

# Document Similarity Analyzer

## Source Code Listing - Rust Implementation

### Table of Contents

#### Entry Points

1. src/main.rs
2. src/lib.rs

#### API Module

3. src/api/mod.rs
4. src/api/server.rs
5. src/api/handlers.rs
6. src/api/error.rs

#### Core Module

7. src/core/mod.rs
8. src/core/normalize.rs
9. src/core/tokenize.rs
10. src/core/tf.rs
11. src/core/idf.rs
12. src/core/vectorize.rs
13. src/core/similarity.rs
14. src/core/matrix.rs
15. src/core/pipeline.rs

#### Models Module

16. src/models/mod.rs
17. src/models/document.rs
18. src/models/request.rs
19. src/models/response.rs

#### Tests

20. tests/integration.rs

## Entry Points

---

## 1. src/main.rs

```

1  //! Document Similarity Analyzer - Main Entry Point
2  //!
3  //! A backend service for analyzing document similarity using TF-IDF
4  //! and Cosine Similarity with parallel processing.
5
6  use tracing_subscriber::{layer::SubscriberExt, util::SubscriberInitExt};
7
8  use document_similarity_analyzer::api::run_server;
9
10 #[tokio::main]
11 async fn main() -> anyhow::Result<()> {
12     // Initialize tracing/logging
13     tracing_subscriber::registry()
14         .with(
15             tracing_subscriber::EnvFilter::try_from_default_env()
16                 .unwrap_or_else(|_| "document_similarity_analyzer=debug,tower_http=debug".into()),
17         )
18         .with(tracing_subscriber::fmt::layer())
19         .init();
20
21     // Get port from environment or use default
22     let port: u16 = std::env::var("PORT")
23         .ok()
24         .and_then(|p| p.parse().ok())
25         .unwrap_or(3000);
26
27
28
29 }
```

## 2. src/lib.rs

```

1  //! Document Similarity Analyzer
2  //!
3  //! A backend service for analyzing document similarity using TF-IDF and Cosine Similarity
4  //! with parallel processing powered by Rayon.
5  //!
6  //! ## Architecture
7  //! - `api` - HTTP API handlers and server configuration
8  //! - `core` - Pure functions for text processing and similarity computation
9  //! - `models` - Immutable data structures
10 pub mod api;
11 pub mod core;
12 pub mod models;
13
```

## API Module

---

### 3. src/api/mod.rs

```

1 //! HTTP API module
2
3 mod error;
4 pub mod handlers;
5 mod server;
6
7 pub use error::AppError;
8 pub use handlers::{analyze_handler, health_handler, validate_request};
9 pub use server::{create_router, run_server};
10

```

### 4. src/api/server.rs

```

1 //! HTTP server configuration
2
3 use axum::{
4     routing::{get, post},
5     Router,
6 };
7 use std::net::SocketAddr;
8 use tower_http::cors::{Any, CorsLayer};
9 use tracing::info;
10
11 use super::handlers::{analyze_handler, health_handler};
12
13 /// Creates the Axum router with all routes configured
14 pub fn create_router() -> Router {
15     // Configure CORS
16     let cors = CorsLayer::new()
17         .allow_origin(Any)
18         .allow_methods(Any)
19         .allow_headers(Any);
20
21     Router::new()
22         .route("/health", get(health_handler))
23         .route("/analyze", post(analyze_handler))
24         .layer(cors)
25 }
26
27 /// Runs the HTTP server
28 pub async fn run_server(port: u16) -> anyhow::Result<()> {
29     let app = create_router();
30     let addr = SocketAddr::from(([0, 0, 0, 0], port));
31
32     info!("🚀 Server starting on http://{}{}", addr);
33     info!("📝 POST /analyze - Analyze document similarity");
34     info!("❤️ GET /health - Health check");
35
36     let listener = tokio::net::TcpListener::bind(addr).await?;
37     axum::serve(listener, app).await?;
38
39     Ok(())
40 }

```

## 5. src/api/handlers.rs

```
1 //! API request handlers
2
3 use axum::Json;
4
5 use crate::core::analyze_documents;
6 use crate::models::{AnalyzeRequest, AnalyzeResponse};
7 use super::AppError;
8
9 /// Maximum number of documents allowed per request
10 const MAX_DOCUMENTS: usize = 100;
11
12 /// Minimum number of documents required for comparison
13 const MIN_DOCUMENTS: usize = 2;
14
15 /// Maximum length of a single document in characters
16 const MAX_DOCUMENT_LENGTH: usize = 50_000;
17
18 /// Handler for POST /analyze endpoint
19 ///
20 /// Receives documents and returns their similarity matrix.
21 pub async fn analyze_handler(
22     Json(payload): Json<AnalyzeRequest>,
23 ) -> Result<Json<AnalyzeResponse>, AppError> {
24     // Validate input
25     validate_request(&payload)?;
26
27     // Process documents through the pipeline
28     let result = analyze_documents(&payload.documents);
29
30     // Convert to response format
31     Ok(Json(AnalyzeResponse::from(result)))
32 }
33
34 /// Validates the analyze request
35 pub fn validate_request(request: &AnalyzeRequest) -> Result<(), AppError> {
36     // Check if documents array is empty
37     if request.documents.is_empty() {
38         return Err(AppError::NoDocuments);
39     }
40
41     // Check minimum documents
42     if request.documents.len() < MIN_DOCUMENTS {
43         return Err(AppError::NotEnoughDocuments(request.documents.len()));
44     }
45
46     // Check maximum documents
47     if request.documents.len() > MAX_DOCUMENTS {
48         return Err(AppError::TooManyDocuments(request.documents.len()));
49     }
50
51     // Check each document
52     for (i, doc) in request.documents.iter().enumerate() {
53         if doc.trim().is_empty() {
54             return Err(AppError::EmptyDocument(i));
55         }
56         if doc.len() > MAX_DOCUMENT_LENGTH {
57             return Err(AppError::DocumentTooLong(i, MAX_DOCUMENT_LENGTH));
58         }
59     }
60 }
61
62 Ok(())
63 }
```

```
56
57 // Handler for GET /health endpoint
58 pub async fn health_handler() -> &'static str {
59     "OK"
60 }
61
62
63
```

## 6. src/api/error.rs

```
1 //! Application error types
2
3 use axum::{
4     http::StatusCode,
5     response::{IntoResponse, Response},
6     Json,
7 };
8 use serde::Serialize;
9 use thiserror::Error;
10
11 /// Application error types
12 #[derive(Debug, Error)]
13 pub enum AppError {
14     #[error("Too many documents: {0}, maximum allowed is 100")]
15     TooManyDocuments(usize),
16
17     #[error("Empty document at index {0}")]
18     EmptyDocument(usize),
19
20     #[error("No documents provided")]
21     NoDocuments,
22
23     #[error("Not enough documents: minimum 2 required for comparison, got {0}")]
24     NotEnoughDocuments(usize),
25
26     #[error("Document at index {0} exceeds maximum length of {1} characters")]
27     DocumentTooLong(usize, usize),
28
29     #[error("Internal server error: {0}")]
30     Internal(#[from] anyhow::Error),
31 }
32
33 // Error response body
34 #[derive(Debug, Serialize)]
35 struct ErrorResponse {
36     error: String,
37     code: String,
38 }
39
40 impl IntoResponse for AppError {
41     fn into_response(self) -> Response {
42         let (status, code) = match &self {
43             AppError::TooManyDocuments(_) => (StatusCode::BAD_REQUEST, "TOO_MANY_DOCUMENTS"),
44             AppError::EmptyDocument(_) => (StatusCode::BAD_REQUEST, "EMPTY_DOCUMENT"),
45             AppError::NoDocuments => (StatusCode::BAD_REQUEST, "NO_DOCUMENTS"),
46             AppError::NotEnoughDocuments(_) => (StatusCode::BAD_REQUEST, "NOT_ENOUGH_DOCUMENTS"),
47             AppError::DocumentTooLong(_, _) => (StatusCode::BAD_REQUEST, "DOCUMENT_TOO_LONG"),
48             AppError::Internal(_) => (StatusCode::INTERNAL_SERVER_ERROR, "INTERNAL_ERROR"),
49         };
50
51         let body = ErrorResponse {
52             error: self.to_string(),
53             code: code.to_string(),
54         };
55
56         (status, Json(body)).into_response()
57     }
58 }
59
60 }
```

## Core Module

### 7. src/core/mod.rs

```
1 //! Core processing functions - all pure functions with no side effects
2
3 mod normalize;
4 mod tokenize;
5 mod tf;
6 mod idf;
7 mod vectorize;
8 mod similarity;
9 mod matrix;
10 mod pipeline;
11
12 pub use normalize::normalize_text;
13 pub use tokenize::tokenize;
14 pub use tf::compute_tf;
15 pub use idf::compute_idf;
16 pub use vectorize::vectorize;
17 pub use similarity::cosine_similarity;
18 pub use matrix::compute_similarity_matrix;
19 pub use pipeline::analyze_documents;
```

## 8. src/core/normalize.rs

```
1  //! Text normalization - pure function
2
3  /// Normalizes text by converting to lowercase, removing punctuation,
4  /// and collapsing multiple whitespace into single space.
5  ///
6  /// # Arguments
7  /// * `text` - The input text to normalize
8  ///
9  /// # Returns
10 /// A new String with normalized text
11 ///
12 /// # Example
13 /// ```
14 /// use document_similarity_analyzer::core::normalize_text;
15 /// let result = normalize_text("Hello, World!");
16 /// assert_eq!(result, "hello world");
17 /// ```
17 pub fn normalize_text(text: &str) -> String {
18     text.chars()
19         .map(|c| {
20             if c.is_ascii_punctuation() {
21                 ' '
22             } else {
23                 c.to_ascii_lowercase()
24             }
25         })
26         .collect::<String>()
27         .split_whitespace()
28         .collect::<Vec<&str>>()
29         .join(" ")
30 }
```

## 9. src/core/tokenize.rs

```
1 //! Tokenization - pure function
2
3 /// Tokenizes text into a vector of words by splitting on whitespace.
4 /**
5  * `text` - The input text to tokenize (should be pre-normalized)
6 /**
7 /**
8 /**
9 /**
10 /**
11 /**
12 /**
13 /**
14 /**
15 /**
16 pub fn tokenize(text: &str) -> Vec<String> {
17     text.split_whitespace()
18         .filter(|s| !s.is_empty())
19         .map(|s| s.to_string())
20         .collect()
21 }
22
23
```

## 10. src/core/tf.rs

## 11. src/core/idf.rs

## 12. src/core/vectorize.rs

```
1 //! TF-IDF Vectorization - pure function
2
3 use std::collections::HashMap;
4
5 /// Converts TF and IDF into a TF-IDF vector based on vocabulary order.
6 ///
7 /// # Arguments
8 /// * `tf` - Term Frequency map for a document
9 /// * `idf` - Inverse Document Frequency map (global)
10 /// * `vocab` - Ordered list of all terms (defines vector dimensions)
11 ///
12 /// # Returns
13 /// A Vec<f32> representing the document's TF-IDF vector
14 ///
15 /// # Example
16 /// ```
17 /// use document_similarity_analyzer::core::vectorize;
18 /// use std::collections::HashMap;
19 ///
20 /// let tf: HashMap<String, f32> = [("hello".to_string(), 0.5)].into_iter().collect();
21 /// let idf: HashMap<String, f32> = [("hello".to_string(), 0.693)].into_iter().collect();
22 /// let vocab = vec!["hello".to_string()];
23 /// let vector = vectorize(&tf, &idf, &vocab);
24 /// ```
25 pub fn vectorize(
26     tf: &HashMap<String, f32>,
27     idf: &HashMap<String, f32>,
28     vocabulary: &[String],
29 ) -> Vec<f32> {
30     vocabulary
31         .iter()
32         .map(|term| {
33             let tf_value = tf.get(term).copied().unwrap_or(0.0);
34             let idf_value = idf.get(term).copied().unwrap_or(0.0);
35             tf_value * idf_value
36         })
37         .collect()
38 }
```

## 13. src/core/similarity.rs

```

1  ///! Cosine Similarity calculation - pure function
2
3  /// Computes cosine similarity between two vectors.
4  /// Formula:  $(A \cdot B) / (||A|| * ||B||)$ 
5  ///
6  /// # Arguments
7  /// * `vec_a` - First vector
8  /// * `vec_b` - Second vector
9  ///
10 /// # Returns
11 /// Cosine similarity value between -1.0 and 1.0
12 /// Returns 0.0 if either vector has zero magnitude
13 ///
14 /// # Example
15 /// ``
16 /// use document_similarity_analyzer::core::cosine_similarity;
17 ///
18 /// let a = vec![1.0, 0.0];
19 /// let b = vec![1.0, 0.0];
20 /// assert_eq!(cosine_similarity(&a, &b), 1.0);
21 /// ``
22 pub fn cosine_similarity(vec_a: &[f32], vec_b: &[f32]) -> f32 {
23     if vec_a.len() != vec_b.len() || vec_a.is_empty() {
24         return 0.0;
25     }
26
27     // Compute dot product: A · B
28     let dot_product: f32 = vec_a
29         .iter()
30         .zip(vec_b.iter())
31         .map(|(a, b)| a * b)
32         .sum();
33
34     // Compute magnitudes: ||A|| and ||B||
35     let magnitude_a: f32 = vec_a.iter().map(|x| x * x).sum():<f32>().sqrt();
36     let magnitude_b: f32 = vec_b.iter().map(|x| x * x).sum():<f32>().sqrt();
37
38     // Handle zero magnitude case
39     if magnitude_a == 0.0 || magnitude_b == 0.0 {
40         return 0.0;
41     }
42
43     dot_product / (magnitude_a * magnitude_b)
44 }

```

## 14. src/core/matrix.rs

```

1 //! Similarity Matrix generation - parallel computation
2
3 use rayon::prelude::*;
4 use super::cosine_similarity;
5
6 /// Computes a full NxN similarity matrix from TF-IDF vectors.
7 /// Utilizes parallel processing for efficiency.
8 ///
9 /// # Arguments
10 /// * `vectors` - Slice of TF-IDF vectors, one per document
11 ///
12 /// # Returns
13 /// NxN matrix where matrix[i][j] = cosine_similarity(vectors[i], vectors[j])
14 ///
15 /// # Properties
16 /// - Diagonal is always 1.0 (document is identical to itself)
17 /// - Matrix is symmetric: matrix[i][j] == matrix[j][i]
18 pub fn compute_similarity_matrix(vectors: &[Vec<f32>]) -> Vec<Vec<f32>> {
19     let n = vectors.len();
20
21     if n == 0 {
22         return vec![];
23     }
24
25     // Parallel computation of similarity matrix
26     (0..n)
27         .into_par_iter()
28         .map(|i| {
29             (0..n)
30                 .map(|j| {
31                     if i == j {
32                         1.0 // Diagonal is always 1.0
33                     } else if i < j {
34                         // Compute similarity for upper triangle
35                         cosine_similarity(&vectors[i], &vectors[j])
36                     } else {
37                         // For lower triangle, compute it too
38                         cosine_similarity(&vectors[i], &vectors[j])
39                     }
40                 })
41                 .collect()
42             })
43             .collect()
44         })
45         .collect()
46     }
47 }
```

## 15. src/core/pipeline.rs

```

1  //! Document analysis pipeline - parallel processing orchestration
2
3  use rayon::prelude::*;
4  use std::sync::Arc;
5
6  use crate::models::SimilarityMatrix;
7  use super::{
8      normalize_text,
9      tokenize,
10     compute_tf,
11     compute_idf,
12     vectorize,
13     compute_similarity_matrix
14 };
15
16 /// Analyzes multiple documents and computes their similarity matrix.
17 /// Uses parallel processing for all possible stages.
18 ///
19 /// # Pipeline
20 /// 1. Normalize all documents (parallel)
21 /// 2. Tokenize all documents (parallel)
22 /// 3. Compute TF for each document (parallel)
23 /// 4. Compute IDF globally (single-threaded, needs all TFs)
24 /// 5. Build vocabulary from IDF
25 /// 6. Vectorize each document (parallel)
26 /// 7. Compute similarity matrix (parallel)
27 ///
28 /// # Arguments
29 /// * `documents` - Slice of document strings to analyze
30 ///
31 /// # Returns
32 /// A SimilarityMatrix containing the NxN similarity values and document labels
33 pub fn analyze_documents(documents: &[String]) -> SimilarityMatrix {
34     if documents.is_empty() {
35         return SimilarityMatrix::new(vec![], vec![]);
36     }
37
38     // Generate document labels
39     let labels: Vec<String> = (0..documents.len())
40         .map(|i| format!("doc{}", i))
41         .collect();
42
43     // Stage 1: Normalize (parallel)
44     let normalized: Vec<String> = documents
45         .par_iter()
46         .map(|doc| normalize_text(doc))
47         .collect();
48
49     // Stage 2: Tokenize (parallel)
50     let tokenized: Vec<Vec<String>> = normalized
51         .par_iter()
52         .map(|doc| tokenize(doc))
53         .collect();
54
55     // Stage 3: Compute TF for each document (parallel)
56     let tfs: Vec<_> = tokenized
57         .par_iter()
58         .map(|tokens| compute_tf(tokens))
59         .collect();
60
61     // Stage 4: Compute IDF (single-threaded, needs all TFs)
62     let idf = compute_idf(&tfs);

```

```
56 // Stage 5: Build vocabulary (sorted for consistency)
57 let mut vocabulary: Vec<String> = idf.keys().cloned().collect();
58 vocabulary.sort();
59
60 // Share vocabulary and IDF across threads
61 let vocab_arc = Arc::new(vocabulary);
62 let idf_arc = Arc::new(idf);
63
64 // Stage 6: Vectorize each document (parallel)
65 let vectors: Vec<Vec<f32>> = tfs
66     .par_iter()
67     .map(|tf| vectorize(tf, &idf_arc, &vocab_arc))
68     .collect();
69
70 // Stage 7: Compute similarity matrix (parallel)
71 let matrix = compute_similarity_matrix(&vectors);
72
73 SimilarityMatrix::new(matrix, labels)
74 }
```

## Models Module

---

### 16. src/models/mod.rs

```
1 //! Immutable data models for document similarity analysis
2
3 mod document;
4 mod request;
5 mod response;
6
7 pub use document::*;
8 pub use request::*;
9 pub use response::*;


```

## 17. src/models/document.rs

```
1 //! Document and processing data structures
2
3 use serde::{Deserialize, Serialize};
4 use std::collections::HashMap;
5
6 /// Represents a raw document with its ID
7 #[derive(Debug, Clone, Serialize, Deserialize)]
8 pub struct Document {
9     pub id: String,
10    pub content: String,
11 }
12
13 impl Document {
14     pub fn new(id: impl Into<String>, content: impl Into<String>) -> Self {
15         Self {
16             id: id.into(),
17             content: content.into(),
18         }
19     }
20
21     /// Represents a tokenized document
22 #[derive(Debug, Clone)]
23 pub struct TokenizedDoc {
24     pub id: String,
25     pub tokens: Vec<String>,
26 }
27
28 impl TokenizedDoc {
29     pub fn new(id: impl Into<String>, tokens: Vec<String>) -> Self {
30         Self {
31             id: id.into(),
32             tokens,
33         }
34     }
35 }
36
37     /// Term Frequency map for a single document
38 pub type TermFrequency = HashMap<String, f32>;
39
40     /// Inverse Document Frequency map across all documents
41 pub type InverseDocumentFrequency = HashMap<String, f32>;
42
43     /// TF-IDF Vector representation of a document
44 #[derive(Debug, Clone)]
45 pub struct TfIdfVector {
46     pub id: String,
47     pub vector: Vec<f32>,
48 }
49
50 impl TfIdfVector {
51     pub fn new(id: impl Into<String>, vector: Vec<f32>) -> Self {
52         Self {
53             id: id.into(),
54             vector,
55         }
56     }
57 }
58
59     /// Similarity matrix result
60 #[derive(Debug, Clone, Serialize, Deserialize)]
61 pub struct SimilarityMatrix {
```

```
56     /// NxN similarity matrix where matrix[i][j] is similarity between doc i and doc j
57     pub matrix: Vec<Vec<f32>>,
58     /// Document indices/labels
59     pub index: Vec<String>,
60 }
61 impl SimilarityMatrix {
62     pub fn new(matrix: Vec<Vec<f32>>, index: Vec<String>) -> Self {
63         Self { matrix, index }
64     }
65 }
66
```

## 18. src/models/request.rs

```
1 //! API request models
2
3 use serde::{Deserialize, Serialize};
4
5 /// Request payload for document analysis
6 #[derive(Debug, Clone, Serialize, Deserialize)]
7 pub struct AnalyzeRequest {
8     /// List of document texts to analyze
9     pub documents: Vec<String>,
10 }
11 impl AnalyzeRequest {
12     pub fn new(documents: Vec<String>) -> Self {
13         Self { documents }
14     }
15 }
16
```

## 19. src/models/response.rs

```
1 //! API response models
2
3 use serde::{Deserialize, Serialize};
4
5 /// Response payload for document analysis
6 #[derive(Debug, Clone, Serialize, Deserialize)]
7 pub struct AnalyzeResponse {
8     /// NxN similarity matrix
9     pub similarity_matrix: Vec<Vec<f32>>,
10    /// Document indices/labels
11    pub index: Vec<String>,
12 }
13
14 impl AnalyzeResponse {
15     pub fn new(similarity_matrix: Vec<Vec<f32>>, index: Vec<String>) -> Self {
16         Self {
17             similarity_matrix,
18             index,
19         }
20     }
21
22     impl From<crate::models::SimilarityMatrix> for AnalyzeResponse {
23         fn from(matrix: crate::models::SimilarityMatrix) -> Self {
24             Self {
25                 similarity_matrix: matrix.matrix,
26                 index: matrix.index,
27             }
28         }
29     }
}
```

## Tests

---

## 20. tests/integration.rs

```

1  ///! Integration tests for Document Similarity Analyzer
2  //!
3  //! These tests verify the full pipeline and API endpoints work correctly.
4
5  use document_similarity_analyzer::core::analyze_documents;
6  use document_similarity_analyzer::models::AnalyzeRequest;
7  use document_similarity_analyzer::api::validate_request;
8
9  #[test]
10 fn test_pipeline_with_known_documents() {
11     let documents = vec![
12         "The quick brown fox jumps over the lazy dog".to_string(),
13         "A quick brown dog outpaces a lazy fox".to_string(),
14         "Hello world this is a test".to_string(),
15     ];
16
17     let result = analyze_documents(&documents);
18
19     // Verify matrix dimensions
20     assert_eq!(result.matrix.len(), 3);
21     assert_eq!(result.matrix[0].len(), 3);
22
23     // Verify diagonal is 1.0
24     for i in 0..3 {
25         assert!(
26             (result.matrix[i][i] - 1.0).abs() < 0.001,
27             "Diagonal element [{i}][{i}] should be 1.0"
28         );
29     }
30
31     // Verify symmetry
32     for i in 0..3 {
33         for j in 0..3 {
34             assert!(
35                 (result.matrix[i][j] - result.matrix[j][i]).abs() < 0.001,
36                 "Matrix should be symmetric at [{i}][{j}]"
37             );
38         }
39     }
40
41     // First two documents should have higher similarity than with the third
42     let sim_0_1 = result.matrix[0][1];
43     let sim_0_2 = result.matrix[0][2];
44     let sim_1_2 = result.matrix[1][2];
45
46     assert!(sim_0_1 > sim_0_2, "Documents 0 and 1 should be more similar than 0 and 2");
47     assert!(sim_0_1 > sim_1_2, "Documents 0 and 1 should be more similar than 1 and 2");
48 }
49
50 /// Test deterministic output - same input should give same output
51 #[test]
52 fn test_deterministic_output() {
53     let documents = vec![
54         "machine learning is fascinating".to_string(),
55         "deep learning neural networks".to_string(),
56         "natural language processing".to_string(),
57     ];
58
59     let result1 = analyze_documents(&documents);
60     let result2 = analyze_documents(&documents);
61
62     assert_eq!(result1, result2);
63 }

```

```

56     assert_eq!(result1.matrix.len(), result2.matrix.len());
57
58     for i in 0..result1.matrix.len() {
59         for j in 0..result1.matrix[i].len() {
60             assert!(
61                 (result1.matrix[i][j] - result2.matrix[i][j]).abs() < 0.0001,
62                 "Results should be deterministic at [{i}][{j}]"
63             );
64         }
65     }
66
67     /// Test with identical documents
68 #[test]
69 fn test_identical_documents() {
70     let documents = vec![
71         "hello world".to_string(),
72         "hello world".to_string(),
73     ];
74
75     let result = analyze_documents(&documents);
76
77     for i in 0..2 {
78         for j in 0..2 {
79             assert!(
80                 (result.matrix[i][j] - 1.0).abs() < 0.001,
81                 "Identical documents should have similarity 1.0"
82             );
83         }
84     }
85
86     /// Test with completely different documents
87 #[test]
88 fn test_completely_different_documents() {
89     let documents = vec![
90         "apple banana cherry".to_string(),
91         "xyz uvw rst".to_string(),
92     ];
93
94     let result = analyze_documents(&documents);
95
96     assert!(
97         result.matrix[0][1].abs() < 0.001,
98         "Documents with no common terms should have ~0 similarity"
99     );
100
101     /// Test with many documents for parallel processing
102 #[test]
103 fn test_many_documents_parallel() {
104     let documents: Vec<String> = (0..50)
105         .map(|i| format!("document number {} with some content here", i))
106         .collect();
107
108     let result = analyze_documents(&documents);
109
110     assert_eq!(result.matrix.len(), 50);
111     assert_eq!(result.matrix[0].len(), 50);
112
113     for i in 0..50 {
114         for j in 0..50 {
115             let sim = result.matrix[i][j];
116             assert!(sim >= 0.0 && sim <= 1.0);
117         }
118     }
119 }

```

```

114     }
115 }
116
117 // Test pipeline with special characters
118 #[test]
119 fn test_special_characters() {
120     let documents = vec![
121         "Hello, World! How are you?".to_string(),
122         "HELLO WORLD HOW ARE YOU".to_string(),
123     ];
124
125     let result = analyze_documents(&documents);
126
127     assert!(
128         (result.matrix[0][1] - 1.0).abs() < 0.001,
129         "Documents that normalize to same text should have similarity 1.0"
130     );
131 }
132
133 // Test index labels
134 #[test]
135 fn test_index_labels() {
136     let documents = vec![
137         "doc one".to_string(),
138         "doc two".to_string(),
139         "doc three".to_string(),
140     ];
141
142     let result = analyze_documents(&documents);
143
144     assert_eq!(result.index, vec!["doc0", "doc1", "doc2"]);
145 }
146
147 #[cfg(test)]
148 mod api_tests {
149     use document_similarity_analyzer::models::AnalyzeRequest;
150     use document_similarity_analyzer::api::validate_request;
151
152     #[test]
153     fn test_validate_no_documents() {
154         let request = AnalyzeRequest::new(vec![]);
155         let result = validate_request(&request);
156         assert!(result.is_err());
157     }
158
159     #[test]
160     fn test_validate_single_document() {
161         let request = AnalyzeRequest::new(vec![only one".to_string()]);
162         let result = validate_request(&request);
163         assert!(result.is_err());
164     }
165
166     #[test]
167     fn test_validate_empty_string() {
168         let request = AnalyzeRequest::new(vec!["valid".to_string(), "".to_string()]);
169         let result = validate_request(&request);
170         assert!(result.is_err());
171     }
172
173     #[test]
174     fn test_validate_too_many_documents() {
175         let documents: Vec<String> = (0..101).map(|i| format!("doc {}", i)).collect();
176         let request = AnalyzeRequest::new(documents);
177         let result = validate_request(&request);
178         assert!(result.is_err());
179     }

```

```
172     }
173
174     #[test]
175     fn test_validate_valid_request() {
176         let request = AnalyzeRequest::new(vec![
177             "first doc".to_string(),
178             "second doc".to_string()
179         ]);
180         let result = validate_request(&request);
181         assert!(result.is_ok());
182     }
183
184     #[test]
185     fn test_validate_exactly_100_documents() {
186         let documents: Vec<String> = (0..100).map(|i| format!("doc {}", i)).collect();
187         let request = AnalyzeRequest::new(documents);
188         let result = validate_request(&request);
189         assert!(result.is_ok());
190     }
191
192
193
194
195
196 }
```

---

**End of Source Code Listing**

Document Similarity Analyzer - Rust Implementation

Total: 20 source files