

# Fine-Tuning DistilBERT for Legal Case Prediction

A Step-by-Step Guide with Practical Example

Presented by: Alisha Rath

# Introduction to DistilBERT

- What is DistilBERT?
  - A distilled version of BERT (Bidirectional Encoder Representations from Transformers).
  - Lightweight, faster, and retains 97% of BERT's performance.
  - Uncased: Ignores case sensitivity (e.g., 'law' and 'LAW' are treated the same).
- Why Distillation?
  - Reduces model size and training time.
  - Ideal for resource-constrained environments.

# Importance of Fine-Tuning for Legal Case Prediction

- Why Fine-Tuning?
  - Pre-trained models are trained on general language corpora.
  - Fine-tuning adapts the model to legal-specific language and prediction tasks.
- Use Case: Predicting legal case outcomes.
  - Predict outcomes of court cases based on legal documents, charges, and evidence.

# Understanding DistilBERT's Architecture for Legal Text

- How DistilBERT Works:
  - Transformer architecture with attention mechanisms.
  - Processes legal texts and understands context.
- Why Use DistilBERT?
  - Faster inference for handling long legal documents.
  - Suitable for scaling to large legal datasets.

# Pre-Training vs. Fine-Tuning

- Pre-Training:
  - Trained on general language corpora.
- Fine-Tuning:
  - Adapts the pre-trained model for specific tasks like legal case prediction.
  - Requires a labeled legal dataset to learn task-specific knowledge.

# Dataset for Fine-Tuning DistilBERT

- Dataset Example: Supreme Court Judgment Prediction Dataset (Kaggle).
  - Contains data on legal cases (charges, descriptions, outcomes).
- Data Preprocessing:
  - Clean legal text.
  - Tokenize using DistilBERT's tokenizer.

# Example: Fine-Tuning DistilBERT for Legal Case Prediction

- Steps:
  1. Load the pre-trained DistilBERT model using Hugging Face.
  2. Tokenize legal documents.
  3. Fine-tune the model to predict case outcomes.

# Code Example: Loading DistilBERT

```
from transformers import DistilBertTokenizer,  
DistilBertForSequenceClassification, Trainer, TrainingArguments  
  
# Load pre-trained DistilBERT tokenizer and model  
tokenizer = DistilBertTokenizer.from_pretrained('distilbert-base-  
uncased')  
model =  
DistilBertForSequenceClassification.from_pretrained('distilbert-base-  
uncased', num_labels=2)  
  
# Tokenize dataset  
train_encodings = tokenizer(train_texts, truncation=True,  
padding=True, max_length=512)  
test_encodings = tokenizer(test_texts, truncation=True, padding=True,  
max_length=512)
```



# Code Example: Fine-Tuning with Trainer API

```
# Define training arguments
training_args = TrainingArguments(
    output_dir='./results',
    num_train_epochs=3,
    per_device_train_batch_size=16,
    per_device_eval_batch_size=16,
    warmup_steps=500,
    weight_decay=0.01,
    logging_dir='./logs',
    logging_steps=10,
)

# Trainer for fine-tuning
trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=train_dataset,
    eval_dataset=eval_dataset
)

# Fine-tune the model
trainer.train()
```

# Model Evaluation for Legal Case Prediction

- Evaluating the Model:
  - Predict the outcome of unseen legal cases.
  - Assess performance metrics: accuracy, precision, recall, F1-score.

- Example Code for Evaluation:

```
predictions = trainer.predict(test_dataset)
preds = np.argmax(predictions, axis=-1)
accuracy = np.mean(preds == test_labels)
print(f"Test Accuracy: {accuracy:.4f}")
```

# Challenges in Fine-Tuning Legal Models

- Domain-Specific Language:
  - Legal texts contain complex language.
- Imbalanced Datasets:
  - Legal outcomes may be skewed.
- Data Privacy:
  - Legal datasets may include sensitive data.

# Key Considerations for Fine-Tuning Legal Models

- Hyperparameter Tuning:
  - Optimize model performance.
- Preprocessing Legal Text:
  - Clean and tokenize documents.
- Transfer Learning Challenges:
  - Domain shift between general data and legal texts.

# Conclusion

- Fine-tuning DistilBERT adapts it for legal case prediction.
- Next Steps:
  - Optimize with more datasets and tasks (e.g., verdict prediction).
  - Explore real-world deployment in legal firms.