# NAACP: Racial Bias in News Coverage - Technical Project Document

Jackson Fisk, Jack Li, Li Xi  -  2024-02-06 v0.0.1-dev

## Overview

_In this document, based on the available project outline and summary of the project pitch, to the best of your abilities, you will come up with the technical plan or goals for implementing the project such that it best meets the stakeholder requirements._

## A. Provide a solution in terms of human actions to confirm if the task is within the scope of automation through AI.

*To assist in outlining the steps needed to achieve our final goal, outline the AI-less process that we are trying to automate with Machine Learning. Provide as much detail as possible.*

As humans, we can read articles and determine if any article has racial bias and give positive or negative reactions to them. Our approach is to identify the sentiment of race-related words or slang, which indicates probable racial bias. We want to add another layer into the workflow by specifying where the sentiment is pointing towards and which entities are being portrayed. We will display the sentiment outcome on a dashboard and address the broad academic question about racial bias.

## B. Problem Statement:

*In as direct terms as possible, provide the "Data Science" or "Machine Learning" problem statement version of the overview. Think of this as translating the above into a more technical definition to execute on. eg: a classification problem to segregate users into one of three groups on based on the historical user data available from a publicly available database.*

- We want to perform various forms of natural language processing including sentiment analysis on the news text of known racist news articles, in order to ascertain language patterns in these kinds of texts. We will train the machine learning model with new data and determine if the sentiment is positive or negative.

- We will pre-process the text data to highlight the words related to sentiments and specific racial groups.
- We'll also find some pre-trained embeddings on sentiment words and words regarding race

## C. Checklist for project completion

*Provide a bulleted list to the best of your current understanding, of the concrete technical goals and artifacts that, when complete, define the completion of the project. This checklist will likely evolve as your project progresses.*

1. Deliverable 1: Exploratory Data Analysis (EDA) results for analyzing and investigating data sets and summarizing their main characteristics.
2. Deliverable 2: An Excel sheet listing all of the news articles and a few analysis columns, including keywords caught by an ML model, positive or negative sentiment, and a score of confidence.
3. Final deliverable: A dashboard to display our analysis of a specific article. Users will be able to identify the sentiment (positive or negative) of a specific racial group according to our ML model. Users will also be able to view the training set and the model details in the dashboard.

## D. Outline a path to operationalization.

*Data Science Projects should have an operationalized endpoint in mind from the onset. Briefly describe how you see the tool produced by this project being used by the end-user beyond a jupyter notebook or proof of concept. If possible, be specific and call out the relevant technologies that will be useful when making this available to the stakeholders as a final deliverable.*

First, instead of just using Jupyter Notebook as a demo or proof of concept, we aim to deploy the model in a web app with a UI so that users can visualize what our model can achieve. At the same time, our final deliverable should be able to be retrained on new datasets and fine-tuned easily through some API. We are thinking of embedding the model into a browser extension so users could conveniently access our model when reading an online news article, which can be built with Javascript. Another possible way we can deliver is by deploying with Gradio: a simple input output interface for ML models that can be created very quickly. The Gradio app could potentially be delivered earlier as the work is less intensive.

# Resources

## Data Sets

1. [Census Geocoder](#)
2. [Census Data (demographic data)](#)
3. [Census Guide and Census API Tutorial](#)
4. News Articles in sparkgrp SCC project - `/projectnb/sparkgrp/naacp-new-data/data-out`

## References

1. https://naacp.org/about
2. https://www.wgbh.org/
3. https://github.com/BU-Spark/ml-naacp-sentiment

# Weekly Meeting Updates

*Keep track of ongoing meetings in the Project Description document prepared by Spark staff for your project.*

Note: Once this markdown is finalized and merge, the contents of this should also be appended to the Project Description document.