

Bangla Article Similarity and recommendation using Cosine Similarity

Abstract

Automated Bangla article similarity is basically finding the most relevant date that resembles with search text.. Despite several comprehensive textual datasets are available for different languages, a few small datasets are curated on Bangla language. As a result, a few works address Bangla document recommendation problem, and due to the lack of enough data. In this work, we created a large dataset of Bangla articles from **The Daliy Ittefaq**, which contains around 70k articles. This huge diverse dataset helps us to create a system by utilizing TF-IDF features, which finally gives recommendation.

Methodology

We have used various features in this project. They would be clarified by following segments:

Dataset

For dataset, We have used the articles from famous Bengali newspaper “The Daily Ittefaq”. Almost 54K documents have been retrieved from the website of the newspaper for the years 2019 and 2020. Almost 9,00,000 sentences are present in the corpus.

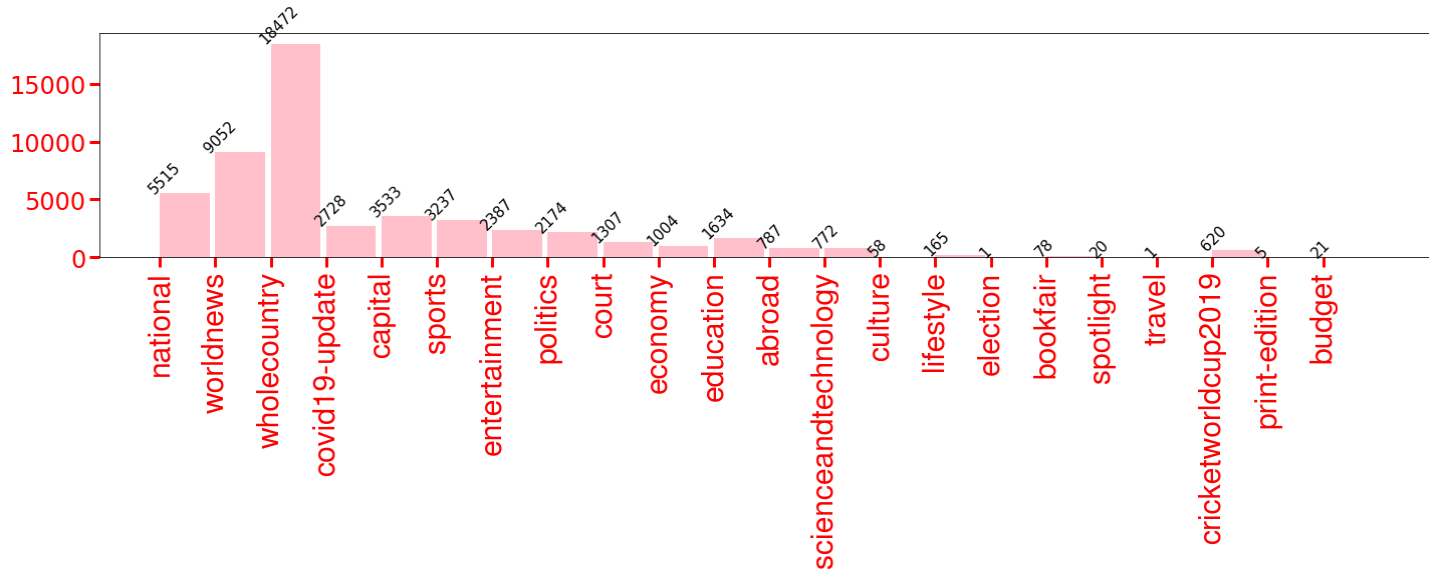
Summary

#Total Article count =====> 53571
#unique unigram =====> 252073
#Total Unigram =====> 13326446
#unique bigram =====> 3718409
#Total Bigram =====> 13579033
#unique trigram =====> 8218523
#Total Trigram =====> 24030358

Category Summary: Total 22 Categories

<u>Type</u>	<u>News Count</u>	<u>Type (Bengali Version)</u>
-------------	-------------------	-------------------------------

'national':	5515,	জাতীয়
'worldnews':	9052,	আন্তর্জাতিক
'wholecountry':	18472,	সারাদেশ
'covid19-update':	2728,	কোভিড-১৯
'capital':	3533,	রাজধানী
'sports':	3237,	খেলা
'entertainment':	2387,	বিনোদন
'politics':	2174,	রাজনীতি
'court':	1307,	কোর্ট
'economy':	1004,	অর্থনীতি
'education':	1634,	শিক্ষা
'abroad':	787,	বিদেশ
'Scienceandtechnology':	772,	বিজ্ঞান ও টেক
'culture':	58,	সংস্কৃতি
'lifestyle':	165,	জীবনযাপন
'bookfair':	78,	বইমেলা
'cricketworldcup2019':	620	ক্রিকেট বিশ্বকাপ -১৯



So, the methodology is consisted of following procedures:

A. Preprocessing

Preprocessing is very important in terms of raw data is concerned and especially in Bengali language where different punctuation marks are there.

a. Punctuation Removal

First of all, I have ensured the proper cleaning of each document by cleaning the punctuation marks. Again many Unicode were not detected, so they were needed to be cleaned as well.

b. Duplicate Sentence Removal

There were many duplicate sentences in the corpus as well. So I have made sure that they unique sentences will only be there clearing the duplicate ones.

c. Stop words Removal

There are many stop words in Bengali languages like “ও”, “এবং”, “আর” etc. These stopwords are identified by computing frequency based unigrams. The most frequent words were the stop words so many stopwords were identified from the frequency based unigrams and from Bengali language point of view.

d .Stemming

There are several words in different formations like শহর, শহরে, শহরের etc. They basically denote the same meaning শহর but with the inclusion of terms like 's', 'es' it shows the different form. To remove redundancy, the conversion of several formations into an actual word is also needed.

e. Tokenization

Word tokenization is done here for each document to compute the tf-idfs of each word of each document which is described in the coming sections.

f.Small documents removal

Article of smaller lengths are removed in this case.

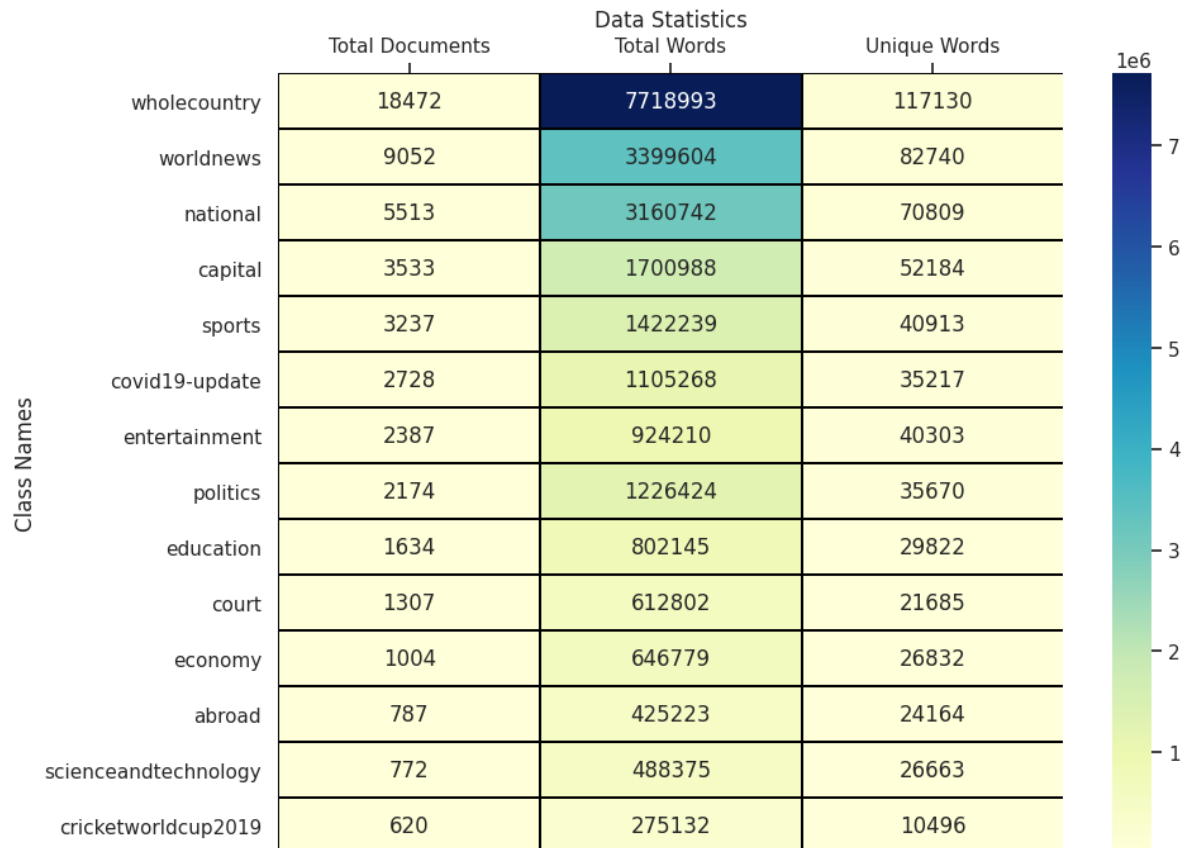


Fig:After pre-processing dataset

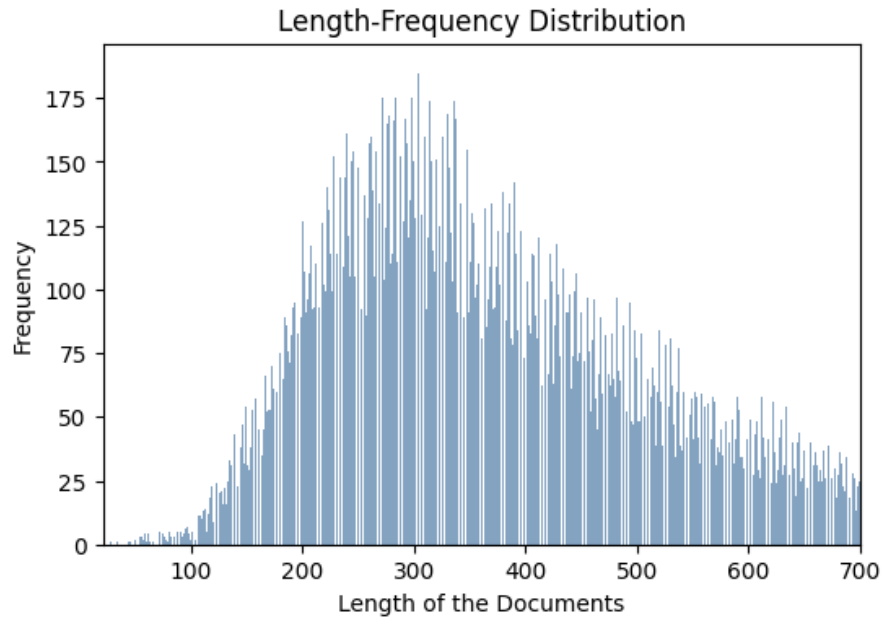


Fig:Lengths of documents

B. Bangla Article Similarity and recommendation using Cosine Similarity

a.TF-IDF

TF-IDF is a measurement that symbolizes how important a word is to a document in a association of documents. Two metrics are multiplied in this regard. They are the term frequency and the inverse document frequency. Term frequency is how many times a word appears in a document divided by the length of the document in terms of word count and the inverse document frequency is the number of total documents divided by the number of documents that the word has appeared in.

$$TF*IDF = \frac{\text{No of perticular term in that document}}{\text{Total no of term in that document}} * \frac{\text{Total No of document}}{\text{No of document where the perticular word appears}}$$

.....(1)

Fig: TF-IDF

b. Cosine Similarity

Cosine similarity measures the similarity between two vectors of an inner product space. It is measured by the cosine of the angle between two vectors and determines whether two vectors are pointing in roughly the same direction. It is often used to measure document similarity in text analysis.

$$\text{similarity}(A,B) = \frac{A \cdot B}{\|A\| \times \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n A_i^2} \times \sqrt{\sum_{i=1}^n B_i^2}} \dots\dots\dots(2)$$

Using the search text the similarity with articles are measured and top 5 articles are recommended.

Shresthos's trigram search Recommender!!!

সাকিব আল হাসান

সাকিব আল হাসান
সাকিব আল হাসানের
সাকিব আল হাসানকে

Cosine Similarity

Search

67% Similarity

করোনায় অসহায়দের মাঝে অর্থ দানের জন্য নিজের প্রিয় ব্যাটটি নিলামে তুলেছিলেন সাকিব আল হাসান। তার এই ব্যাটটির সর্বোচ্চ দাম উঠল ২০ লাখ টাকা। যিনি ব্যাটটি কিনতে চেয়েছেন তিনি যুক্তরাষ্ট্র প্রবাসী বাংলাদেশি। তার নাম রাজ। তবে পেমেন্ট পাওয়ার পর, আনুষ্ঠানিকভাবে বিজয়ীর নাম ঘোষণা করা হবে জানিয়েছেন নিলাম পরিচালনাকারী কর্তৃপক্ষ। বুধবার (২২ এপ্রিল) বিকালেই ফে সবুকে অকশন ফর অ্যাকশন নামক পেইজে ব্যাটটি নিলামে তোলা হয়। ব্যাটটির ভিত্তিমূল্য ধরা হয় ৫ লাখ টাকা। বাংলাদেশ সময় রাত দশটায় ভিডিও কনফারেন্সিংয়ের মাধ্যমে অনলাইনে নিলাম পরিচালনা করে অকশন ফর অ্যাকশন কর্তৃপক্ষ। সাকিব আল হাসান নিজে ভিডিও কনফারেন্সিংয়ে যুক্ত ছিলেন। রাত সোয়া ১১টায় নিলাম শেষ হয়। ব্যাট নিয়ে সাকিব আল হাসান বলেছেন, 'মানুষের জীবনের মূল্যের চেয়ে নিশ্চয়ই ব্যাটের মূল্য বেশি না। মানুষের কল্যাণের জন্য কাজ করতে পারব বলে ভালো লাগছে। সম্পূর্ণ টাকাই করোনা মোকাবেলার জন্য ব্যয় করা হবে।' ব্যাট নিয়ে সাকিব আল হাসান বলেছেন, 'মানুষের জীবনের মূল্যের চেয়ে নিশ্চয়ই ব্যাটের মূল্য বেশি না।' করোনা ভাইরাস মোকাবেলার জন্য ২০১৯ বিশ্বকাপে খেলা নিজের প্রিয় ব্যাট নিলামে বিক্রির সিদ্ধান্ত নেন টাইগার অলরাউন্ডার সাকিব আল হাসান। মঙ্গলবার সামাজিক যোগাযোগমাধ্যম ফেসবুকে লাইভে এসে সাকিব তার ব্যাটটি নিলামে তোলার কথা তুলে ধরেন। করোনা ভাইরাস মোকাবেলার জন্য ২০১৯ বিশ্বকাপে খেলা নিজের প্রিয় ব্যাট নিলামে বিক্রির সিদ্ধান্ত নেন টাইগার অলরাউন্ডার সাকিব আল হাসান। করোনায় মোকাবেলার জন্য সাকিব আল হাসান একটি ফাউন্ডেশন চালু করেছেন। তার নাম সাকিব আল হাসান ফাউন্ডেশন। এই ফাউন্ডেশনের মাধ্যমেই সাকিব অর্থ সংগ্রহ করছেন। তিনি সকলকে ফাউন্ডেশনে অর্থ দেয়ার জন্য অনুরোধ করেছেন। ফাউন্ডেশনের তহবিলে যত টাকা জমা হবে তার সবই মানুষের কল্যাণের জন্য ব্যয় করা হবে। ইন্তেফাক/এসআইইন্তেফাক/এসআই।

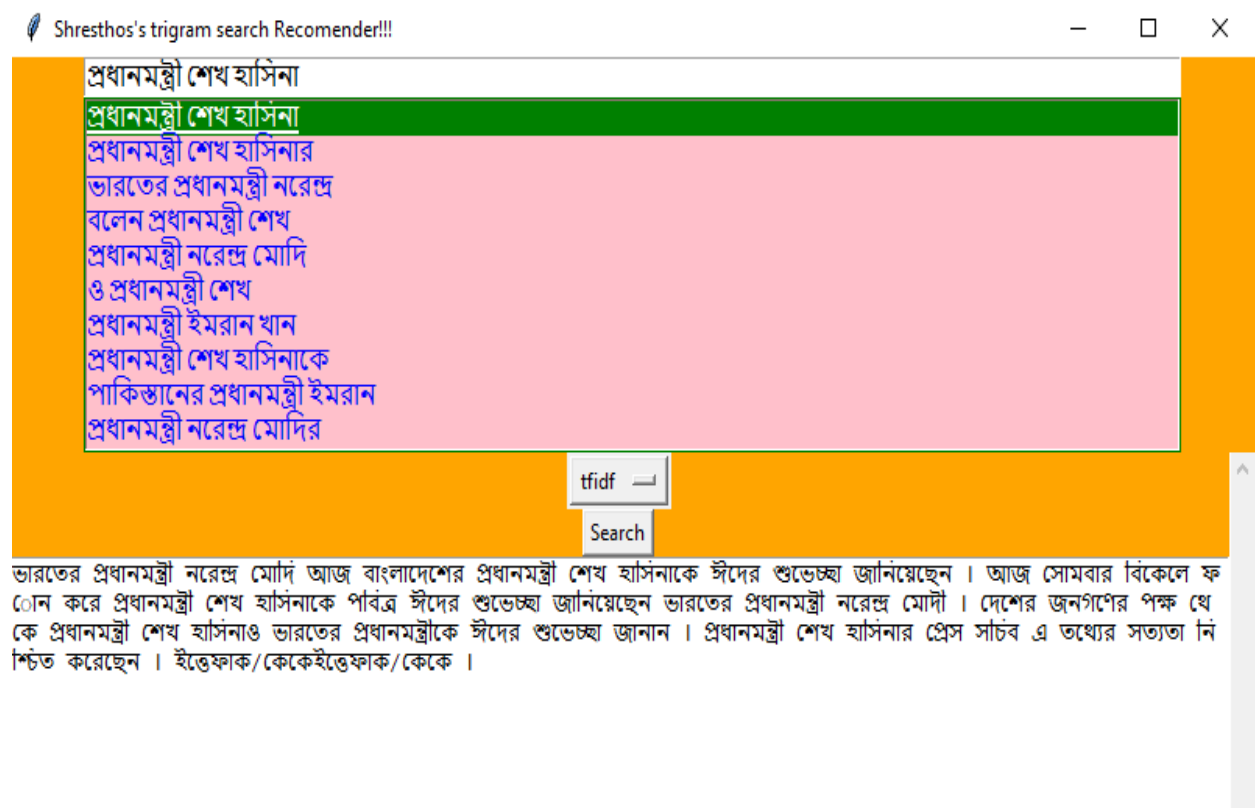
-----Article-----

57% Similarity

করোনা ভাইরাস আতঙ্কে সম্প্রতি যুক্তরাষ্ট্রে একটি হোটেলে কোয়ারেন্টাইন শেষে স্বাভাবিক জীবনে ফিরেছেন সাবেক বিশ্ব সেরা অলরাউন্ডার সাকিব আল হাসান। আর পরিবারের কাছে ফিরেই একটি সসংবাদ দিলেন তিনি। মঙ্গলবার বাংলাদেশ সময় দুপুরে নিজ আফিস

C. Bangla Article Similarity and recommendation using TF-IDF

In this method ,the search strings are tokenized and for those tokens the tfidf values are calculated and the highest valued article is returned.



D. Bangla Article Similarity and recommendation using Direct Search

In this method, the searched string is searched in all the documents and the document that has the highest frequency of the search text is returned. The following figure shows it.

প্রধানমন্ত্রী শেখ হাসিনা

প্রধানমন্ত্রী শেখ হাসিনা

প্রধানমন্ত্রী শেখ হাসিনার

ভারতের প্রধানমন্ত্রী নরেন্দ্র

বলেন প্রধানমন্ত্রী শেখ

প্রধানমন্ত্রী নরেন্দ্র মোদি

ও প্রধানমন্ত্রী শেখ

প্রধানমন্ত্রী ইমরান খান

প্রধানমন্ত্রী শেখ হাসিনাকে

পাকিস্তানের প্রধানমন্ত্রী ইমরান

প্রধানমন্ত্রী নরেন্দ্র মোদির

Direct

Search

প্রধানমন্ত্রী শেখ হাসিনা বলেছেন, আজকের শিশুরাই আগামী দিনে বাংলাদেশকে নেতৃত্ব দেবে। তাই তাদের দেশ প্রেমের আদর্শ নিয়ে বেড়ে উঠতে হবে। জাতির পিতা বঙ্গবন্ধু শেখ মুজিবুর রহমানের আদর্শে উজ্জীবিত হয়ে শিশুদের জীবন গঠনের আহ্বান জানান তিনি। টাঙ্গাপাড়ায় আজ রবিবার দুপুরে জাতির পিতার শততম জন্ম দিন ও জাতীয় শিশু দিবস উপলক্ষে বঙ্গবন্ধুর সমাধিসৌধ কমপ্লেক্সের পা বালক প্লাজায় মাইলা ও শিশু বিষয়ক মন্ত্রণালয় এবং গোপালগঞ্জ জেলা প্রশাসন আয়োজিত শিশু সমাবেশে প্রধান অতিথির ভাষণে তিনি এসব কথা বলেন। প্রধানমন্ত্রী আরো বলেন, শিশুরাই একাদিন দেশের নেতৃত্ব দেবে। উন্নত সমৃদ্ধ বাংলাদেশ গড়তে এখন থেকেই নিজেদের যোগ্য করে গড়ে তুলতে হবে। জাতির পিতার স্বপ্নের সোনার বাংলা প্রতিষ্ঠার মাধ্যমে শিশুদের জন্য একটি সুন্দর বাংলাদেশ গড়তে কাজ করছে তার সরকার। প্রধানমন্ত্রী শেখ হাসিনা বলেন, 'জাতির পিতা বঙ্গবন্ধু শেখ মুজিবুর রহমান মানবদরদী ছিলেন। নিজের বই গারব ছাত্রদের মাঝে বিলিয়ে দিতেন। ফুলে যাওয়ার সময় নিজের ছাতা অন্যকে দিয়ে দিতেন। নিজের গোলার ধান বের করে নিদিয়ায় গরীব-দুঃখী মানুষের মাঝে বিলিয়ে দিতেন তিনি।' প্রধানমন্ত্রী শেখ হাসিনা বলেন, 'জাতির পিতা বঙ্গবন্ধু শেখ মুজিবুর রহমান মানবদরদী ছিলেন।' প্রধানমন্ত্রী শেখ হাসিনা বলেন, বঙ্গবন্ধু নিজের সবকিছু বিলিয়ে দিয়েছেন মানুষের জন্য। মানুষের অধিকারের কথা বলতে গিয়েই বারবার কারাবরণ করেছেন। সবার জন্য শিক্ষা, স্বাস্থ্য নিশ্চিত করার স্বপ্ন দেখেছিলেন। জাতিসংঘের আগেই ১৯৭৪ সালে শিশু অধিকার রক্ষায় শিশু আইন করেন বঙ্গবন্ধু। প্রধানমন্ত্রী শেখ হাসিনা বলেন, বঙ্গবন্ধু নিজের সবকিছু বিলিয়ে দিয়েছেন মানুষের জন্য। তিনি বলেন, ৭৫ এর পর বঙ্গবন্ধুর ৭ মার্চের ভাষণ নিষিদ্ধ করা হয়েছিলো। আর সেই ভাষণ আজ বিশ্বের শ্রেষ্ঠ ভাষণের একটি। তিনি আরো বলেন, শিশুদের সুনাগরিক হিসেবে গড়ে তুলতে আধুনিক শিক্ষায় শিক্ষিত করার কাজ করছে সরকার। বঙ্গবন্ধু যে সুখী সমৃদ্ধ বাংলাদেশ দেখতে চেয়েছিলেন। সেটি বাস্তবায়ন করাই এখন লক্ষ্য। বঙ্গবন্ধু কন্যা শেখ হাসিনা বলেন, 'জাতির পিতা একটি সুন্দর দেশ গড়ে তুলতে চেয়েছিলেন। বাংলাদেশের মানুষ একটি উন্নত জীবন পাবে, এটাই তার লক্ষ্য ছিল। কিন্তু সে কাজ তিনি করে যেতে পারেননি। ১৫ আগস্ট তাঁকে মেরে ফেলা হলো। আমি পরিবার হারালাম, আপনজন হারালাম কিন্তু বাংলাদেশের মানুষ হারায়োঁল তাদের স্বাধীনতার চেতনা, উন্নত জীবন পাওয়ার সম্ভাবনা।' বঙ্গবন্ধু কন্যা শেখ হাসিনা বলেন, 'জাতির পিতা একটি সুন্দর দেশ গড়ে তুলতে চেয়েছিলেন।' প্রধানমন্ত্রী আরো বলেন, ৭৫ এ আমার পিতাকে স্বপরিবারে হত্যার ৬ বছর পর দেশে ফিরে প্রতিজ্ঞা নিয়ে ছিলাম, জাতির পিতার স্বপ্ন বাস্তবায়ন করব। দেশের জন্য কাজ করে যাচ্ছি। দেশে যখন ফিরে আসি জানতাম, যে কোন সময় আমাকেও এই পরিপাত ভোগ করতে হতে পারে। দেশ এগিয়ে যাচ্ছে, এগিয়ে যাবে। কেউ ক্ষুধার্ত থাকবে না, বিনা চিকিৎসায় মারা যাবে না, উন্নত জীবন পাবে। এটা নিশ্চিত করাই এখন লক্ষ্য। প্রধানমন্ত্রী আরো বলেন, ৭৫ এ আমার পিতাকে স্বপরিবারে হত্যার ৬ বছর পর দেশে ফিরে প্রতিজ্ঞা নিয়ে ছিলাম, জাতির পিতার স্বপ্ন বাস্তবায়ন করব। শেখ হাসিনা বলেন, ২০২০ সালে বঙ্গবন্ধুর জন্মশতবার্ষিকী পালন করা হবে। ২০২০ সাল থেকে ২১ সাল পর্যন্ত মুজিব বর্ষ হিসেবে ঘোষণা করা হয়েছে। 'শেখ হাসিনা বলেন, ২০২০ সালে বঙ্গবন্ধুর জন্মশতবার্ষিকী পালন করা হবে।' তিনি পরিবারের ছাড়াছাড়া দেশের মানবপন্থিত জনগণের এ সম্মানে