b)	195. a)	249. c)	303. b)	357. c)	411. b)
34. b)	88. c)	142. b)	196. b)	250. c)	304. c)	358. b)	412. b)
35. b)	89. b)	143. a)	197. b)	251. b)	305. b)	359. b)	413. a)
36. b)	90. b)	144. b)	198. b)	252. b)	306. c)	360. d)	414. b)
37. b)	91. c)	145. b)	199. b)	253. b)	307. a)	361. b)	415. c)
38. b)	92. b)	146. a)	200. a)	254. a)	308. b)	362. d)	416. b)
39. c)	93. b)	147. b)	201. c)	255. b)	309. b)	363. b)	417. b)
40. b)	94. a)	148. b)	202. b)	256. b)	310. b)	364. c)	418. c)
41. b)	95. b)	149. b)	203. b)	257. a)	311. a)	365. b)	419. b)
42. b)	96. b)	150. b)	204. c)	258. b)	312. b)	366. b)	420. b)
43. b)	97. b)	151. a)	205. b)	259. b)	313. b)	367. c)	421. b)
44. b)	98. b)	152. b)	206. c)	260. b)	314. a)	368. b)	422. b)
45. b)	99. c)	153. b)	207. a)	261. b)	315. b)	369. c)	423. b)
46. c)	100. a)	154. b)	208. c)	262. c)	316. c)	370. a)	424. a)
47. b)	101. b)	155. c)	209. b)	263. d)	317. b)	371. b)	425. c)
48. b)	102. b)	156. b)	210. a)	264. b)	318. b)	372. b)	426. b)
49. b)	103. b)	157. b)	211. a)	265. b)	319. b)	373. c)	427. b)
50. b)	104. b)	158. c)	212. a)	266. b)	320. b)	374. a)	428. b)
51. c)	105. b)	159. b)	213. b)	267. b)	321. b)	375. c)	
52. b)	106. b)	160. b)	214. b)	268. b)	322. c)	376. b)	
53. b)	107. b)	161. b)	215. c)	269. b)	323. c)	377. a)	
54. b)	108. b)	162. c)	216. b)	270. b)	324. b)	378. a)	