

DATA COLLECTION AND LABELING

CS 203: Software Tools and Techniques for AI

Prof. Nipun Batra



MODULE OVERVIEW



FOUR CORE COMPONENTS

- 1. Data Collection** - Tools and techniques for gathering data from diverse sources
- 2. Data Validation** - Ensuring data quality and reliability
- 3. Data Labeling** - Annotating datasets with ground truth
- 4. Data Augmentation** - Expanding datasets strategically



Key Insight

Quality data is the foundation of successful AI systems. Garbage in, garbage out!

PART 1: DATA COLLECTION

Instrumenting and Logging Data Sources



WHY DATA COLLECTION MATTERS

- Real-world AI **systems** depend on continuous data flow
- User behavior changes over time → models need fresh data
- Production **systems** require automated collection pipelines
- Debugging often requires understanding what data was seen

COMMON DATA SOURCES

DIGITAL SOURCES

- Web applications
- Mobile apps
- IoT devices
- APIs and databases

PHYSICAL SOURCES

- Sensors
- Cameras
- Microphones
- Manual entry



INSTRUMENTATION: THE FOUNDATION

Instrumentation = Adding code to collect data about system behavior

KEY PRINCIPLES

- 1. Minimal Performance Impact** - Don't slow down production systems
- 2. Comprehensive Coverage** - Capture all relevant events
- 3. Privacy-Aware** - Respect user privacy and regulations (GDPR, CCPA)
- 4. Structured Logging** - Consistent formats for easy parsing



BASIC INSTRUMENTATION EXAMPLE

```
1 import logging
2 import json
3 from datetime import datetime
4
5 def log_user_action(user_id, action, metadata):
6     event = {
7         "timestamp": datetime.utcnow().isoformat(),
8         "user_id": user_id,
9         "action": action,
10        "metadata": metadata
11    }
12    logging.info(json.dumps(event))
```



WEB ANALYTICS TOOLS

GOOGLE ANALYTICS 4

```
1 // Track custom events
2 gtag('event', 'search', {
3   'search_term': query,
4   'results_count': results.length
5 });
```

Pros: Free, Rich ecosystem,
Real-time dashboards **Cons:**
Privacy concerns, Limited
customization

MIXPANEL

```
1 // Track with properties
2 mixpanel.track('Video Played', {
3   'video_id': video.id,
4   'duration': video.length,
5   'quality': '1080p'
6 });
```

Pros: Detailed analytics,
Cohort analysis **Cons:** Costly
at scale



SELF-HOSTED ANALYTICS

WHY SELF-HOST?

- **Data ownership** - Complete control over your data
- **Privacy compliance** - GDPR-friendly by design
- **No tracking cookies** - Lightweight and fast
- **Cost-effective** - Predictable infrastructure costs



PLAUSIBLE & UMAMI EXAMPLES

```
1 // Plausible: Simple, privacy-focused
2 <script defer data-domain="yourdomain.com"
3         src="https://plausible.io/js/script.js"></script>
4
5 plausible('signup', {props: {plan: 'premium'}});
```

```
1 // Umami: Open-source alternative
2 <script async src="https://analytics.yourdomain.com/umami.js"
3         data-website-id="your-website-id"></script>
4
5 umami.track('button-click', {button: 'subscribe'});
```



APPLICATION PERFORMANCE MONITORING

SENTRY

```
1 import sentry_sdk
2
3 sentry_sdk.init(
4     dsn="your-dsn",
5     traces_sample_rate=1.0
6 )
7
8 try:
9     process_data()
10 except Exception as e:
11     sentry_sdk.capture_exception(e)
```

DATADOG

```
1 from datadog import initialize, st
2
3 initialize(
4     api_key='your-key',
5     app_key='your-app-key'
6 )
7
8 statsd.increment('api.requests')
9 statsd.histogram('api.latency', 24
```



DATABASE EVENT LOGGING

```
1 -- PostgreSQL: Track all changes to users table
2 CREATE TABLE user_audit (
3     audit_id SERIAL PRIMARY KEY,
4     user_id INT,
5     action VARCHAR(10),
6     old_data JSONB,
7     new_data JSONB,
8     changed_at TIMESTAMP DEFAULT NOW()
9 );
10
11 CREATE OR REPLACE FUNCTION audit_user_changes()
12 RETURNS TRIGGER AS $$$
13 BEGIN
14     IF (TG_OP = 'UPDATE') THEN
15         INSERT INTO user_audit(user_id, action, old_data, new_data)
16             VALUES (NEW.id, 'UPDATE', row_to_json(OLD), row_to_json(NEW));
17     END IF;
18     RETURN NEW;
19 END;
```



CHANGE DATA CAPTURE (CDC)

CDC = Capturing and streaming database changes in real-time

DEBEZIUM EXAMPLE

```
1 {
2   "name": "inventory-connector",
3   "config": {
4     "connector.class": "io.debezium.connector.postgresql.PostgresConnector"
5     "database.hostname": "postgres",
6     "database.port": "5432",
7     "table.include.list": "public.orders,public.customers",
8     "topic.prefix": "dbserver1"
9   }
10 }
```

Benefits: Real-time streaming, No app changes, Complete history



OPENTELEMETRY: INDUSTRY STANDARD

```
1 from opentelemetry import trace
2 from opentelemetry.sdk.trace import TracerProvider
3
4 trace.set_tracer_provider(TracerProvider())
5 tracer = trace.get_tracer(__name__)
6
7 @tracer.start_as_current_span("process_request")
8 def process_request(user_id):
9     span = trace.get_current_span()
10    span.set_attribute("user.id", user_id)
11    # Your logic here
```



WEB SCRAPING WITH SCRAPY

```
1 import scrapy
2
3 class ProductSpider(scrapy.Spider):
4     name = 'products'
5     start_urls = ['https://example.com/products']
6
7     custom_settings = {
8         'DOWNLOAD_DELAY': 2, # Respectful scraping
9         'CONCURRENT_REQUESTS': 1
10    }
11
12    def parse(self, response):
13        for product in response.css('div.product'):
14            yield {
15                'name': product.css('h2::text').get(),
16                'price': product.css('span.price::text').get(),
17                'rating': product.css('div.rating::attr(data-rating)').get()
18            }
```



ETHICAL WEB SCRAPING

1. **Check robots.txt** - Respect website's scraping policy
2. **Rate Limiting** - Don't overwhelm servers
3. **User-Agent** - Identify yourself honestly
4. **Terms of Service** - Read and comply with ToS
5. **Personal Data** - Respect privacy laws (GDPR, CCPA)



Warning

Scraping can violate ToS or copyright. Always consult legal counsel for commercial use.

STREAMING DATA WITH KAFKA

```
1 from kafka import KafkaProducer, KafkaConsumer
2 import json
3
4 # Producer: Collect data
5 producer = KafkaProducer(
6     bootstrap_servers=['localhost:9092'],
7     value_serializer=lambda v: json.dumps(v).encode('utf-8')
8 )
9
10 event = {
11     'user_id': '12345',
12     'event_type': 'page_view',
13     'page': '/products'
14 }
15 producer.send('user-events', event)
16
17 # Consumer: Process data
18 consumer = KafkaConsumer('user-events')
19 for message in consumer:
```



DATA COLLECTION BEST PRACTICES

SCHEMA DESIGN WITH PYDANTIC

```
1 from pydantic import BaseModel, Field
2 from datetime import datetime
3
4 class UserEvent(BaseModel):
5     user_id: str
6     event_type: str
7     timestamp: datetime = Field(default_factory=datetime.utcnow)
8     properties: dict
9     session_id: Optional[str] = None
```



SAMPLING STRATEGIES

```
1 import random
2
3 class SamplingCollector:
4     def __init__(self, sample_rate=0.1):
5         self.sample_rate = sample_rate
6
7     def should_collect(self, event):
8         if random.random() < self.sample_rate:
9             return True
10        if event.get('priority') == 'high':
11            return True
12        return False
```

When to use: High-traffic systems where full collection is expensive



PRIVACY AND COMPLIANCE

```
1 import hashlib
2
3 class PrivacyAwareCollector:
4     def __init__(self, pii_fields=['email', 'phone']):
5         self.pii_fields = pii_fields
6
7     def anonymize(self, data):
8         anonymized = data.copy()
9         for field in self.pii_fields:
10             if field in anonymized:
11                 value = str(anonymized[field])
12                 hashed = hashlib.sha256(value.encode()).hexdigest()
13                 anonymized[field] = hashed
14
15     return anonymized
```



PART 2: DATA VALIDATION

Ensuring Data Quality and Reliability



WHY DATA VALIDATION MATTERS

THE COST OF BAD DATA

- Garbage In, Garbage Out
- Silent Failures
- Expensive Debugging
- Lost Trust

REAL-WORLD IMPACT

E-commerce: Missing IDs → 30% drop in CTR
Medical: Wrong units → Misdiagnoses
Fraud Detection: Duplicates → 45% false positives



TYPES OF DATA QUALITY ISSUES

- 1. Completeness** - Missing or null values
- 2. Accuracy** - Incorrect or outdated values
- 3. Consistency** - Contradictory data
- 4. Timeliness** - Stale or outdated data
- 5. Uniqueness** - Duplicate records
- 6. Validity** - Violates business rules

DATA VALIDATION WITH PYDANTIC

```
1 from pydantic import BaseModel, Field, validator, EmailStr
2
3 class UserEvent(BaseModel):
4     user_id: str = Field(..., min_length=1, max_length=100)
5     email: EmailStr
6     age: int = Field(..., ge=0, le=150)
7     event_type: str
8
9     @validator('event_type')
10    def validate_event_type(cls, v):
11        allowed = ['click', 'view', 'purchase', 'signup']
12        if v not in allowed:
13            raise ValueError(f'event_type must be one of {allowed}')
14        return v
15
16    @validator('timestamp')
17    def timestamp_not_future(cls, v):
18        if v > datetime.utcnow():
19            raise ValueError('timestamp cannot be in the future!')
```



USING PYDANTIC

```
1 try:
2     event = UserEvent(
3         user_id="user_123",
4         email="user@example.com",
5         age=25,
6         event_type="click",
7         timestamp=datetime.utcnow()
8     )
9     print("✓ Valid event")
10 except ValueError as e:
11     print("✗ Validation errors:", e.json())
```



GREAT EXPECTATIONS FRAMEWORK

```
1 import great_expectations as gx
2 from great_expectations.dataset import PandasDataset
3
4 df = pd.read_csv('user_events.csv')
5 dataset = PandasDataset(df)
6
7 # Define expectations
8 dataset.expect_column_values_to_not_be_null('user_id')
9 dataset.expect_column_values_to_be_unique('event_id')
10 dataset.expect_column_values_to_be_in_set(
11     'event_type',
12     ['click', 'view', 'purchase']
13 )
14 dataset.expect_column_values_to_be_between('age', 0, 150)
15
16 # Validate
17 results = dataset.validate()
```



PANDERA: STATISTICAL VALIDATION

```
1 import pandera as pa
2 from pandera import Column, Check
3
4 schema = pa.DataFrameSchema({
5     "user_id": Column(str, checks=[
6         Check.str_length(min_value=1, max_value=100)
7     ]),
8     "age": Column(int, checks=[
9         Check.in_range(min_value=0, max_value=150),
10        Check(lambda s: s.mean() > 18)
11    ]),
12    "purchase_amount": Column(float, checks=[
13        Check.greater_than_or_equal_to(0)
14    ], nullable=True)
15 })
16
17 validated_df = schema.validate(df)
```



CUSTOM VALIDATION PIPELINES

```
1 from dataclasses import dataclass
2 from typing import List
3
4 @dataclass
5 class ValidationResult:
6     passed: bool
7     errors: List[str]
8     warnings: List[str]
9
10 class DataValidationPipeline:
11     def __init__(self):
12         self.validators = []
13
14     def add_validator(self, validator):
15         self.validators.append(validator)
16         return self
17
18     def validate(self, df):
19         all_errors = []
```



STATISTICAL DISTRIBUTION VALIDATION

```
1 class DistributionValidator:
2     def __init__(self, reference_df, columns):
3         self.reference_stats = {
4             col: {
5                 'mean': reference_df[col].mean(),
6                 'std': reference_df[col].std()
7             }
8             for col in columns
9         }
10
11    def validate(self, df):
12        errors = []
13        for col, ref_stats in self.reference_stats.items():
14            current_mean = df[col].mean()
15            mean_shift = abs(current_mean - ref_stats['mean']) / ref_stats['std']
16
17            if mean_shift > 2:
18                errors.append(
19                    f'!{col} - Mean shifted by {mean_shift * 100} std deviation')
20
```



DATA QUALITY MONITORING

```
1 import sentry_sdk
2 from datadog import statsd
3
4 class DataQualityMonitor:
5     def monitor(self, df, dataset_name):
6         result = self.pipeline.validate(df)
7
8         # Send metrics
9         statsd.gauge(f'{dataset_name}.row_count', len(df))
10        statsd.gauge(f'{dataset_name}.validation.errors', len(result.errors))
11
12        # Alert on failures
13        if not result.passed:
14            sentry_sdk.capture_message(
15                f"Data validation failed for {dataset_name}",
16                extras={"errors": result.errors}
17            )
```



PART 3: DATA LABELING

Annotating Datasets with Ground Truth



WHAT IS DATA LABELING?

Data Labeling: Adding meaningful tags or annotations to raw data

WHY IT'S CRITICAL

- **Supervised Learning** requires labeled examples
- **Quality labels** → Better models
- **Consistent labels** → Reliable evaluation
- **Cost:** Often 60-80% of ML project time and budget



TYPES OF LABELING TASKS

CLASSIFICATION

- Image: “cat” vs “dog”
- Text: “spam” vs “not spam”
- Audio: speaker identification

OBJECT DETECTION

- Bounding boxes
- Keypoint annotation
- Polygon segmentation

SEQUENCE LABELING

- Named Entity Recognition
- Part-of-speech tagging
- Time series anomalies

STRUCTURED PREDICTION

- Semantic segmentation
- Dependency parsing
- Relationship extraction



LABEL STUDIO

Label Studio: Open-source, multi-modal annotation platform

KEY FEATURES

- Multi-modal: Images, text, audio, video, time series
- Custom interfaces with XML config
- ML-assisted labeling and active learning
- Collaboration features
- Export: JSON, CSV, COCO, Pascal VOC, YOLO

INSTALLATION



1

LABEL STUDIO: IMAGE CLASSIFICATION

```
1 <View>
2   <Image name="image" value="$image_url"/>
3   <Choices name="choice" toName="image">
4     <Choice value="Cat"/>
5     <Choice value="Dog"/>
6     <Choice value="Other"/>
7   </Choices>
8 </View>
```



LABEL STUDIO: OBJECT DETECTION

```
1 <View>
2   <Image name="image" value="$image_url"
3     zoom="true" zoomControl="true"/>
4   <RectangleLabels name="label" toName="image">
5     <Label value="Person" background="red"/>
6     <Label value="Car" background="blue"/>
7     <Label value="Bicycle" background="green"/>
8   </RectangleLabels>
9 </View>
```

LABEL STUDIO: NER

```
1 <View>
2   <Text name="text" value="$text"/>
3   <Labels name="label" toName="text">
4     <Label value="Person" background="red"/>
5     <Label value="Organization" background="blue"/>
6     <Label value="Location" background="green"/>
7     <Label value="Date" background="orange"/>
8   </Labels>
9 </View>
```

Example: “Apple Inc. CEO Tim Cook announced the new iPhone in Cupertino”



ALTERNATIVE LABELING TOOLS

LABELBOX

- Enterprise features
- Model-assisted labeling
- **Cons:** Expensive, vendor lock-in

PRODIGY

- By spaCy creators
- Active learning built-in
- **Cons:** License cost

CVAT

- Open-source by Intel
- Video annotation
- **Cons:** Setup complexity

SCALE AI / SAGEMAKER

- Managed services
- Human workforce included
- **Cons:** Very expensive



INTER-ANNOTATOR AGREEMENT

Problem: Different annotators may label differently

Solution: Measure agreement to ensure quality

COHEN'S KAPPA (K)

Measures agreement between **two annotators** accounting for chance

$$\kappa = \frac{p_o - p_e}{1 - p_e}$$

Interpretation:

- 0.0-0.20: Slight



COHEN'S KAPPA IMPLEMENTATION

```
1 from sklearn.metrics import cohen_kappa_score
2
3 annotator1 = ['spam', 'ham', 'spam', 'ham', ...]
4 annotator2 = ['spam', 'ham', 'ham', 'ham', ...]
5
6 kappa = cohen_kappa_score(annotator1, annotator2)
7 print(f"Cohen's Kappa: {kappa:.3f}")
```

Example Output:

Cohen's Kappa: 0.745
→ Substantial agreement ✓



FLEISSL' KAPPA: MULTIPLE ANNOTATORS

```
1 from statsmodels.stats.inter_rater import fleiss_kappa
2 import numpy as np
3
4 # Matrix: rows = items, columns = category counts
5 data = np.array([
6     [3, 0],  # Item 1: all 3 annotators said "ham"
7     [2, 1],  # Item 2: 2 said "ham", 1 said "spam"
8     [0, 3],  # Item 3: all 3 said "spam"
9 ])
10
11 kappa = fleiss_kappa(data, method='fleiss')
12 print(f"Fleiss' Kappa: {kappa:.3f}")
```



IMPROVING AGREEMENT

1. CLEAR GUIDELINES

- Define each label precisely
- Provide examples and counter-examples
- Document edge cases
- Create decision trees for ambiguous cases

2. TRAINING PHASE

- Use gold standard datasets
- Review disagreements together
- Iterative calibration sessions



SAMPLING STRATEGIES

RANDOM SAMPLING

```
1 import random  
2 random.sample(dataset, n=100)
```

Pros: Unbiased Cons: May miss rare classes

STRATIFIED SAMPLING

```
1 from sklearn.model_selection import train_test_split  
2 _, sampled, _, _ = train_test_split(  
3     dataset, labels,  
4     train_size=len(dataset)-100,  
5     stratify=labels  
6 )
```

Pros: Ensures class representation



SNOWBALL SAMPLING

```
1 from sklearn.neighbors import NearestNeighbors
2
3 class SnowballSampler:
4     def sample(self, labeled_data, unlabeled_data, n=100):
5         # Vectorize
6         labeled_vectors = self.vectorizer.fit_transform(labeled_data)
7         unlabeled_vectors = self.vectorizer.transform(unlabeled_data)
8
9         # Find nearest neighbors
10        nn = NearestNeighbors(n_neighbors=5)
11        nn.fit(unlabeled_vectors)
12        distances, indices = nn.kneighbors(labeled_vectors)
13
14        return [unlabeled_data[i] for i in indices.flatten()[:n]]
```

Use case: Finding more examples of a rare class



ACTIVE LEARNING

Active Learning: Intelligently select which examples to label next

PROCESS

1. Train model on small labeled set
2. Find **most informative** unlabeled examples
3. Label those examples
4. Retrain model
5. Repeat

Goal: Achieve good performance with minimal labeling



UNCERTAINTY SAMPLING

```
1 class UncertaintySampler:  
2     def sample(self, X_unlabeled, n=10):  
3         probs = self.model.predict_proba(X_unlabeled)  
4  
5         # Higher entropy = more uncertain  
6         uncertainties = -np.sum(probs * np.log(probs + 1e-10), axis=1)  
7  
8         # Select top-n most uncertain  
9         uncertain_indices = np.argsort(uncertainties)[-n:]  
10        return uncertain_indices
```



ACTIVE LEARNING LOOP

```
1 def active_learning_loop(X_labeled, y_labeled, X_unlabeled, budget=100):
2     model = MultinomialNB()
3     sampler = UncertaintySampler(model)
4
5     for iteration in range(budget // 10):
6         # Train model
7         model.fit(X_labeled, y_labeled)
8
9         # Select informative examples
10        indices = sampler.sample(X_unlabeled, n=10)
11
12        # Get human labels
13        new_X = [X_unlabeled[i] for i in indices]
14        new_y = get_human_labels(new_X)
15
16        # Update datasets
17        X_labeled.extend(new_X)
18        y_labeled.extend(new_y)
19
```



ACTIVE LEARNING BENEFITS

COST REDUCTION

Random Sampling: 10,000 labels → 85% accuracy

Active Learning: 2,000 labels → 85% accuracy

Savings: 80% fewer labels!

Key Insight: Active learning reaches target accuracy with 5-10× fewer labels

REAL-WORLD SAVINGS

- $10k \text{ labels} \times \$0.50 = \$5,000$
- $2k \text{ labels} \times \$0.50 = \$1,000$
- **Saved: \$4,000**

WEAK SUPERVISION WITH SNORKEL

Weak Supervision: Use noisy/heuristic labels instead of manual annotation

```
1 from snorkel.labeling import labeling_function
2
3 SPAM = 1
4 HAM = 0
5 ABSTAIN = -1
6
7 @labeling_function()
8 def lf_contains_money(x):
9     money_keywords = ['$','money','cash','prize']
10    return SPAM if any(kw in x.text.lower() for kw in money_keywords) else
11
12 @labeling_function()
13 def lf_short_message(x):
14    return SPAM if len(x.text.split()) < 10 else ABSTAIN
15
16 @labeling_function()
17 def lf_known_sender(x):
18    trusted_domains = ['company.com', 'university.edu']  
CS203: Software Tools and Techniques for AI  
nature HAM is an email sender for a trusted domain). else ABSTAIN
```



COMBINING WEAK LABELS

```
1 from snorkel.labeling.model import LabelModel  
2  
3 # Apply labeling functions  
4 lfs = [lf_contains_money, lf_short_message, lf_known_sender]  
5 L_train = applier.apply(df=df)  
6  
7 # Combine noisy labels  
8 label_model = LabelModel(cardinality=2)  
9 label_model.fit(L_train=L_train, n_epochs=500)  
10  
11 # Get predictions  
12 predicted_labels = label_model.predict(L=L_train)
```

Advantages: Fast, Flexible, No manual labeling needed



PRE-LABELING WITH MODELS

```
1 class PreLabelingPipeline:
2     def pre_label(self, unlabeled_data, confidence_threshold=0.9):
3         predictions = self.model.predict_proba(unlabeled_data)
4
5         auto_labeled = []
6         needs_review = []
7
8         for i, probs in enumerate(predictions):
9             max_prob = max(probs)
10
11         if max_prob >= confidence_threshold:
12             auto_labeled.append({
13                 'data': unlabeled_data[i],
14                 'label': np.argmax(probs),
15                 'source': 'model'
16             })
17         else:
18             needs_review.append(unlabeled_data[i])
19
```



PART 4: DATA AUGMENTATION

Expanding Training Datasets Strategically



WHY DATA AUGMENTATION?

THE PROBLEM

- Limited labeled data is expensive
- Class imbalance leads to biased models
- Overfitting on small datasets
- Rare events underrepresented

THE SOLUTION

Data Augmentation: Creating new training examples by applying transformations

Benefits: Increase dataset size, Improve generalization,

CS 203: Software Tools and Techniques for AI

Reduce overfitting Balance classes



IMAGE DATA AUGMENTATION

```
1 from torchvision import transforms  
2  
3 transform = transforms.Compose([  
4     # Geometric  
5     transforms.RandomHorizontalFlip(p=0.5),  
6     transforms.RandomRotation(degrees=15),  
7     transforms.RandomResizedCrop(size=224, scale=(0.8, 1.0)),  
8  
9     # Color  
10    transforms.ColorJitter(  
11        brightness=0.2, contrast=0.2,  
12        saturation=0.2, hue=0.1  
13    ),  
14  
15    # Blurring  
16    transforms.GaussianBlur(kernel_size=3),  
17  
18    # Normalization  
19    transforms.ToTensor()
```



ADVANCED: ALBUMENTATIONS

```
1 import albumentations as A
2
3 transform = A.Compose([
4     A.RandomResizedCrop(height=224, width=224),
5     A.HorizontalFlip(p=0.5),
6     A.ShiftScaleRotate(shift_limit=0.0625, scale_limit=0.1, rotate_limit=15),
7     A.ElasticTransform(p=0.3),
8
9     # Color
10    A.RandomBrightnessContrast(p=0.5),
11    A.HueSaturationValue(p=0.5),
12
13    # Blur and noise
14    A.OneOf([
15        A.MotionBlur(p=0.5),
16        A.GaussianBlur(p=0.5),
17    ], p=0.3),
18    A.GaussNoise(p=0.3),
19])
```



CUTOUT, MIXUP, CUTMIX

CUTOUT

Random
rectangular masks

```
1 A.CoarseDropout(  
2     max_holes=8,  
3     max_height=32,  
4     p=0.5  
5 )
```

MIXUP

Blend two images

```
1 def mixup(x1, y1, x2  
2     lam = np.random.  
3     x = lam * x1 + (   
4     y = lam * y1 + (   
5         return x, y
```

CUTMIX

Paste regions

```
1 def cutmix(x1, y1, x  
2     # Cut and paste  
3     # region from x2  
4     ...
```



TEXT DATA AUGMENTATION

SYNONYM REPLACEMENT

```
1 from nltk.corpus import wordnet
2
3 def get_synonyms(word):
4     synonyms = set()
5     for syn in wordnet.synsets(word):
6         for lemma in syn.lemmas():
7             synonyms.add(lemma.name())
8     return list(synonyms)
9
10 def synonym_replacement(text, n=2):
11     words = text.split()
12     random_indices = random.sample(range(len(words)), n)
13
14     for idx in random_indices:
15         synonyms = get_synonyms(words[idx])
16         if synonyms:
17             words[idx] = random.choice(synonyms)
18
19     return ' '.join(words)
```



EDA: EASY DATA AUGMENTATION

```
1 class EDA:  
2     @staticmethod  
3     def random_insertion(words, n=1):  
4         """Insert n random synonyms"""  
5         ...  
6  
7     @staticmethod  
8     def random_swap(words, n=1):  
9         """Swap n pairs of words"""  
10        ...  
11  
12    @staticmethod  
13    def random_deletion(words, p=0.1):  
14        """Delete each word with probability p"""  
15        return [w for w in words if random.random() > p]
```



BACK-TRANSLATION

```
1 from transformers import MarianMTModel, MarianTokenizer
2
3 class BackTranslator:
4     def __init__(self, intermediate_lang='fr'):
5         # English → French
6         self.forward_model = MarianMTModel.from_pretrained(
7             f'Helsinki-NLP/opus-mt-en-{intermediate_lang}')
8
9         # French → English
10        self.backward_model = MarianMTModel.from_pretrained(
11            f'Helsinki-NLP/opus-mt-{intermediate_lang}-en')
12
13
14    def augment(self, text):
15        # Translate to French and back
16        intermediate = self.forward_model.translate(text)
17        return self.backward_model.translate(intermediate)
```

Result: “The weather is beautiful” → “The weather is nice”



CONTEXTUAL AUGMENTATION WITH BERT

```
1 from transformers import pipeline
2
3 class ContextualAugmentor:
4     def __init__(self):
5         self.unmasker = pipeline('fill-mask', model='bert-base-uncased')
6
7     def augment(self, text):
8         words = text.split()
9         mask_idx = random.randint(0, len(words) - 1)
10
11     masked_words = words.copy()
12     masked_words[mask_idx] = '[MASK]'
13
14     predictions = self.unmasker(' '.join(masked_words))
15
16     return [pred['sequence'] for pred in predictions[:3]]
```



AUDIO AUGMENTATION

```
1 import librosa
2
3 class AudioAugmentor:
4     @staticmethod
5     def add_noise(audio, noise_factor=0.005):
6         noise = np.random.randn(len(audio))
7         return audio + noise_factor * noise
8
9     @staticmethod
10    def time_stretch(audio, rate=1.0):
11        return librosa.effects.time_stretch(audio, rate=rate)
12
13    @staticmethod
14    def pitch_shift(audio, sr=22050, n_steps=2):
15        return librosa.effects.pitch_shift(audio, sr=sr, n_steps=n_steps)
```



TIME SERIES AUGMENTATION

```
1 class TimeSeriesAugmentor:  
2     @staticmethod  
3     def jittering(x, sigma=0.03):  
4         """Add Gaussian noise"""  
5         return x + np.random.normal(0, sigma, x.shape)  
6  
7     @staticmethod  
8     def scaling(x, sigma=0.1):  
9         """Multiply by random factor"""  
10        factor = np.random.normal(1, sigma, size=(x.shape[0], 1))  
11        return x * factor  
12  
13    @staticmethod  
14    def time_warping(x, sigma=0.2):  
15        """Warp time dimension"""  
16        # Use cubic spline interpolation  
17        ...
```



SMOTE FOR TABULAR DATA

SMOTE: Synthetic Minority Over-sampling Technique

```
1 from imblearn.over_sampling import SMOTE  
2  
3 # Imbalanced dataset: {0: 900, 1: 100}  
4 smote = SMOTE(sampling_strategy='auto', random_state=42)  
5 X_resampled, y_resampled = smote.fit_resample(X_train, y_train)  
6  
7 # Now balanced: {0: 900, 1: 900}
```

How it works:

1. For each minority sample, find k nearest neighbors
2. Randomly select one neighbor
3. Create synthetic sample along line segment



SMOTE VARIANTS

ADASYN

Adaptive Synthetic Sampling

```
1 from imblearn.over_sampling import  
2  
3 adasyn = ADASYN(random_state=42)  
4 X_res, y_res = adasyn.fit_resample
```

Advantage: Generates more samples in harder-to-learn regions

BORDERLINESMOTE

Focus on decision boundary

```
1 from imblearn.over_sampling import  
2  
3 bsmote = BorderlineSMOTE(random_st  
4 X_res, y_res = bsmote.fit_resample
```

Advantage: More conservative, focuses on boundary



GENERATIVE MODELS

```
1 from diffusers import StableDiffusionPipeline
2
3 pipe = StableDiffusionPipeline.from_pretrained(
4     "stabilityai/stable-diffusion-2-1"
5 )
6
7 def generate_augmented_images(prompt, n=10):
8     images = []
9     for i in range(n):
10         image = pipe(prompt, num_inference_steps=50).images[0]
11         images.append(image)
12     return images
13
14 # Generate synthetic training data
15 synthetic = generate_augmented_images(
16     "A high-quality photo of a defective product",
17     n=100
18 )
```

Use cases: Rare defects, medical imaging, artistic styles



BEST PRACTICES

1. DOMAIN-APPROPRIATE

✓ Good: Horizontal flip for general images ✗ Bad: Horizontal flip for text/digits (changes meaning)

2. PRESERVE SEMANTICS

✗ Bad: Extreme rotation for digits ($6 \rightarrow 9$) ✓ Good: Slight rotation ($\pm 15^\circ$)

3. VALIDATION SPLIT FIRST

```
1 # ✓ Correct
2 train, val = train_test_split(data)
3 aug_train = augment(train) # Only augment training
4
5 # ✗ Wrong (data leakage!)
6 aug_data = augment(data)
7 train, val = train_test_split(aug_data)
```



MONITOR AUGMENTATION IMPACT

```
1 def evaluate_augmentation(model, train_data, val_data, aug_levels):
2     results = []
3
4     for strength in aug_levels:
5         aug_train = augment(train_data, strength=strength)
6         model.fit(aug_train)
7         val_acc = model.evaluate(val_data)
8         results.append({'strength': strength, 'accuracy': val_acc})
9
10    # Plot results to find optimal augmentation
11    plt.plot([r['strength'] for r in results],
12              [r['accuracy'] for r in results])
13    plt.xlabel('Augmentation Strength')
14    plt.ylabel('Validation Accuracy')
```



CLASS-SPECIFIC AUGMENTATION

```
1 class ClassBalancedAugmentor:
2     def __init__(self, target_samples_per_class=1000):
3         self.target = target_samples_per_class
4
5     def augment(self, X, y):
6         X_aug, y_aug = [], []
7
8         for cls in np.unique(y):
9             X_cls = X[y == cls]
10            current_count = len(X_cls)
11
12            if current_count < self.target:
13                n_to_generate = self.target - current_count
14                augmented = self.generate_samples(X_cls, n_to_generate)
15                X_aug.extend(augmented)
16                y_aug.extend([cls] * len(augmented))
17
18        return np.array(X_aug), np.array(y_aug)
```



SUMMARY



KEY TAKEAWAYS

DATA COLLECTION

- Instrument systems for automatic capture
- Use appropriate tools (analytics, APM, CDC)
- Implement buffering, batching, error handling
- Respect privacy and comply with regulations

DATA VALIDATION

- Validate early and often
- Use schema validation (Pydantic, Pandera)
- Monitor data quality metrics



KEY TAKEAWAYS (CONT.)

DATA LABELING

- Use appropriate tools (Label Studio, CVAT)
- Measure inter-annotator agreement (Cohen's Kappa)
- Apply smart sampling (active learning)
- Leverage weak supervision when possible

DATA AUGMENTATION

- Choose domain-appropriate transformations
- Preserve label semantics
- Augment only training data



HANDS-ON LAB

COMPLETE DATA PIPELINE

Tasks:

1. Data Collection (30 min) - Scrape reviews, set up logging
2. Data Validation (30 min) - Create schemas, build validation pipeline
3. Data Labeling (45 min) - Label Studio, calculate agreement, active learning
4. Data Augmentation (30 min) - Apply EDA, back-translation, measure impact



ADDITIONAL RESOURCES

DOCUMENTATION

- Label Studio: <https://labelstud.io/guide/>
- Great Expectations: <https://docs.greatexpectations.io/>
- Albumentations: <https://albumentations.ai/docs/>
- Snorkel: <https://www.snorkel.org/>

KEY PAPERS

- “SMOTE: Synthetic Minority Over-sampling Technique”
(2002)
- “Snorkel: Rapid Training Data Creation with Weak
Supervision” (2017)



QUESTIONS?



THANK YOU!

Next: Reproducibility & Versioning



Contact: nipun.batra@iitgn.ac.in

