# SUNUM İÇİN HAZIRLANMIŞ SLAYTTAN EKRAN GÖRÜNTÜLERİ

# Giriş

### Problem Tanımı

- -Sosyal medya ve dijital platformlarda nefret söyleminin yaygınlığı artmaktadır.
- -Nefret söylemi; bireyler ve topluluklar üzerinde ciddi psikolojik ve sosyal etkiler yaratabilir.
- -Bu projede, nefret söylemini otomatik olarak tespit eden bir model geliştirilmektedir.

### Amaç

Nefret söylemini, hakaret içeren söylemleri ve tarafsız ifadeleri ayırabilen bir sistem oluşturmak.

Derin öğrenme yaklaşımlarıyla başarımlı bir sınıflandırma modeli geliştirmek

## Kullanılan Veri Seti

- -Kaynak: Kaggle hate-speech-and-offensive-language-master
- -Veri Seti İçeriği:
  - -Toplam 25.000+ tweet.
  - -İçerikler 3 sınıfa ayrılmıştır:
    - -Nefret Söylemi (Hate Speech)
    - -Hakaret (Çirkin Dil) (Offensive)
    - -Tarafsız (Çoklu Sınıf Neutral)
- -Sınıf Dağılımı dengelenmiştir (oversampling kullanılmıştır).

# Yöntem ve Teknolojiler

BERT (Bidirectional Encoder Representations from Transformers) modeli Model

Transformers kütüphanesi yardımıyla BERT tümleştirildi.

Metin temizleme: Metin İşleme

-URL'ler, özel karakterler ve büyük/küçük harf farklılıkları giderildi.

Tokenization: BERT tokenizer kullanıldı.

Modelleme Model, 10 epoch boyunca eğitildi(Son versiyon için(7'de zorla durduruldu)).

AdamW optimizasyonu ve linear scheduler kullanıldı.

## Koddaki Önemli Noktalar

BERT Modelinin Kullanımı:

- Transformer tabanlı modellerin gücü, metinlerin anlamını anlamada ve sınıflandırmada önemli bir rol oynar. BertForSequenceClassification bu projeyi özel kılıyor.
- Oversampling ve Ağırlıklı Kayıp Fonksiyonu: 2 Veri dengesizliğini dengelemek için kullanılan oversampling ve sınıf ağırlıklarının kayıp fonksiyonuna eklenmesi, modelin nadir sınıfları öğrenmesini sağlıyor.
- Erken Durdurma ve Model Kaydetme: 3 Eğitim sırasında aşırı öğrenmenin (overfitting) önlenmesi ve en iyi modelin kaydedilmesi, uygulama açısından kritik.

# Sonuçlar

### Model Performansi

Eğitim ve Doğrulama Doğrulukları

En iyi epoch sonucunda doğrulama doğruluğu: 96.37%

Ortalama eğitim kaybı: 0.1520

Modelin tahmin doğruluğu, normalize edilmiş karışıklık matrisi ile görsel olarak sunuldu.

### Karışıklık Matrisi Sınıf Performansı

Precision, Recall ve F1-Score

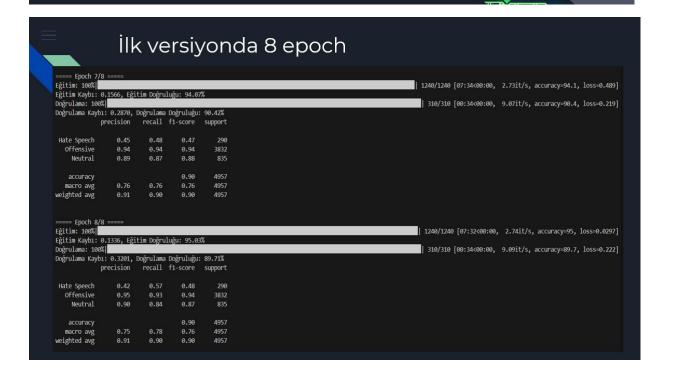
Hate Speech: F1-Score: 0.95

Offensive: F1-Score: 0.93

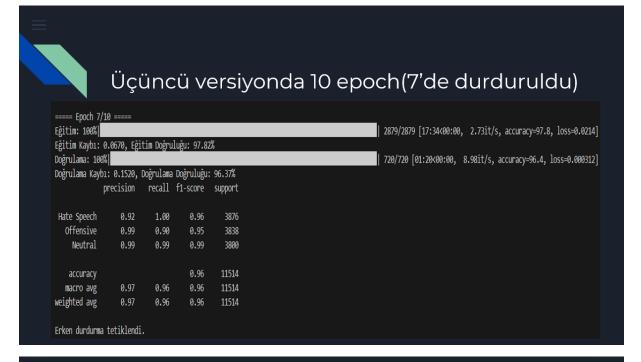
Neutral: F1-Score: 0.98



# Celecek Çalışmalar Veri seti daha fazla dil ve kültüre genelleştirilebilir. Daha karmaşık modeller (örneğin, BERT Large veya RoBERTa) kullanılabilir. Gerçek zamanlı analiz için optimizasyon yapılabilir.



### İkinci versiyonda 16 epoch (9'da durduruldu) Eğitim: 100% | Eğitim Doğruluğu: 97.44% | Doğrulama: 100% | Poğrulama: 100% | Eğitim Doğruluğu: 97.44% | Doğrulama: 100% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim Doğruluğu: 07.44% | Eğitim 2098/2098 [12:43<00:00, 2.75it/s, accuracy=97.4, loss=0.00462] | 525/525 [00:57<00:00, 9.09it/s, accuracy=92.5, loss=0.00212] Doğrulama Kaybı: 0.3014, Doğrulama Doğruluğu: precision recall f1-score Hate Speech Offensive 3679 3870 840 Neutral accuracy | 2098/2098 [12:43<00:00, 2.75it/s, accuracy=97.9, loss=0.00141] Eğitim Kaybı: 0.0646, Eğitim Doğruluğu: 97.93% Doğrulama: 100% Doğrulama Doğrulama Doğruluğu: | 525/525 [00:57<00:00, 9.09it/s, accuracy=93.1, loss=0.00301] 0.94 0.93 0.87 3679 3870 840 Hate Speech Offensive Neutral 8389 8389



# Genel Sonuç

Sonuç olarak, denenmiş 7 değiştirilmiş hiperparametreli örneklerden sonra ulaşılan en iyi,örnek;

Eğitim Kaybı: 0.0670, Eğitim Doğruluğu: 97.82%

Doğrulama Kaybı: 0.1520, Doğrulama Doğruluğu: 96.37%

değerlerine sahiptir.

# İŞLEME İÇİN KULLANILAN KODLAR VE PARÇALANMIŞ İŞ BÖLÜMLERİ

```
# Gerekli kütüphaneleri yükleyin
import pandas as pd
import numpy as np
import nltk
import re
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt
import torch
from torch.utils.data import DataLoader, Dataset
import torch.nn as nn
from transformers import BertTokenizer, BertForSequenceClassification, get_scheduler
from tqdm import tqdm
```

```
nltk.download('stopwords')
from nltk.corpus import stopwords
data path = "C:\\Users\\mamie\\Downloads\\hate-speech-and-offensive-language-master\\data\\labeled data.csv"
df = pd.read_csv(data_path)
def clean_text(text):
    text = text.lower()
    text = re.sub(r'http\S+', '', text) # URL'leri kaldır
    text = re.sub(r'[^a-zA-Z\s]', '', text) # Özel karakterleri kaldır
    text = text.strip()
    return text
df['tweet'] = df['tweet'].apply(clean text)
# Sınıf dağılımını dengeleme (oversampling)
class counts = df['class'].value counts()
df balanced = pd.concat([df[df['class'] == cls].sample(class counts.max(), replace=True) for cls in class counts.index])
df = df_balanced.reset_index(drop=True)
texts = df['tweet'].values
labels = df['class'].values
tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
```

```
# Dataset sinifi
class HateSpeechDataset(Dataset):
   def __init__(self, texts, labels, tokenizer, max_len):
       self.texts = texts
       self.labels = labels
       self.tokenizer = tokenizer
       self.max len = max len
   def __len__(self):
       return len(self.texts)
   def __getitem__(self, item):
        text = str(self.texts[item])
       label = self.labels[item]
        encoding = self.tokenizer.encode_plus(
           text,
           add special tokens=True,
           max length=self.max len,
           return token type ids=False,
           padding="max_length",
           truncation=True,
           return_attention_mask=True,
           return_tensors='pt',
            'input_ids': encoding['input_ids'].squeeze(0),
            'attention mask': encoding['attention mask'].squeeze(0),
            'label': torch.tensor(label, dtype=torch.long)
```

# Parametreler

```
EPOCHS = 10
MAX LEN = 128
BATCH SIZE = 16
LEARNING_RATE = 2e-5
X_train, X_test, y_train, y_test = train_test_split(texts, labels, test_size=0.2, random_state=42)
train_dataset = HateSpeechDataset(X_train, y_train, tokenizer, MAX_LEN)
test_dataset = HateSpeechDataset(X_test, y_test, tokenizer, MAX_LEN)
train_loader = DataLoader(train_dataset, batch_size=BATCH_SIZE, shuffle=True)
test loader = DataLoader(test dataset, batch_size=BATCH_SIZE)
# Model vükleme
model = BertForSequenceClassification.from_pretrained("bert-base-uncased", num_labels=3, hidden_dropout_prob=0.4)
device = torch.device("cuda" if torch.cuda.is available() else "cpu")
model = model.to(device)
class_counts = df['class'].value_counts().sort_index().values
class_weights = torch.tensor([sum(class_counts) / c for c in class_counts]).to(device)
loss fn = nn.CrossEntropyLoss(weight=class weights)
optimizer = torch.optim.AdamW(model.parameters(), lr=LEARNING_RATE, weight_decay=0.01)
scheduler = get_scheduler("linear", optimizer=optimizer, num_warmup_steps=0, num_training_steps=len(train_loader) * EPOCHS)
best_val_loss = float('inf')
early_stopping_patience = 3
early_stopping_counter = 0
```

```
def train_epoch(model, data_loader, loss_fn, optimizer, device, scheduler):
   model.train()
   total loss = 0
   correct = 0
    total = 0
    loop = tqdm(data_loader, leave=True, desc="Eğitim")
    for batch in loop:
        input_ids = batch['input_ids'].to(device)
        attention_mask = batch['attention_mask'].to(device)
        labels = batch['label'].to(device)
        optimizer.zero grad()
        outputs = model(input ids=input ids, attention mask=attention mask, labels=labels)
        loss = outputs.loss
        total_loss += loss.item()
       preds = torch.argmax(outputs.logits, dim=1)
        correct += (preds == labels).sum().item()
        total += labels.size(0)
        loss.backward()
        optimizer.step()
        scheduler.step()
        loop.set_postfix(loss=loss.item(), accuracy=100 * correct / total)
   return total_loss / len(data_loader), 100 * correct / total
```

```
def eval model(model, data loader, loss fn, device):
   model.eval()
   total loss = 0
   correct = 0
   total = 0
   all_preds = []
   all_labels = []
   loop = tqdm(data loader, leave=True, desc="Doğrulama")
   with torch.no_grad():
       for batch in loop:
            input_ids = batch['input_ids'].to(device)
            attention_mask = batch['attention_mask'].to(device)
           labels = batch['label'].to(device)
           outputs = model(input_ids=input_ids, attention_mask=attention_mask, labels=labels)
            loss = outputs.loss
           total_loss += loss.item()
            preds = torch.argmax(outputs.logits, dim=1)
            correct += (preds == labels).sum().item()
            total += labels.size(0)
            all_preds.extend(preds.cpu().numpy())
            all_labels.extend(labels.cpu().numpy())
            # İlerleme çubuğunu güncelle
            loop.set_postfix(loss=loss.item(), accuracy=100 * correct / total)
   return total_loss / len(data_loader), 100 * correct / total, all_preds, all_labels
```

```
accuracies = []
history = {'train_acc': [], 'val_acc': []}
for epoch in range(EPOCHS):
    print(f"\n==== Epoch {epoch + 1}/{EPOCHS} =====")
    train_loss, train_acc = train_epoch(model, train_loader, loss_fn, optimizer, device, scheduler)
   print(f"Eğitim Kayba: {train loss:.4f}, Eğitim Doğruluğu: {train acc:.2f}%")
   val_loss, val_acc, val_preds, val_labels = eval_model(model, test_loader, loss fn, device)
   print(f"Doğrulama Kaybı: {val_loss:.4f}, Doğrulama Doğruluğu: {val_acc:.2f}%")
    accuracies.append((epoch + 1, val_acc))
   history['train_acc'].append(train_acc)
   history['val acc'].append(val acc)
   print(classification report(val labels, val preds, target names=['Hate Speech', 'Offensive', 'Neutral']))
   # Early Stopping kontrolü
   if val_loss < best_val_loss:</pre>
       best val loss = val loss
        early stopping counter = 0
        torch.save(model.state_dict(), "best_model.pt")
        early_stopping_counter += 1
    if early stopping counter >= early stopping patience:
        print("Erken durdurma tetiklendi.")
        break
```

```
model.load_state_dict(torch.load("best_model.pt"))
def plot_confusion_matrix(preds, labels):
    cm = confusion_matrix(labels, preds, normalize='true')
    plt.figure(figsize=(8, 6))
    sns.heatmap(cm, annot=True, fmt='.2f', cmap='Blues', xticklabels=['Hate Speech', 'Offensive', 'Neutral'], yticklabels=['Hate Speech',
    'Offensive', 'Neutral'])
plt.xlabel('Tahmin Edilen')
    plt.ylabel('Gerçek')
    plt.title('Karişiklik Matrisi')
    plt.show()
val_loss, val_acc, val_preds, val_labels = eval_model(model, test_loader, loss_fn, device)
print(f"Test Kayba: {val loss:.4f}, Test Doğruluğu: {val acc:.2f}%")
plot_confusion_matrix(val_preds, val_labels)
def plot_history(history):
    plt.figure(figsize=(10, 6))
    plt.plot(history['train_acc'], label='Eğitim Doğruluğu')
    plt.plot(history['val_acc'], label='Doğrulama Doğruluğu')
    plt.xlabel('Epochs')
plt.ylabel('Doğruluk (%)')
    plt.title('Model Eğitimi ve Doğrulama Doğruluğu')
    plt.legend()
    plt.show()
plot history(history)
```