# Article Analysis Report:

**Prepared by: Vedang Mandhana**
**https://github.com/mvedang**

**Objective: The objective of this project is to collect all articles related to HIV from the Times Of India archives from January 2010 onwards and classify them. Analysis and visualization of the articles should be carried out and presented accordingly.**

# Collection of data:

Articles containing the word 'HIV' in the title were extracted from the Times Of India (TOI) archives dating from January 1st, 2010 to December 18th, 2018.
The archive can be found here: https://timesofindia.indiatimes.com/archive.cms

BeautifulSoup was used to create a web crawler Python script and for the extraction of data from webpages of TOI.
A total of such 1382 articles titles were found containing 'HIV' in the title.
The month of the article, year of the article, title of the article, link to the article and contents of the article were saved to a .csv file.
The script used to collect the data can be found here:
https://github.com/mvedang/Article-Analysis/blob/master/data/Scraping.ipynb

# Preprocessing:

The titles of all the articles were compared and duplicate articles were removed. The count of number of articles reduces to 1348.
All the text data was converted to lowercase for efficient preprocessing.
All the stop words (e.g. the, is, a, an, etc.) present in the Natural Language Toolkit (NLTK) database of stopwords were removed from all the texts.
Regular expressions were used to clean the text data containing special characters.
NLTK's SnowballStemmer was used to replace all words in text data with their stems (if available for particular word).
After preprocessing the data, we observe that all the words are converted into their stems. All the proper nouns are left as it is.

# Data before preprocessing:

NEW DELHI: A vaccine to protect HIV patients from contracting tuberculosis (TB) and ultimately dying of it has finally become a reality. After a seven-year-long trial in Africa, scientists have for the first time developed a vaccine that was succesful in reducing the rate of definite TB infection by almost 39% among 2,000 HIV-infected patients in Tanzania. TB is the biggest killer of HIV-infected patients in the world. In India, over 60% of HIV-infected patients with a weak immunity system get infected with TB and ultimately die of it. Scientists from Dartmouth Medical School (DMS) have reported results of their clinical trial of this new vaccine against TB -- Mycobacterium vaccae (MV) -- in the January 29 online issue of the journal AIDS. The study will be published in the March print issue of the journal. Principal investigator Ford von Reyn from DMS said, "Since development of a new vaccine against TB is a major international health priority, especially for patients with HIV infection, we and our Tanzanian collaborators are very encouraged by the results of the study." The vaccine is a type known as an inactivated, whole-cell mycobacterial vaccine and is expected to be economical to produce and distribute. Von Reyn described the trial as a "significant milestone -- the first to demonstrate that any type of vaccine can prevent an infectious complication of HIV in adults". He added that the next steps are to improve the manufacturing methods to support the production of the larger quantities of the TB vaccine needed for further studies and subsequent clinical use. Health ministry officials in Delhi say the vaccine gives tremendous hope to India which has a huge burden of both TB and HIV patients. Of the over two million HIV-infected Indians, over 10% are expected to have full blown AIDS. Every AIDS patient has 15% chance every year of developing TB, which shows that every AIDS patient will develop TB some time in his life. "That\'s why under India\'s latest AIDS control programme, we are dealing with TB and HIV simultaneously. When people become infected with TB and AIDS, it is almost always an irreversible formula," ministry officials said. Since newly-infected HIV patients risk contracting TB almost immediately, investigators are targeting a strategy for immunization with MV before patients need to start taking antiretroviral drugs. The scientific team at Dartmouth began Phase-I human studies with MV in the United States in 1994 and demonstrated that a multiple-dose series of MV was safe in both healthy subjects and patients with HIV infection. The group then conducted Phase-II studies in larger groups of adults in Zambia and in Finland. In the Zambian trial, researchers found that MV boosted immune responses against TB that had first been primed in childhood with the current TB vaccine, BCG. Subsequently, the DarDar group received NIH funding to conduct the large Phase-III efficacy trial among HIV-infected patients with prior BCG immunization in Tanzania. HIV patients are particularly vulnerable to TB because their immune systems are compromised. The vaccine works by boosting the immune responses of patients who have already been given the BCG vaccine earlier in life.
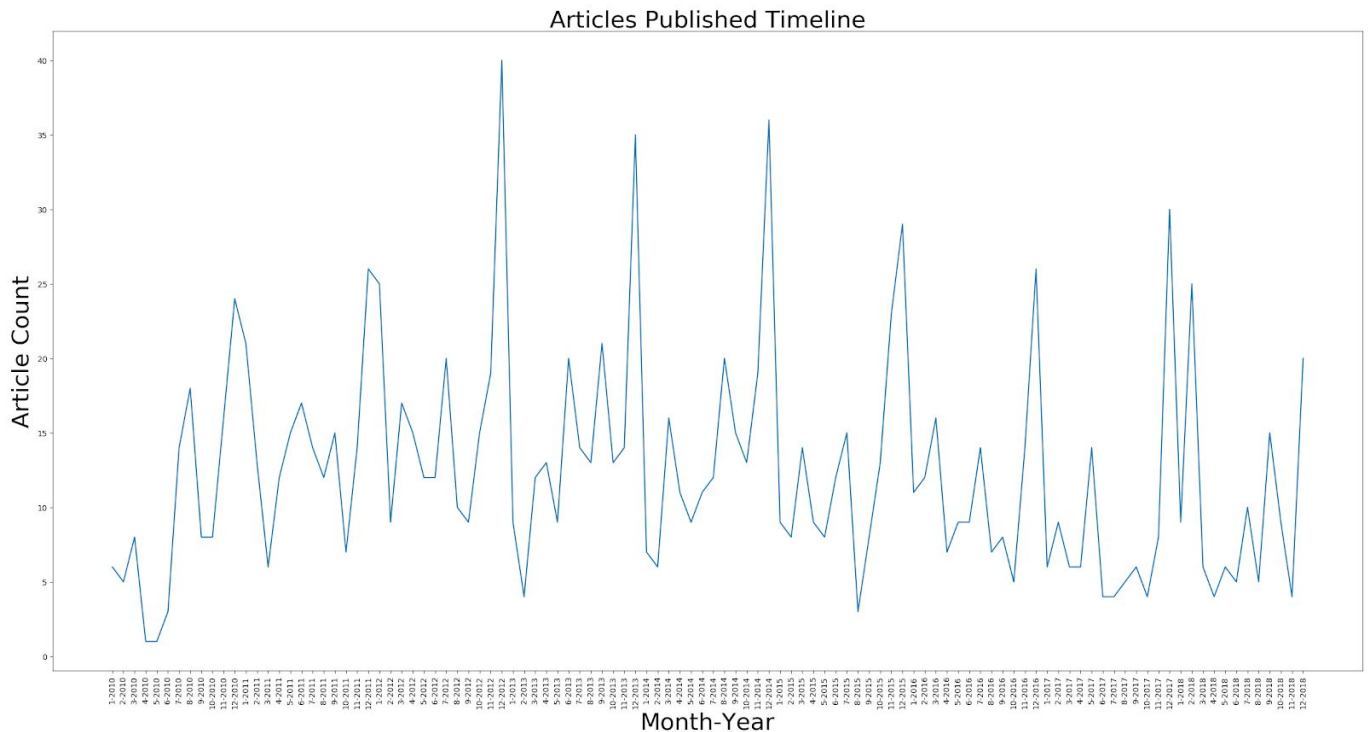
# Data after preprocessing:

new delhi : vaccin protect hiv patient contract tuberculosi tb ultim die final becom realiti seven - year - long trial africa scientist first time develop vaccin succes reduc rate definit tb infect almost 39 among 2000 hiv - infect patient tanzania tb biggest killer hiv - infect patient world india 60 hiv - infect patient weak immun system get infect tb ultim die it scientist dartmouth medic school dms report result clinic trial new vaccin tb - - mycobacterium vacca mv - - januari 29 onlin issu journal aid studi publish march print issu journal princip investig ford von reyn dms said sinc develop new vaccin tb major intern health prioriti especi patient hiv infect tanzanian collabor encourag result studi vaccin type known inactiv whole - cell mycobacteri vaccin expect econom produc distribut von reyn describ trial signific mileston - - first demonstr type vaccin prevent infecti complic hiv adult ad next step improv manufactur method support product larger quantiti tb vaccin need studi subsequ clinic use health ministri offici delhi say vaccin give tremend hope india huge burden tb hiv patient two million hiv - infect indian 10 expect full blown aid everi aid patient 15 chanc everi year develop tb show everi aid patient develop tb time life that s india s latest aid control programm deal tb hiv simultan peopl becom infect tb aid almost alway irrevers formula ministri offici said sinc newli - infect hiv patient risk contract tb almost immedi investig target strategi immun mv patient need start take antiretrovir drug scientif team dartmouth began phase - i human studi mv unit state 1994 demonstr multipl - dose seri mv safe healthi subject patient hiv infect group conduct phase - ii studi larger group adult zambia finland zambian trial research found mv boost immun respons tb first prime childhood current tb vaccin bcg subsequ dardar group receiv nih fund conduct larg phase - iii efficaci trial among hiv - infect patient prior bcg immun tanzania hiv patient particular vulner tb immun system compromis vaccin work boost immun respons patient alreadi given bcg vaccin earlier life

# Analysis of articles:

## Publishing of articles with respect to time:

Plotted is the variations in the number of articles published per month by TOI on HIV. The month which saw the most articles published on HIV was December 2012. The full scale image can be found here: https://github.com/mvedang/Article-Analysis/blob/master/results/articlegraph.png



We observe from this graph that the trend of publishing articles by the TOI on HIV is quite varying and there is no clear increment or decrement in the number of articles published every subsequent month.

## CountVectorizer:

CountVectorizer from scikit-learn library is used to create a table of 1348 rows containing the count of 14980 columns for unique words. The count of each word occuring in an article is recorded in the row dedicated to the article. CountVectorizer is used to indicate the frequency of usage of a particular word in a document, and prepares a table containing the word count in each document.

The top 10 words with the highest occurrence in 1348 articles are:
1) hiv - 8681
2) said - 4388
3) aid - 2856
4) posit - 2579
5) test - 2299
6) infect - 2246
7) patient - 2227
8) peopl - 2052
9) year - 2017
10) state - 1861

The full resolution of the image can be found here:
https://github.com/mvedang/Article-Analysis/blob/master/results/wordcount.png

Usage Of Words

# k-means on count vector table:

The count vector table is clustered using k-means clustering from the scikit-learn library.
The number of clusters selected were 10.

Top terms per cluster:
Cluster 0: hiv said aid peopl posit year patient state test art
Cluster 1: hiv said posit aid year test hospit infect also polic
Cluster 2: hiv test posit said aid peopl women year infect health
Cluster 3: hiv virus infect said cell research drug aid new studi
Cluster 4: children hiv said school posit year parent infect home aid
Cluster 5: patient hiv hospit said treatment centr art posit medic doctor
Cluster 6: blood hiv test said infect bank transfus hospit state posit
Cluster 7: hiv state aid said sex among year case preval worker
Cluster 8: hiv infect peopl aid test india treatment said year new
Cluster 9: tb patient hiv test infect dr said among studi treatment

In the above clusters, we can observe a clustering of topics which are similar to each other.

Let us take Cluster 4:
This cluster contains the words 'children', 'school', 'parent' and 'home' among other stem words.
'children' and 'school' relate to each other as children go to school.
'parent' and 'children' relate to each other as parents are responsible for children.
'parent', 'children', and 'home' are related as parents and children live together at home.

Let us take Cluster 5:
This cluster contains the words 'hospit', 'patient', 'doctor' and 'treatment' among other stem words.
'hospit' and 'doctor' are related as doctors are found in hospitals.
'doctor' and 'patient' are related as doctors treat patients.
'treatment', 'hospital', 'patient', 'doctor' are related as patient goes to hospital to visit doctor for treatment.

Observations:
However, using CountVectorizer for k-means clustering does not always make sense.
We can observe 'hiv' is present in every cluster of words. The reason is that as 'hiv' is mentioned in almost every article, it is treated as an important term.
The flaw in using CountVectorizer for k-means clustering is that the k-means clustering considers a term important based only on the number of occurrences.
By this logic, if we were to include stop words, they would have also appeared in every cluster.
Hence, we must also establish importance of each term based on the importance with respect to the document.

# TF-IDF vector table:

TF-IDF stands for term frequency-inverse document frequency which is used to determine the importance of a word with respect to the document of the corpus.

Term frequency ( x ) = ( total number of times the word x appears in the document ) / ( total number of terms in the document )
Inverse document frequency ( x ) = log ( total number of documents / total number of documents with term x in it )

TF-IDF is calculated as:
TF-IDF = TF(x) * IDF(x)

# k-means on TF-IDF vector table:

The TF-IDF vector table is clustered using k-means clustering from the scikit-learn library.
The number of clusters selected were 10.
We have set the max_df and min_df parameters of the TfidfVectorizer to 0.9 and 0.1 respectively and have taken the 100 most important words to form clusters.

Top terms per cluster:
Cluster 0: polic said case famili posit home year old children report
Cluster 1: test women posit said case district mother aid year state
Cluster 2: research virus infect drug studi said use new india patient
Cluster 3: aid state said sex preval awar among district peopl control
Cluster 4: peopl posit live said aid life year say also famili
Cluster 5: hospit patient doctor said dr medic posit treatment report govern
Cluster 6: blood test hospit infect said report posit doctor govern state
Cluster 7: children school parent said posit home child year live aid
Cluster 8: art patient centr treatment said aid drug peopl medicin test
Cluster 9: woman said hospit posit polic old year medic famili child

In the above clusters, we can observe a clustering of topics which are similar to each other.

Let us take Cluster 6:
This cluster contains the words 'test', 'hospit', 'report', 'govern' and 'state' among other stem words.
'govern' and 'state' relate to each other as state is under the jurisdiction of government.
'report', 'test and 'hospit' relate to each other as hospitals prepare reports of various tests conducted.
'hospit' and 'govern' are related as there are hospitals managed by the government.

Observations:
Here, we observe that even though all the articles are about HIV, and 'hiv' is the most used term in the corpus, no cluster of words contains the term 'hiv'.
Generating a TF-IDF vector table is more efficient as more intricate relationships between words can be identified and such words can be clustered.
Using TF-IDF has proven to be more accurate in many studies.

# Visualization:

A sample of 10% of the dataset is taken for better visual results and to avoid cluttering. The 10% is taken randomly from the dataset.

# Plotting of k-means clusters using TF-IDF vector tables:

Distance between the terms in the TF-IDF vector table is calculated using cosine_similarity from scikit-learn library to calculate cosine similarity.
Distance = 1 - cosine_similarity ( TF-IDF vector table )

Multidimensional Scaling is applied to the distance calculated. It is used to visualize the level of similarity of individual cases in a dataset.

Articles are plotted according to the clusters they are classified into. The colors of each cluster are indicated by the legend. The full scale image of the cluster can be found here:
https://github.com/mvedang/Article-Analysis/blob/master/results/cluster.png



Observations:
The TFIDF cluster displays the similarity between the content of two articles.
The less the distance between two articles, the more similar its contents are.
The more the distance between two article, the less similar its contents are.
We observe that plots of some clusters are very close to plots of other clusters. This demonstrates that these articles may also be quite similar and have been clustered on the basis of a very small characteristic.

# Hierarchical Clustering:

Articles are clustered hierarchically using dendrogram.
Matrix is generated on the distance calculated between terms of TF-IDF vector table using ward in scipy library.
ward is used to perform Ward's linkage on condensed or redundant distance matrix.
Dendrogram is a tree diagram showing arrangement of clusters.

Observations:
Significantly different articles are grouped at the top of the dendrogram. As the hierarchical clustering progresses, articles under one cluster tend to be more similar and are differentiated in the order of less important characteristics of the document.

Let us compare two articles in the hierarchical cluster placed next to each other:

1) Scientists proffer new method to weaken HIV

   Though the symptoms of Acquired Immuno Deficiency Syndrome (AIDS) can be controlled with a cocktail of anti-retroviral drugs, the disease, caused by the Human Immunodeficiency Virus, or HIV, does not have a permanent cure as yet. The treatment of AIDS, however, has got a major boost with scientists from Bengaluru reporting in a study that if the mutation rate of HIV virus is increased by two to six times, it could become ineffective against its host/human body in 10 years. Lead researcher Narendra Dixit, however, cautions that by increasing the mutation rate, such a person will not get cured, but their findings reveal how long it will take for the virus to become non-infectious, and the disease non-progressive.Some of these drugs to increase the mutation levels are already undergoing clinical trials. According to the researchers, a major advantage of these drugs is that they are not susceptible to drug resistance. Hence, Dixit, who is an associate professor at Department of Chemical Engineering, Indian Institute of Science (IISc), Bengaluru, said that in the coming years, with better drug combinations, the treatment window could be reduced further. Their findings have been reported in "Physical Biology" journal.Scientists have so far failed in their efforts to find a permanent cure as HIV virus mutates or changes rapidly, thus making drugs and potential vaccines ineffective. This implies that by the time a drug tries to make the virus ineffective, another mutant that is resistant to the drug, has already developed. Hence, this latest finding is important because it establishes that a threshold level exists for HIV and the virus loses its potency level if its mutation rate is increased by over six times its current or natural mutation rate."One reason why researchers have not been successful in their endeavour to find a permanent cure for HIV/AIDS is because of the high rate at which the virus mutates, causing it to overcome the selection pressures imposed by drugs and potential vaccines. From the virus's point of view, however, there is a flipside to having very high mutation rates. If it increases beyond a tipping point, called error threshold, then it leads to low fitness in the virus population, making it an ineffective pathogen," said Dixit.Prof Dixit and his team have been studying the process of infection caused by HIV and its evolution using certain computer simulation models in the laboratory, which allow them to closely track the fitness levels of the genetically diverse populations within a host and the factors that contribute to it.The finding also has implications for how we treat AIDS, said the research team. High mutation rates can be induced in HIV using a special class of anti-retroviral drugs. If the mutation rate can be raised to the error threshold, then, in principle, the virus can be prevented from causing progressive infection. This research shows that, in fact, such a benign state can be achieved by using mutation-inducing drugs for a period of about 10 years. For their study, the research team used information from the drugs which are currently under trial. They demonstrated in their simulation model that by administering these drugs for 10 years, the replicating HIV will become harmless. This, according to experts, can be an alternative method or strategy to attacking and eventually killing the virus.

2) Kick-and-kill strategy boosts HIV cure effort

MELBOURNE: The elusive quest for an HIV cure received a boost at the world AIDS conference on Tuesday as scientists said they had forced the virus out of a hiding place where it had lurked after being suppressed by drugs. The experiment, carried out with six HIV-infected volunteers, is an important advance in the so-called kick-and-kill approach for a cure, they said. The technique aims to force the Human Immunodeficiency Virus (HIV) from its last redoubt after it is beaten back by antiretroviral drugs. These drugs can bring HIV in the blood to below detectable levels, enabling sick patients to return almost miraculously to normal life. But the therapy has to be taken every day, is costly and carries potential side effects. If the drugs are stopped, HIV usually rebounds within a few weeks and starts once more to infect other immune cells, exposing the body to opportunistic microbes. For the past three years, scientists have focused on ways to kick HIV out of its bolthole and then kill the hideaway cells. In a presentation at the International AIDS Conference in Melbourne, researchers from Aarhus University in Denmark described a step forward in the first stage of this process. Six patients on antiretrovirals took an anti-cancer drug called romidepsin, which prompted virus production in HIV-infected cells to crank up to between 2.1 and 3.9 times above normal. In five patients, the level of virus in the blood increased to measurable levels, an important threshold. agencies

Observations:
In the above two articles, we can observe that the articles focus on the topic of efforts taken to cure HIV. The discussion of drugs being used to take a step towards curing HIV is discussed in the documents. This is the similarity which has hierarchically distinguished both the documents in the last level of the tree, upto which they were treated as documents of the same type.
Hence, we can see that TF-IDF helps us to classify documents with satisfying results.

The hierarchical clustering of 10% of the documents has been carried out, The full scale image can be found here: https://github.com/mvedang/Article-Analysis/blob/master/results/dendrogram.png

['HIV claimed 60,000 less lives last year in India']
['HIV/AIDS sets alarm bells ringing in Chhattisgarh']
['Scientists report breakthrough in reducing risk of acquiring HIV infection']
['\u200bOver 40% living with HIV in India are women']
['In news over HIV, village bans media']
['51-year-old HIV+ man infects thousands in US']
['In 7 years, HIV cases in India dropped by 27%: UN report']
['Scientists proffer new method to weaken HIV']
['Kick-and-kill strategy boosts HIV cure effort']
['From cold to HIV, this drug can fight any viral infection']
['Human trials for HIV vaccine get under way']
['Stock-out hits HIV treatment across India']
['Cipla gets USFDA approval for HIV drug']
['After brief delay, Mysuru gets its quota of anti-HIV drug']
['Hepatitis more lethal than HIV, malaria & dengue']
['State ushers in safe injection norms to fight HIV, hepatitis']
['200 HIV positive babies are born in Mizoram']
['HC reserves order on Junagadh HIV infections']
['CBI to investigate Junagadh kids HIV case']
['HIV from transfusion: Parents of affected kids seek probe']
['HIV+ couple struggle for power connection']
['Karnataka villagers drain lake after HIV+ woman drowns in it']
['Cops escort HIV+ widow home']
['Notice to Punjab on probe into HIV in jails']
['Secunderabad teacher looking for HIV+ bride cheated of Rs 16L']
['HIV+ couple yet to come out of trauma']
['Man ends life after testing positive for HIV']
['No number plate SHIVSHHI Dapoli-Mumbai']
['VRINDAVAN ROAD, SHIVSHISHAKTI COMPLEX, MARKED']
['Anna threatened with HIV+ needles in anonymous letter, FIR lodged']
['HIV + man made four futile suicide attempts after killing his wife, daughters']
['17 HIV+ children break out of home, three staffers arrested']
['Submit report on action taken against lab in false HIV case, says Maharashtra legislative council chairman']
['Legislators to adopt 5 HIV+ kids each']
['Students make docu to ease stigma faced by HIV+ kids']
['Ten-day training for HIV-affected']
['74% male sex workers in Chennai at risk of HIV: Study']
['Day after TOI report on eviction of HIV+ children, landlord takes them back']
['Ministers will adopt HIV-positive children during Belgaum session']
['"Denying HIV positive kids education violates SC notice"']
['Upendra adopts HIV-affected children']
['HIV, HIB: Certificate from lab must for overseas students']
['HIV-positive girl is back in college']
['HIV+ boy driven out of school in Kolkata']
['"Create awareness on HIV, AIDS"']
['Roadshows bust HIV myths']
['HIV Congress awards French countess for service']
['TN first state to develop HIV stigma index']
['Conclave on HIV concludes at Kohima']
['Share HIV-control cost: NACO, state non-committal']
['HIV programme concludes']
['HIV+ people gherao Bihar State AIDS Control Society office']
['DK to roll out single window system for people with HIV']
['Poor HIV patients can travel free to ART centres']
['Programmes to sensitise people about HIV patients soon']
['IMS-BHU takes out HIV/AIDS awareness rally']
['Legal aid centre lends a helping hand to HIV infected patients']
['200 nuns running centre for HIV patients']
['Soya-wheat mix can help protect immunity of HIV patients: Study']
['8% of TB patients in Karnataka are HIV positive, third highest in India']
['State sees decline in HIV positive cases']
['Startups help HIV patients access affordable healthcare']
['GMC reveals identity of HIV+ man, kin in trauma']
['HIV not necessarily a death warrant']
['Devoted to helping HIV+ kids']
['USACS to provide free medicines to all HIV + patients']
['"Health secretary directs release of HIV undertrial's body"']
['Adolescents with HIV will now have a reason to smile']
['"NGO blamed for HIV patient's death"']
['New HIV infections down, but AIDS deaths rise 35% in 3 yrs']
['"57% drop in HIV cases but illiteracy, unprotected sex still biggest hurdles"']
['Over 500 patients in Maharashtra contracted HIV via infected needles in 5 years: RTI']
['"International Women's Day 2018: HIV+ woman who's spreading positivity"']
['Insurance a mirage for poor HIV positives']
['No funds, no condoms, 3L in Maharashtra at HIV risk']
['Project to treat HIV+ drug addicts in northeast to start soon']
['HIV test for babies to start early']
['Maharashtra makes good progress in dealing with HIV']
['Positive news for Mumbai: HIV prevalence dips']
['HIV incidences decline in AP, but so does central funds']
['No baby of HIV +ve moms in last 2 years tested positive']
['No place in morgue freezer for HIV+ kid']
['Lifeline for 10,000 poor HIV patients']
['HIV cases on a decline']
['HIV positive women turn changemakers']
['Most HIV-hit kids lose parents to AIDS']
['Trusts set up auction for benefit of HIV+ children']
['HIV people with positive attitude']
['There are 60 HIV+ kids in Varanasi']
['Fear of having HIV among the most common phobias in Kolkata']
['Warden unleashing HIV attack: Inmate']
['New act addresses discrimination, assists persons living with HIV']
['"Dominic made me promise that I would work for PLHIV as a lawyer"']
['"59% of migrant populace tests HIV positive in Rajasthan"']
['Maharashtra to unveil HIV-friendly private sector workplace policy']
['Thane collector extends aid to poor HIV+ve persons']
['HIV positive rate lowest ever']
['HIV patients left in the dark due to delay in issuing results']
['HIV+ve cases on the rise in Burhanpur']
['HIV+ patients in Palamu to get government aid']
['Over 16,000 HIV cases found in Chhattisgarh in last one decade']
['HIV cases decline in Gurgaon this year']
['Number of HIV cases down in Indore']
['Fresh HIV cases drop 60% in three years']
['"Taking HIV to zero target may now remain a dream"']
['Sharp rise in HIV deaths in Punjab']
['"Guntur HIV girl's condition worsens"']
['"HIV-positive baby's parents pass paternity test, probe given to APSACS"']
['"Kerala girl who 'tested' HIV+ after blood transfusion dies"']
['Angul man suffers unnecessary stigma in HIV test flip flop']
['HIV-positive woman returns to old job as lab technician during deluge']
['AIDS control trust to reach out to HIV+ blood donors']
['Govt to introduce nucleic acid test to detect HIV infection']
['HIV due to transfusion: Kid tests negative in Chennai']
['"Thalassemic boy from Vadodara tests HIV+, laxity alleged"']
['40 normal patients in Ratlam declared HIV+']
['Civil surgeon faces culpable homicide charge as thalassaemic kids contract HIV']
['Woman gets HIV from 'infected' blood']
['CBI court orders further probe in thalassaemic kids contracting HIV']
['CBI court orders further probe in kids contracting HIV']
['In 16 months, 127 patients contracted HIV due to inappropriate blood transfusion']
['27 e-blood banks in state soon to check HIV donors']
['Refusing treatment to HIV+ woman: Inquiry finds doc guilty']
['HIV patient claims Bengaluru hospital denied treatment']
['Cardiac miracle as HIV positive man gets 7 stents']
['PMCH lab report declares HIV positive case as negative']
['Surgeons replace HIV+ man's food pipe with part of stomach']
['"Charlie Sheen Says He's Off HIV Meds in Favor of Mexican Alternative Treatment"']
['Hospital gridlock after HIV panic in Bhopal']
['HIV+ve dumped outside LNJP']
['Probe into HIV+ woman's ordeal']
['HIV positive woman refused delivery at hospital, gives birth to stillborn']
['Doctor, who refused to attend to HIV positive pregnant woman, gets notice']
['HIV+ woman found dead outside SMS']
['Faulty HIV test on pregnant woman at Civil Hospital']