Name: Swathi C

Registration Number:22BDS0387

StatKeyEval

<u>A Statistical Framework for Dynamic Keyword Extraction, Evaluation, and Assessment</u> Automation

Aim:

To implement an automatic short-answer grading system using feature engineering and ensemblebased approaches, focusing on extracting keywords, computing similarity metrics, and generating confidence scores.

Algorithm:

1. Text Preprocessing

- Convert all text to lowercase.
- Remove punctuation marks and numbers.
- Remove common stop words (e.g., "the", "is", "and").
- Strip extra whitespace.

2. Keyword Extraction (Using IGRKE - Information Gain Ratio Keyword Extraction)

- Split preprocessed text into individual words.
- Compute Information Gain Ratio (IGR) for each word based on its entropy and conditional probability within reference answers and student responses.
- Adjust for word significance using the Frequency Adjustment Factor (FAF) to balance importance across reference and response texts.
- Rank words based on IGR × FAF score and extract the most informative keywords.
- Store extracted keywords for both reference answers and student responses.

3. Keyword Mutation (Using SCM - Statistical Co-occurrence Mutation)

- Group responses by question.
- Construct a co-occurrence matrix of word pairs based on document frequency.
- Compute Pointwise Mutual Information (PMI) between words to identify semantically related terms.
- Identify frequently occurring words (≥65% of responses) with high PMI similarity to existing reference keywords.
- Apply Uniqueness Filtering to ensure new keywords add distinct meaning.

• Add these mutated keywords to reference keyword sets for more flexible grading.

4. Vector Representation

- Create a universal keyword list combining all extracted and mutated keywords.
- Represent each reference answer and student response as a binary vector:
- 1 if the keyword is present, 0 if absent.
- Normalize vectors to account for varying answer lengths.

5. Similarity Calculation (Using Multi-Dimensional Similarity Function)

- Compute four similarity metrics between reference and student answer vectors:
- Cosine similarity (Simcos)
- Normalized Euclidean distance (Simeuc)
- Normalized Manhattan distance (Simman)
- Adjusted Pearson correlation (Simpearson)
- Compute a hybrid similarity score:
- Similarity(A,S) = $0.4 \times \text{Simcos} + 0.3 \times \text{Simeuc} + 0.2 \times \text{Simman} + 0.1 \times \text{Simpearson}$

6. Score Generation

- Compute a weighted composite similarity score based on the hybrid similarity function.
- Scale the similarity score to match the original grading scale.
- Apply rounding to get the final predicted score.

7. Performance Evaluation

- Calculate error metrics:
- Root Mean Square Error (RMSE)
- Mean Absolute Error (MAE)
- Mean Absolute Percentage Error (MAPE)
- Generate correlation statistics (Pearson R, Spearman ρ).
- Compute coefficient of determination (R2) to assess predictive accuracy.
- Perform error analysis across different score ranges to identify potential grading inconsistencies.

Research Paper:

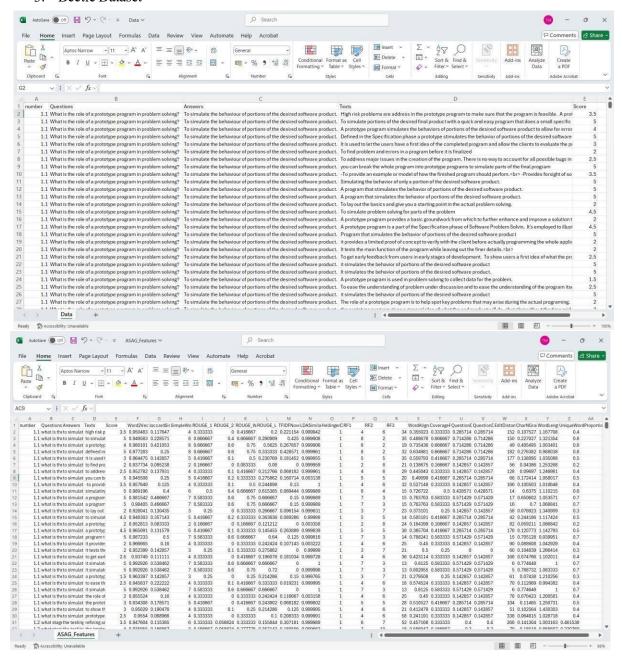
Title: Feature Engineering and Ensemble-Based Approach for Improving Automatic Short-Answer Grading Performance

Authors: Archana Sahu and Plaban Kumar Bhowmick.

Conference/Journal: Educational Data Mining Conference (2018)

Datasets:

- 1. UNT Dataset
- 2. SciEntsBank Dataset
- Beetle Dataset



1. Information Gain Ratio Keyword Extraction (IGRKE)

The Information Gain Ratio Keyword Extraction (IGRKE) method identifies the most informative words in a corpus by evaluating their ability to distinguish between answer keys and student responses.

1.1 Entropy and Information Theory Foundation

For any word W in the corpus, its probability distribution is defined as:

- P(W): Probability of word W occurring.
- $P(\neg W) = 1 P(W)$: Probability of W not occurring.
- C: The document category (Answer Key vs. Student Response).

The total entropy of the word occurrence is given by:

$$H(W) = -P(W) * log_2 P(W) - P(\neg W) * log_2 P(\neg W)$$

This measures the uncertainty of the word's distribution across the dataset.

1.2 Conditional Entropy and Information Gain

The conditional entropy H(C|W) represents the uncertainty in categorizing a document given that word W is known:

$$H(C|W) = -P(A|W) * log_2 P(A|W) - P(R|W) * log_2 P(R|W)$$
 where:

- P(A|W) is the probability of the document being an Answer Key given W.
- P(R|W) is the probability of the document being a Student Response given W.

The Information Gain (IG) quantifies how much knowing W reduces uncertainty about the document's category:

IG(W) = H(C) - H(C|W) where H(C) is the entropy of the category distribution.

1.3 Normalization Using Split Information

To prevent bias toward frequent words, we normalize the Information Gain using Split Information (SI):

$$SI(W) = -P(W) * log_2 P(W) - P(\neg W) * log_2 P(\neg W)$$

This normalizes IG to Information Gain Ratio (IGR):

$$IGR(W) = IG(W) \div SI(W)$$

1.4 Frequency Adjustment Factor (FAF) for Balancing Word Significance

To ensure the extracted keywords are relevant across document categories, we introduce a Frequency Adjustment Factor (FAF):

$$FAF = (F \text{ answer * } F \text{ response}) \div (Total F \text{ answer * Total } F \text{ response}) \text{ where:}$$

• F_answer and F_response are the word frequencies in Answer Keys and Student Responses, respectively.

• Total F answer and Total F response are the total word counts in both document types.

1.5 Final IGRKE Scoring Function

The final scoring function combines IGR and FAF using a logarithmic transformation:

$$Score(W) = IGR(W) * (1 + log(1 + FAF * 1000))$$

2. Statistical Co-occurrence Mutation (SCM)

2.1 Co-occurrence Matrix Construction

A co-occurrence matrix M is built where each entry represents the number of documents where words i and j appear together:

M[i, j] = count of documents where both words i and j appear

2.2 Pointwise Mutual Information (PMI) for Semantic Association

$$PMI(i, j) = log_2 [P(i, j) \div (P(i) * P(j))]$$

2.3 Hybrid Similarity Measure for Keyword Comparison

$$Sim(K1, K2) = 0.6 \times Jaccard_Sim(K1, K2) + 0.4 \times PMI_Sim(K1, K2)$$

3. Similarity Scoring Function for Answer Matching

Similarity(A, S) = $0.4 \times \text{Jaccard_Sim}(A, S) + 0.4 \times \text{Coverage}(A, S) + 0.2 \times \text{Position_Score}(A, S)$ where:

- Jaccard Similarity measures word overlap.
- Coverage measures the proportion of answer key words found.
- Position Score captures structural similarity:

Position Score(A, S) = Average (1 - | Pos A(w)
$$\div$$
 |A| - Pos S(w) \div |S| |)

Code:

IGRKE Function:

IGRKE <- function(word_occurrences, answer_docs, response_docs) { total_docs <- length(answer_docs) + length(response_docs) docs_with_word <- sum(word_occurrences > 0)
P_W <- docs_with_word / total_docs
P not W <- 1 - P W H W

```
<-0 if (P W > 0 && P_W
< 1) {
  H_{-}W < - -P_{-}W * log2(P_{-}W) - P_{-}not_{-}W * log2(P_{-}not_{-}W)
 }
 answer docs with word <- sum(word occurrences[answer docs] > 0) response docs with word
<- sum(word occurrences[response docs] > 0) answer docs without word <-
length(answer docs) - answer docs with word response docs without word <-
length(response docs) - response docs with word
 P A given W <- answer docs with word / docs with word
 P R given W <- response docs with word / docs with word
 P A given not W <- answer docs without word / (total docs - docs with word)
P R given not W <- response docs without word / (total docs - docs with word)
H C given W < 0 if (P A given W > 0 \&\& P R given W > 0) {
  H C given W present <- -P A given W * log2(P A given W) - P R given W *
log2(P R given W)
 } else
  { H_C_given_W_present <-
  0
 H C given not W \le 0 if (P A given not W \ge 0 &&
P R given not W > 0) {
  H C given not W <- -P A given not W * log2(P A given not W) - P R given not W *
log2(P R given not W)
 }
 H_C_given_W <- P_W * H_C_given_W_present + P_not_W * H_C_given_not_W
 P A <- length(answer docs) / total docs
 P R <- length(response docs) / total docs
 H C < -P A * log 2(P A) - P R * log 2(P R)
 IG <- H C - H C given W
 SI <- H W
 IGR \leftarrow if(SI > 0) IG / SI else 0
 F answer <- sum(word occurrences[answer docs])
```

```
F_response <- sum(word_occurrences[response_docs])
 Total_F_answer <- sum(sapply(answer_docs, function(doc) sum(word_occurrences[doc])))
 Total_F_response <- sum(sapply(response_docs, function(doc) sum(word_occurrences[doc])))
 FAF <- (F_answer * F_response) / (Total_F_answer * Total_F_response)
Score \leftarrow IGR * (1 + \log(1 + \text{FAF} * 1000)) return(Score) }
SCM Function:
SCM <- function(corpus, answer keywords, student keywords, threshold = 0.7)
{ all_words <- unique(c(unlist(corpus))) n_words <- length(all_words)
word indices \leq- setNames(1:n words, all words) M \leq- matrix(0, nrow = n words,
ncol = n \ words) for (doc in corpus) {
                                         doc words <- unique(doc)
                                                                      for (i in
1:(length(doc_words) - 1)) {
                                for (j in (i+1):length(doc_words)) {
                                                                        word i <-
                  word_j <- doc_words[j]
doc words[i]
                                               idx_i <- word_indices[[word_i]]
idx_j <- word_indices[[word_j]]</pre>
     M[idx_i, idx_j] <- M[idx_i, idx_j] + 1
M[idx j, idx i] \le M[idx j, idx i] + 1
    }
 doc_count <- length(corpus) word_probs <- numeric(n_words) for (i in 1:n_words)
    word <- all_words[i]
                          word_probs[i] <- sum(sapply(corpus, function(doc) word</pre>
%in% doc)) / doc count
 PMI \le matrix(0, nrow = n words, ncol = n words)
for (i in 1:n words) {
                                                 if (i
                        for (j in 1:n words) {
!= j \&\& M[i,j] > 0) {
                          joint_prob <- M[i,j] /
doc_count
               expected_prob <- word_probs[i] *
word_probs[j]
                   if (expected\_prob > 0) {
      PMI[i,j] <- log2(joint prob / expected prob)
```

```
}
 } mutation candidates <- list() for (word in student keywords)
{ if (!(word %in% answer keywords) && word %in% all words) {
word idx <- word indices[[word]]
                                    pmi scores <- numeric()</pre>
(ans word in answer keywords) {
                                     if (ans word %in% all words)
       ans idx <- word indices[[ans word]]
                                                 pmi scores <-
c(pmi scores, PMI[word idx, ans idx])
    }
   }
   avg pmi <- mean(pmi scores)
                                    position <- match(word, student keywords) /</pre>
                            position weight <-1 - (0.5 * position)
length(student keywords)
                                                                    score <- avg pmi *
                   max pmi <- max(pmi_scores)
                                                   normalized max pmi < -(max pmi + 10)
position weight
/20
       uniqueness <- 1 - normalized max pmi
                                                if (uniqueness >= threshold)
{ mutation candidates[[word]] <- list(word = word, score = score, uniqueness =
uniqueness)
 sorted candidates <- mutation candidates [order(sapply(mutation candidates, function(x) x$score),
decreasing = TRUE)]
 jaccard sim <- length(intersect(answer keywords, student keywords))/
length(union(answer keywords, student keywords))
pmi_scores <- numeric() for (ans_word in answer keywords)
{ if (ans word %in% all words) {
                                      ans idx <-
word indices[[ans word]] for (stud word in
student keywords) {
                        if (stud word %in% all words)
{ stud idx <- word_indices[[stud_word]]
                                            pmi scores
<- c(pmi scores, PMI[ans idx, stud idx])
```

```
pmi_sim <- (mean(pmi_scores) + 10) / 20 hybrid sim <- (0.6 * jaccard sim)
+ (0.4 * pmi sim) matched keywords <- intersect(answer keywords,
student keywords) position scores <- numeric() for (word in
matched keywords) {
                                    pos a <- match(word,
answer keywords) / length(answer keywords) pos s <- match(word,
student keywords) / length(student keywords) position scores <-
c(position scores, 1 - abs(pos a - pos s))
 position score <- if (length(position scores) > 0) mean(position scores) else 0 coverage
<- length(matched keywords) / length(answer keywords) similarity score <- (0.4 *
jaccard sim) + (0.4 * coverage) + (0.2 * position score)
return(list( mutation candidates = sorted candidates,
                                                          similarity score =
similarity score, jaccard similarity = jaccard sim,
                                                          coverage = coverage,
                                   position score = position score
 ))
}
```

Similarity Scoring Function:

```
calculate_similarity <- function(answer_keywords, student_keywords) { jaccard_sim <- length(intersect(answer_keywords, student_keywords)) / length(union(answer_keywords, student_keywords)) / length(answer_keywords)) / length(answer_keywords) matched_keywords <- intersect(answer_keywords, student_keywords) position_scores <- numeric() for (word in matched_keywords) { pos_a <- match(word, answer_keywords) / length(answer_keywords) pos_s <- match(word, student_keywords) / length(student_keywords) position_scores <- c(position_scores, 1 - abs(pos_a - pos_s)) } position_score <- if (length(position_scores) > 0) mean(position_scores) else 0 similarity_score <- (0.4 * jaccard_sim) + (0.4 * coverage) + (0.2 * position_score) return(similarity_score) }
```

Result:

This novel framework integrates:

- **IGRKE for statistical keyword extraction** using information theory and frequency normalization.
- SCM for identifying semantically related keywords based on co-occurrence and PMI.
- A similarity function that combines keyword overlap, coverage, and positional alignment.

By relying purely on statistical principles without pre-trained models, this framework ensures robustness and adaptability across domains.