

For each multiple-choice question below, mark the **single best** answer by completely filling in the circle.

1) JSON is an example of

B

- ☐ A programming language
- ☐ A data format standard
- ☐ An unstructured data application
- ☐ None of the above

2) Which of the following transformations would **not** be made by a stemmer?

C

- ☐ going -> go
- ☐ goes -> go
- ☐ went -> go
- ☐ All of the above transformations could be made by a stemmer

3) Consider the following sentence:

~~A~~

A writer is a person who cares what words mean, what they say, how they say it.

B.

Suppose we use a simple tokenizer that transforms to lowercase and removes punctuation. Which of the following is a sparse bag of words representation of the sentence?

- ☐ {a:1, writer:1, is:1, a:1, person:1, who:1, cares:1, what:1, words:1, mean:1, what:1, they:1, say:1, how:1, they:1, say:1, it:1}
- ☐ {a:2, writer:1, is:1, person:1, who:1, cares:1, what:2, words:1, mean:1, they:2, say:2, how:1, it:1}
- ☐ {writer:1, person:1, cares:1, words:1, mean:1, say:2}
- ☐ None of the above is a sparse bag-of-words representation of the sentence.

D

4) Which of the following are characteristics of applications built using the UIMA standard?

~~B~~

- ☐ Annotation-oriented processing of data streams
- ☐ Use XML for data communication
- ☐ Use a pipeline-like architecture where analyses engines may be chained together
- ☐ All of the above

5) Suppose you have the matrix V resulting from applying latent semantic analysis to a term-document matrix M . Consider a document d in the corpus that was used to create M .

Then describe in 3-5 sentences why this approach may work better for retrieving similar documents than using the term-document matrix M alone.

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and extend across the width of the page. There are no margins, text, or other markings on the paper.