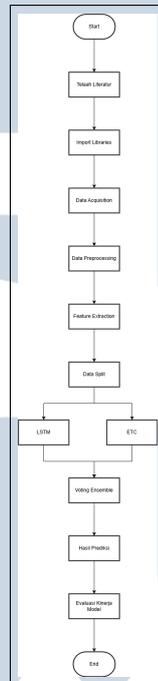


BAB 3 METODOLOGI PENELITIAN

3.1 Alur Penelitian

Alur penelitian merupakan metodologi yang ditampilkan dalam *flowchart* berikut. Langkah alur penelitian secara sistematis dimulai dari studi literatur, *data acquisition* dari *GitHub* penelitian sebelumnya, *data preprocessing*, *feature extraction* dengan *TF-IDF*, *Doc2Vec*, dan *Word2Vec*, *data split*, algoritma *LSTM* dan *ETC*, metode *ensemble learning* dengan cara *soft voting*, dan evaluasi model algoritma.



Gambar 3.1. Flowchart Alur Penelitian

3.2 Metode Penelitian

3.2.1 Telaah Literatur

Telaah literatur dilakukan dengan cara mengumpulkan informasi mengenai penelitian terdahulu yang memiliki hubungan dengan penelitian yang dilakukan. Literatur yang digunakan dalam penelitian ini mengenai lirik lagu, *preprocessing*, algoritma *Extra Trees Classifier*, algoritma *LSTM*, dan *voting classifier*.

3.2.2 Import Libraries

Import Libraries merupakan langkah awal dalam proses pemodelan. *Libraries* merupakan kumpulan *files* yang mengandung *functions* yang bisa digunakan dalam program. *Libraries* dibutuhkan untuk membuat model serta mengevaluasi performa model yang digunakan. Library yang digunakan dalam penelitian adalah *library NLTK (Natural Language Toolkit)*. NLTK berfungsi untuk membantu memproses *NLP (Natural Language Processing)* dalam bahasa *python*.

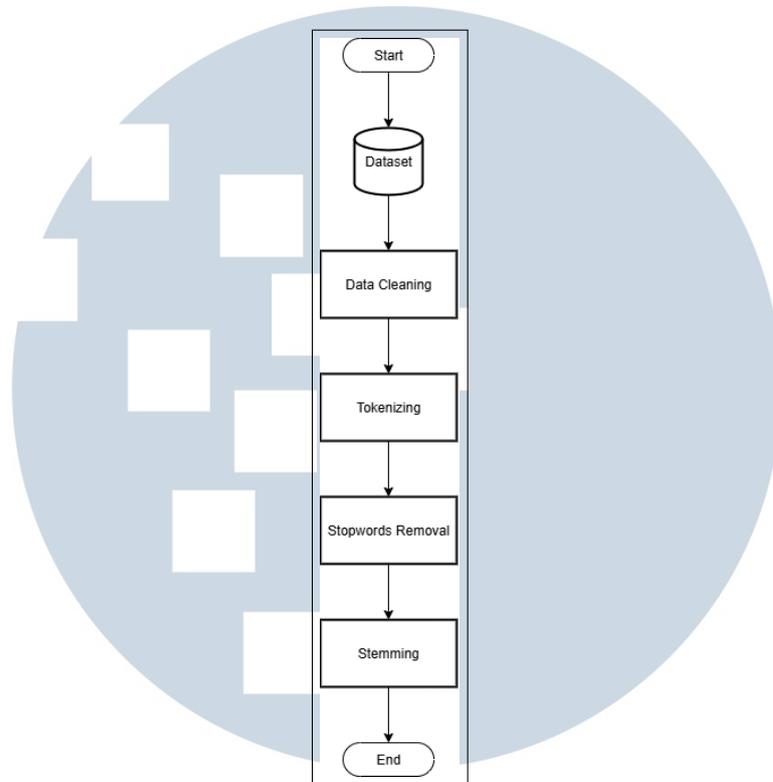
3.2.3 Data Acquisition

Data acquisition merupakan langkah mengunduh *dataset* publik di platform *GitHub*. *Dataset* memiliki kolom seperti nama artis, nama lagu, lirik lagu, dan lagu eksplisit. *Dataset* yang digunakan diambil dari penelitian sebelumnya mengenai deteksi lagu eksplisit dimana terdapat di *GitHub*. Data ini disimpan ke dalam *file* dengan format *Comma Separated Values (CSV)*. *Dataset* memiliki 5566 data dengan 30 atribut, dengan 3 atribut utama yaitu nama lagu (*song*), lirik lagu (*lyrics*), dan klasifikasi lagu eksplisit (*explicit*). Data ini digunakan dalam penelitian untuk menentukan lagu sebagai lagu mengandung konten eksplisit atau tidak.

artist	artist_base	rank	song	year	artist_feat	song_clear	artist_clear	lyrics	acoustic	dnce	abl	energy	explicit	instrumental	key	liveness	loudness	mode	popularity	release_date	speech	tempo	time_signature					
percy	faith	percy	faith	1960	1	theme	from	percy	faith	there's	a	su	0.631	0.466	0.389	0	0.00144	1	0.11	-15.840	1	50	*****	0.0379	81.151	3		
jim	reeves	jim	reeves	1960	2	hell	has	to	im	reves	put	our	sw	0.909	0.554	0.106	0	0.00144	1	0.11	-15.840	1	50	*****	0.0379	81.151	3	
the	everly	the	everly	1960	3	cathy's	clo	the	everly	don't	want			0.412	0.498	0.502	0	0	7	3.72	-8.961	1	50	*****	0.0339	118.609	4	
johnny	pre	johnny	pre	1960	4	running	be	johnny	pre	on	the	bari		0.654	0.772	0.297	0	7.596	06	5	0.125	-14.679	0	40	*****	0.053	118.987	4
mark	dinov	mark	dinov	1960	5	teen	angel	mark	dinov	teen	angel			0.936	0.57	0.0036	0	0	10	0.247	-8.97	1	18	*****	0.0459	101.517	4	
brenda	lee	brenda	lee	1960	6	im	sorry	brenda	lee	im	sorry	oo		0.608	0.529	0.115	0	0.00333	10	0.121	-16.284	1	48	*****	0.0302	101.921	3	
elvis	presle	elvis	presle	1960	7	it's	now	or	elvis	presle	it's	now	or	0.642	0.643	0.491	0	0.00972	4	0.286	-9.312	1	58	*****	0.0344	126.399	4	
jimmy	jane	jimmy	jane	1960	8	handy	mar	jimmy	jane	hey	girls	ga		0.408	0.438	0.659	0	0	10	0.247	-8.97	1	26	*****	0.0497	144.795	4	
elvis	presle	elvis	presle	1960	9	stuck	on	you	can	st				0.758	0.647	0.513	0	0.046	06	7	0.158	-12.372	1	50	*****	0.0421	131.841	4
chubby	che	chubby	che	1960	10	the	twist	chubby	che	come	on	b		0.199	0.526	0.633	0	1.24E-06	4	0.0709	-7.119	1	52	*****	0.0322	156.405	4	
connie	frae	connie	frae	1960	11	everybody		connie	frae	the	tears	i		0.626	0.589	0.643	0	3.62E-05	1	0.237	-8.118	1	37	*****	0.0334	84.306	4	
holly	lyde	holly	lyde	1960	12	wild	one	holly	lyde	oh	wild	one		0.709	0.599	0.775	0	0	7	0.0904	-4.566	1	24	*****	0.0466	148.507	4	
the	brother	the	brother	1960	13	greenfield		the	brother	once	there			0.652	0.485	0.182	0	0	11	0.112	-18.083	0	41	*****	0.0434	118.600	1	
jack	scott	jack	scott	1960	14	what	in	the	jack	scott	what	in	the	0.602	0.506	0.176	0	3.93E-06	4	0.0978	-13.962	1	31	*****	0.0303	76.273	4	
many	robb	many	robb	1960	15	el paso		many	robb	out	in	the	h	0.636	0.654	0.452	0	2.89E-05	2	0.16	-9.709	1	58	*****	0.03	106.662	3	
the	hollywe	the	hollywe	1960	16	alleyoop		the	hollywe	there's	a	m		0.798	0.65	0.352	0	4.10E-06	2	0.114	-10.304	1	10	*****	0.052	64.596	4	
connie	frae	connie	frae	1960	17	my	heart	h	connie	frae	i	old	this	h	0.605	0.338	0.551	0	0	6	0.316	-5.995	1	19	*****	0.0465	108.045	4
brenda	lee	brenda	lee	1960	18	sweet	noth	brenda	lee	alright	my			0.785	0.787	0.389	0	0.00E-05	5	0.158	-12.824	1	37	*****	0.0683	125.216	4	
brian	hyar	brian	hyar	1960	19	itsy	bity	te	brian	hyar	she	was	af	0.682	0.795	0.4	0	0	2	0.0326	-13.01	1	46	*****	0.0749	122.174	4	
roy	orbis	roy	orbis	1960	20	only	the	lor	roy	orbis	dum	dum	c	0.377	0.57	0.529	0	0.00509	5	0.203	-10.769	1	60	*****	0.028	123.273	4	
alon	and	alon	and	1960	21	where	or	a	alon	and	it	seems	wa	0.797	0.454	0.271	0	0	8	0.126	-11.656	1	36	*****	0.0263	110.89	3	
connie	frae	connie	frae	1960	22	where	is	connie	frae	the	saddest	rea		0.865	0.339	0.406	0	1.10E-05	10	0.11	-8.955	1	30	*****	0.0319	109.763	3	
paul	anka	paul	anka	1960	23	puppy	love	paul	anka	and	they	cc		0.662	0.29	0.408	0	0	7	0.272	-9.607	1	46	*****	0.0318	102.229	3	

Gambar 3.2. Import Dataset

3.2.4 Data Preprocessing



Gambar 3.3. Preprocessing Database

Preprocessing merupakan langkah untuk mengubah data menjadi data bersih sehingga data dapat dianalisis. Berikut proses *preprocessing* yang dilakukan.

1. *Data Cleaning*

Data cleaning merupakan proses membersihkan data dimana elemen yang tidak dibutuhkan atau relevan dengan data yang ingin diproses dihilangkan.

2. *Tokenizing*

Tokenizing merupakan proses memisahkan teks menjadi kata tunggal atau frasa untuk lebih mudah melakukan proses terhadap data.

3. *Normalization*

Normalization merupakan proses mengubah kata teks yang tidak baku menjadi kata baku.

4. *Stopword Removal*

Stopword Removal merupakan proses menghapus kata yang tidak memiliki makna penting seperti "yang", "di", atau "dan".

5. *Stemming*

Stemming merupakan proses mengubah kata imbuhan menjadi kata baku tanpa imbuhan.

3.2.5 Feature Extraction

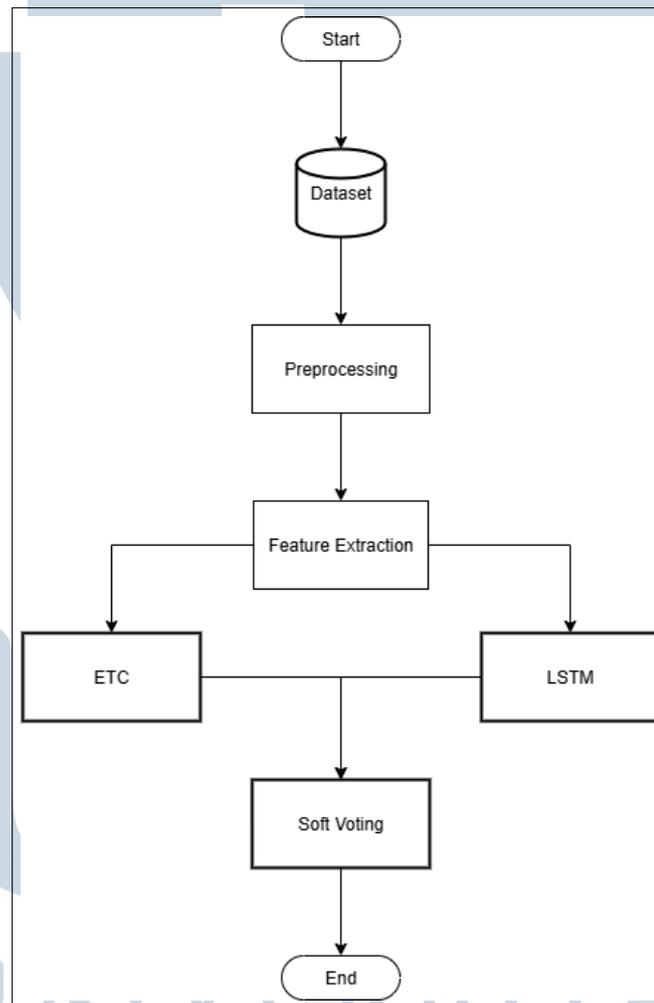
Feature extraction adalah proses yang digunakan untuk mengubah kata dalam teks menjadi nilai *integer* atau *float*. Penelitian ini menggunakan *TF-IDF Vectorizer* dan *Doc2Vec* untuk *machine learning* dan *Word2Vec* untuk *deep learning*. Dalam *feature extraction TF-IDF*, *TF* digunakan dalam mengukur kata muncul dalam dokumen, semakin sering kata muncul maka nilai *TF* akan semakin tinggi. *IDF* merupakan pengukur kata yang dipengaruhi oleh seberapa umum atau langka kata dalam dokumen. Kombinasi *TF* dan *IDF* akan menghasilkan vektor numerik yang menjadi representasi kata dalam dokumen. *Doc2Vec* merupakan *feature extraction* yang bekerja dengan menggunakan *Distributed Memory (DM)* dan *Distributed Bag of Words (DBOW)*. *DM* berfungsi dalam mempertimbangkan konteks dari kata dan ID dokumen sebagai input prediksi kata selanjutnya. *DBOW* berfungsi untuk memprediksi kata dari dokumen dengan menggunakan vektor dokumen sebagai konteks. *Word2Vec* merupakan *feature extraction* yang bekerja dengan menggunakan *Continuous Bag of Words (CBOW)* dan *Skip-gram*. *CBOW* berfungsi dalam memprediksi kata target dengan kata konteks di sekitar. *Skip-gram* berfungsi dalam memprediksi kata konteks dengan kata target.

3.2.6 Algoritma Klasifikasi

Penelitian ini menggunakan algoritma *Extra Trees Classifier* dan *LSTM* untuk mencari klasifikasi terbaik dalam analisis kata seksual dalam lirik lagu. Tahap ini dilakukan untuk mencari nilai hasil dari akurasi algoritma terbaik untuk mengukur kinerja model dari *confusion matrix* dalam melakukan klasifikasi kata seksual dalam lirik lagu. Tahapan yang dilakukan adalah *data acquisition*, *data preprocessing*, *feature extraction* dengan menggunakan *TF-IDF* dan *Doc2Vec* untuk algoritma *ETC* dan *Word2Vec* untuk algoritma *LSTM*, *data split*, *modeling ETC*, dan evaluasi. Di awal, dilakukan *preprocessing* dengan *data cleaning*, *tokenizing*, *normalization*, *stopword removal*, dan *stemming* untuk *dataset*. Lalu data akan dibuat menjadi nilai dengan *feature extraction* dengan *TF-IDF*, *Word2Vec*, dan *Doc2Vec*. Tahap berikutnya adalah melakukan *modeling* dengan menggunakan

ETC dan *LSTM*, dimana digunakan *hyperparameter* dari penelitian sebelumnya dan menggunakan *cross validation* untuk mencari rata-rata akurasi model. Kemudian dilakukan evaluasi *confusion matrix* dari setiap model. Performa model dilihat dari nilai tabel *confusion matrix*.

3.2.7 Metode Ensemble Learning



Gambar 3.4. Metode Ensemble Learning

Ensemble Learning merupakan proses pembelajaran dimana menggabungkan beberapa model dasar untuk mendapatkan model baru dengan hasil yang lebih baik dibandingkan model dasar. Algoritma yang digunakan dalam *ensemble learning* adalah algoritma *LSTM* dan *ETC*. Metode *ensemble learning* yang digunakan dalam penelitian sendiri menggunakan *soft voting* dimana bekerja

dengan cara mengambil hasil dari gabungan model algoritma dan mengambil rata-rata dari probabilitas semua model yang digunakan untuk menghasilkan hasil akhir. Dalam penelitian, *soft voting* yang digunakan merupakan *ETC + LSTM* dengan *soft voting* pertama menggunakan *feature extraction TF-IDF* dalam *ETC* dan *soft voting* kedua menggunakan *feature extraction Doc2Vec* dalam *ETC*.

A K-Fold Cross Validation

K-fold cross validation merupakan metode validasi yang membagi dataset menjadi beberapa bagian (*fold*) dengan ukuran yang sama. Dalam penelitian ini digunakan *10-fold cross validation*, di mana dataset dibagi menjadi 10 bagian. Setiap iterasi menggunakan satu *fold* sebagai *testing set* dan sembilan *fold* lainnya sebagai *training set*. Proses ini dilakukan sebanyak 10 kali untuk memastikan generalisasi model terhadap data yang berbeda pada setiap iterasi.

3.2.8 Testing Model

Setelah model algoritma dibangun menggunakan pendekatan *ensemble voting*, langkah selanjutnya adalah melakukan pengujian terhadap model tersebut. Pengujian dilakukan dengan menggunakan data berupa lirik lagu, dan hasil dari prediksi akan menunjukkan apakah lirik tersebut tergolong *Explicit* atau *Non Explicit*, sesuai dengan klasifikasi model.

3.2.9 Hasil Prediksi

Hasil prediksi merupakan visualisasi evaluasi dari pendekatan *ensemble learning* yang terdiri dari dua model *soft voting* dengan teknik *feature extraction* yang berbeda. Visualisasi ini bertujuan untuk membandingkan kinerja model berdasarkan representasi fitur yang digunakan serta menilai akurasi klasifikasi dalam mendeteksi konten eksplisit pada lirik lagu.