

МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ  
НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ "ЛЬВІВСЬКА ПОЛІТЕХНІКА"

## **ЛЕКЦІЯ 10. Машинне навчання в реальному світі та у виробництві**

---

**Львів -- 2025**

# Лекція зі штучного інтелекту 2025-10

## Вступ

На цьому занятті ми розглянемо практичні аспекти впровадження систем машинного навчання в реальних умовах. Теоретичні знання про алгоритми машинного навчання є важливими, але справжні виклики часто виникають при інтеграції цих алгоритмів у виробничі системи. Ми зосередимося на всьому життєвому циклі проєктів машинного навчання: від правильної підготовки даних, настройки гіперпараметрів, розгортання моделей до моніторингу та вирішення проблем деградації в процесі експлуатації. Ми також розглянемо організаційні та інфраструктурні питання, які є критичними для успішного впровадження ML-систем.

## Теми, що розглядаються

1. Життєвий цикл проєктів машинного навчання
2. Підготовка даних для виробничих систем
3. Налаштування гіперпараметрів
4. Тестування та валідація моделей
5. Розгортання моделей у виробництво
6. Моніторинг та обслуговування ML-систем
7. Деградація моделей та оновлення
8. MLOps: принципи та практики
9. Масштабування ML-систем
10. Етичні питання та відповідальність

## Життєвий цикл проєктів машинного навчання

### Основні етапи життєвого циклу ML

На відміну від традиційного програмного забезпечення, проєкти машинного навчання мають циклічну природу з наступними основними етапами:

1. **Визначення проблеми та цілей:** Формування чіткого розуміння бізнес-задачі та метрик успіху.
2. **Збір та аналіз даних:** Визначення необхідних даних та їх джерел, аналіз якості та репрезентативності.
3. **Підготовка даних:** Очищення, трансформація, аугментація та інші операції для підготовки даних.
4. **Вибір та навчання моделей:** Експериментування з різними алгоритмами та настройка гіперпараметрів.
5. **Оцінка та валідація:** Тестування моделі на різних наборах даних для оцінки її ефективності.

6. **Розгортання моделі:** Інтеграція моделі в існуючу інфраструктуру та процеси.
7. **Моніторинг та обслуговування:** Відстеження продуктивності моделі та оновлення за необхідності.

Важливо розуміти, що життєвий цикл не є лінійним — часто потрібно повертатися до попередніх етапів на основі результатів та нових знань.

## Порівняння з традиційним розвитком програмного забезпечення

Проекти машинного навчання відрізняються від традиційних проєктів розробки ПЗ кількома ключовими аспектами:

1. **Залежність від даних:** ML-системи критично залежать від якості та кількості даних, тоді як традиційні системи більше залежать від правильності алгоритмів.
2. **Ймовірнісна природа:** Результати ML-моделей є ймовірнісними, а не детермінованими, що ускладнює тестування та налагодження.
3. **Концепція технічного боргу:** В ML-системах технічний борг накопичується не лише в коді, але й у даних, конвеєрах обробки та моделях.
4. **Циклічність розробки:** ML-проєкти вимагають частішого повернення до попередніх етапів та повторного навчання моделей.

Ці відмінності підкреслюють необхідність специфічних підходів до управління проєктами машинного навчання та відповідної організаційної структури.

## Підготовка даних для виробничих систем

---

### Важливість якісних даних

Дані є фундаментом будь-якої системи машинного навчання. Якість ML-моделі безпосередньо залежить від якості даних, на яких вона навчається. У виробничих системах неякісні дані можуть призвести до:

1. **Хибних прогнозів:** Моделі, навчені на неякісних даних, дають неточні результати, що може мати серйозні наслідки для бізнесу.
2. **Упереджених рішень:** Якщо тренувальні дані містять упередження, модель відтворюватиме та потенційно посилюватиме ці упередження.
3. **Необхідності перенавчання:** Виявлення проблем з даними після розгортання призводить до необхідності повторного навчання моделей та оновлення систем.

За оцінками експертів, до 80% часу в проєктах машинного навчання витрачається саме на підготовку та обробку даних.

### Основні етапи підготовки даних

1. **Збір та інтеграція даних**

- **Джерела даних:** Визначення та налаштування доступу до релевантних джерел даних (бази даних, API, сховища файлів, потокові джерела).
- **Стратегії вибірки:** Створення репрезентативних вибірок даних, що відображають реальні умови використання.
- **Конвеєри даних (Data Pipelines):** Автоматизація збору даних з різних джерел та їх інтеграція в єдиний набір.

## 2. Очищення та валідація даних

- **Обробка відсутніх значень:** Виявлення та заповнення пропусків (середніми значеннями, медіанами, прогнозованими значеннями тощо).
- **Виявлення та обробка викидів:** Ідентифікація нетипових значень, які можуть негативно впливати на модель.
- **Виправлення структурних помилок:** Корекція неправильних форматів, дублікатів, невідповідностей у кодуванні тощо.
- **Валідація даних:** Перевірка відповідності даних бізнес-правилам та технічним вимогам.

## 3. Трансформація та збагачення

- **Нормалізація та стандартизація:** Приведення числових ознак до спільного масштабу для покращення роботи алгоритмів.
- **Кодування категоріальних змінних:** Перетворення категоріальних змінних у числовий формат (one-hot encoding, target encoding, тощо).
- **Створення нових ознак (Feature Engineering):** Генерація нових ознак на основі існуючих для покращення прогностичної здатності моделі.
- **Зниження розмірності:** Зменшення кількості ознак для прискорення навчання та зменшення проблеми перенавчання.

## Автоматизація підготовки даних

У виробничих системах процес підготовки даних повинен бути максимально автоматизованим та надійним:

1. **Конвеєри обробки даних:** Створення відтворюваних конвеєрів, які автоматично виконують всі етапи підготовки даних.
2. **Моніторинг якості даних:** Впровадження автоматичних перевірок для виявлення змін у розподілі даних або появи аномалій.
3. **Версіонування даних:** Відстеження версій наборів даних для забезпечення відтворюваності результатів та можливості повернення до попередніх версій.
4. **Метадані та документація:** Документування всіх перетворень та їх обґрунтування для полегшення співпраці та налагодження.

## Дрейф даних та його обробка

Дрейф даних (Data Drift) — це зміна характеристик вхідних даних з часом, що може призвести до зниження ефективності моделі. Типи дрейфу даних:

- Коваріаційний дрейф (Covariate Shift):** Зміна розподілу вхідних ознак без зміни відношення між ознаками та цільовою змінною.
- Дрейф домену (Domain Shift):** Зміна в контексті, в якому використовується модель (наприклад, використання моделі, навченої на даних з одного регіону, в іншому регіоні).
- Концептуальний дрейф (Concept Drift):** Зміна відношення між вхідними ознаками та цільовою змінною.

Стратегії обробки дрейфу даних:

- **Регулярне перенавчання моделей:** Оновлення моделей на нових даних за графіком.
- **Адаптивне навчання:** Автоматичне оновлення моделей при виявленні значущих змін у даних.
- **Моніторинг статистичних характеристик:** Відстеження змін у розподілі вхідних та вихідних даних.

## Типові проблеми та виклики

- Незбалансовані дані:** Нерівномірний розподіл класів у задачах класифікації, що може призводити до упереджених моделей.
- Недостатня кількість даних:** Обмежена кількість прикладів, особливо для рідкісних випадків або нових продуктів.
- Конфіденційність та безпека:** Необхідність захисту особистих даних та дотримання законодавчих вимог (GDPR, CCPA тощо).
- Обсяг та швидкість даних:** Виклики, пов'язані з обробкою великих обсягів даних та даних, що надходять у реальному часі.

Подолання цих викликів вимагає комбінації технічних рішень, організаційних підходів та експертизи в конкретній галузі.

## Настройка гіперпараметрів

### Поняття гіперпараметрів та їх важливість

Гіперпараметри — це конфігураційні змінні алгоритмів машинного навчання, які не можуть бути вивчені безпосередньо з даних і мають бути встановлені перед початком навчання. На відміну від параметрів моделі (наприклад, ваги в нейронній мережі), які оптимізуються під час навчання, гіперпараметри визначають структуру та поведінку алгоритму.

Приклади гіперпараметрів для різних алгоритмів:

- **Дерева рішень**: максимальна глибина дерева, мінімальна кількість об'єктів для розділення вузла, критерій розділення.
- **Нейронні мережі**: кількість шарів, кількість нейронів у кожному шарі, швидкість навчання, функції активації, розмір пакету (batch size).
- **Алгоритми на основі ядер**: тип ядра, параметри ядра, параметр регуляризації.
- **Ансамблеві методи**: кількість базових алгоритмів, швидкість навчання, розмір підвибірки.

Правильний вибір гіперпараметрів може значно покращити продуктивність моделі, тоді як невдалий вибір може призвести до поганої продуктивності навіть для теоретично потужних алгоритмів.

## Методи настройки гіперпараметрів

### 1. Пошук за сіткою (Grid Search)

Найпростіший метод, який перебирає всі можливі комбінації значень гіперпараметрів у заданому дискретному просторі.

#### Переваги:

- Гарантоване дослідження всього заданого простору.
- Простота реалізації та паралелізації.

#### Недоліки:

- Обчислювальна складність зростає експоненційно з кількістю гіперпараметрів ("прокляття розмірності").
- Нераціональне використання ресурсів на дослідження регіонів з низькою продуктивністю.

#### Приклад:

```
from sklearn.model_selection import GridSearchCV

param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [3, 5, 7, 10],
    'learning_rate': [0.01, 0.1, 0.2]
}

grid_search = GridSearchCV(
    estimator=model,
    param_grid=param_grid,
    cv=5,
    scoring='accuracy'
)

grid_search.fit(X_train, y_train)
```

### 2. Випадковий пошук (Random Search)

Замість перебору всіх комбінацій, випадковий пошук обирає випадкові комбінації значень гіперпараметрів із заданих розподілів.

### Переваги:

- Ефективніший за пошук за сіткою при високій розмірності.
- Можливість дослідження неперервних просторів гіперпараметрів.
- Зазвичай знаходить хороші рішення за меншу кількість ітерацій.

### Недоліки:

- Не гарантує знаходження оптимального рішення.
- Може пропустити важливі регіони простору гіперпараметрів.

### Приклад:

```
from sklearn.model_selection import RandomizedSearchCV
from scipy.stats import randint, uniform

param_distributions = {
    'n_estimators': randint(50, 500),
    'max_depth': randint(3, 15),
    'learning_rate': uniform(0.01, 0.3)
}

random_search = RandomizedSearchCV(
    estimator=model,
    param_distributions=param_distributions,
    n_iter=100,
    cv=5,
    scoring='accuracy'
)

random_search.fit(X_train, y_train)
```

## 3. Байєсівська оптимізація

Метод, який використовує байєсівський підхід для моделювання цільової функції (продуктивності моделі) та вибору найбільш перспективних комбінацій гіперпараметрів для наступних експериментів.

### Переваги:

- Ефективно використовує результати попередніх експериментів.
- Швидко знаходить хороші рішення з меншою кількістю ітерацій.
- Добре працює для задач з високою обчислювальною складністю.

### Недоліки:

- Складніша реалізація та налаштування.
- Може застрягати в локальних оптимумах.

### Приклад:

```
from skopt import BayesSearchCV
from skopt.space import Real, Integer
```

```
search_space = {
    'n_estimators': Integer(50, 500),
    'max_depth': Integer(3, 15),
    'learning_rate': Real(0.01, 0.3, prior='log-uniform')
}

bayes_search = BayesSearchCV(
    estimator=model,
    search_spaces=search_space,
    n_iter=50,
    cv=5,
    scoring='accuracy'
)

bayes_search.fit(X_train, y_train)
```

## 4. Еволюційні алгоритми

Використовують принципи природного відбору для пошуку оптимальних гіперпараметрів, включаючи мутації, схрещування та відбір найкращих кандидатів.

### Переваги:

- Ефективне дослідження складних просторів гіперпараметрів.
- Здатність уникати локальних оптимумів.
- Добре працюють для проблем з багатьма гіперпараметрами.

### Недоліки:

- Складні в налаштуванні та можуть повільно збігатися.
- Вимагають значних обчислювальних ресурсів.

## 5. Автоматична настройка гіперпараметрів (AutoML)

Сучасні системи автоматичного машинного навчання (AutoML) надають комплексні рішення для автоматичної настройки гіперпараметрів:

- **AutoSKLearn**: Автоматичний вибір алгоритмів та настройка гіперпараметрів.
- **H2O AutoML**: Автоматичне навчання та настройка різних моделей.
- **Google Vertex AI**: Хмарна платформа для автоматичної оптимізації моделей.

## Стратегії ефективної настройки гіперпараметрів

### 1. Багатоетапний підхід:

- Почати з широкого діапазону значень та грубого кроку.
- Зосередитись на перспективних регіонах з дрібнішим кроком.

### 2. Першочергове налаштування важливих гіперпараметрів:

- Спочатку зосередитись на параметрах, які найбільше впливають на продуктивність.



- Наприклад, для нейронних мереж: архітектура, розмір пакету та швидкість навчання.

### **3. Використання знань про алгоритми:**

- Враховувати відомі залежності між гіперпараметрами.
- Застосовувати експертні знання для визначення розумних діапазонів.

### **4. Раннє припинення неперспективних експериментів:**

- Зупиняти навчання моделей, які показують погані результати на ранніх етапах.
- Перенаправляти ресурси на більш перспективні конфігурації.

## **Практичні рекомендації**

### **1. Вибір метрики оцінки:**

- Використовувати метрики, що відповідають бізнес-цілям.
- Враховувати компроміс між різними аспектами продуктивності (наприклад, точність vs. повнота).

### **2. Валідація:**

- Застосовувати перехресну валідацію для надійної оцінки.
- Враховувати часовий аспект для часових рядів (time series).

### **3. Версіонування:**

- Фіксувати всі гіперпараметри та результати експериментів.
- Використовувати системи відстеження експериментів (MLflow, Weights & Biases).

### **4. Обчислювальна ефективність:**

- Враховувати обчислювальні обмеження при плануванні експериментів.
- Використовувати розподілені обчислення для прискорення процесу.

## **Автоматизація процесу в виробничих системах**

У виробничих системах процес настройки гіперпараметрів повинен бути частиною автоматизованих конвеєрів ML:

### **1. Періодична переоцінка гіперпараметрів:**

- Регулярне проведення нових експериментів для пошуку кращих конфігурацій.
- Автоматичне запускання процесу при появі нових даних або зміні умов.

### **2. Інтеграція з CI/CD:**

- Включення настройки гіперпараметрів у конвеєри безперервної інтеграції та розгортання.
- Автоматичне тестування нових конфігурацій перед впровадженням у виробництво.

### **3. Оптимізація за декількома цілями:**

- Врахування не лише точності, але й швидкості виведення, використання пам'яті та інших ресурсних обмежень.
- Пошук компромісу між продуктивністю та операційними витратами.

## Тестування та валідація моделей

---

### Значення ретельного тестування моделей

Тестування ML-моделей відрізняється від тестування традиційного програмного забезпечення. Замість перевірки правильності виконання конкретних функцій, тестування ML-моделей зосереджене на оцінці їхньої здатності узагальнювати закономірності на нових даних.

Ретельне тестування та валідація є критичними з кількох причин:

1. **Виявлення перенавчання (overfitting):** Визначення, чи модель просто "запам'ятовує" тренувальні дані замість виявлення загальних закономірностей.
2. **Оцінка стабільності:** Перевірка стабільності продуктивності моделі в різних умовах та на різних підмножинах даних.
3. **Бізнес-ефективність:** Оцінка відповідності моделі бізнес-цілям та вимогам.
4. **Безпека та етичність:** Виявлення потенційних ризиків, упереджень та небажаної поведінки моделі.

### Розділення даних для тестування

Основні підходи до розділення даних

#### 1. Простий розподіл (Train-Test Split):

- Базовий підхід із розділенням даних на тренувальний та тестовий набори (зазвичай 70-80% на навчання, 20-30% на тестування).
- Простий у реалізації, але чутливий до випадкового розподілу даних.

#### 2. Тренувальний-валідаційний-тестовий розподіл (Train-Validation-Test Split):

- Розширення простого розподілу з додатковим валідаційним набором.
- Тренувальний набір використовується для навчання, валідаційний — для настройки гіперпараметрів, тестовий — для остаточної оцінки.

#### 3. Перехресна валідація (Cross-Validation):

- k-fold: Дані розділяються на k частин, модель навчається k разів, кожного разу використовуючи іншу частину як валідаційну.
- Забезпечує більш надійну оцінку, але вимагає більше обчислювальних ресурсів.

#### 4. Стратифікована перехресна валідація:

- Розширення звичайної перехресної валідації, яке зберігає пропорції класів у кожному розділі.

- Особливо важливо для незбалансованих наборів даних.

## 5. Часове розділення (Time-Based Split):

- Для часових рядів та послідовних даних.
- Тренування на історичних даних, тестування на більш нових даних.
- Імітує реальні умови використання моделі в часі.

## Визначення оптимального розміру тестової вибірки

Розмір тестової вибірки має забезпечувати статистично значущу оцінку продуктивності моделі:

- **Емпіричне правило:** Як мінімум 30 прикладів для кожного класу в задачах класифікації.
- **Статистичний підхід:** Розмір вибірки, що забезпечує прийнятний довірчий інтервал для метрик.
- **Компроміс:** Більший тестовий набір забезпечує надійнішу оцінку, але зменшує кількість даних для навчання.

## Метрики оцінки моделей

Вибір правильних метрик є критичним для оцінки моделей та повинен відповідати бізнес-цілям проєкту.

### Метрики для задач класифікації

#### 1. Точність (Accuracy):

- Частка правильних прогнозів.
- Підходить для збалансованих наборів даних.
- Може бути оманливою для незбалансованих наборів.

#### 2. Точність та повнота (Precision and Recall):

- Precision: Частка правильних позитивних передбачень серед усіх позитивних передбачень.
- Recall: Частка правильно ідентифікованих позитивних прикладів.
- Дають більш детальне розуміння продуктивності для кожного класу.

#### 3. F1-score:

- Гармонічне середнє між точністю та повнотою.
- Хороший компроміс, коли потрібно балансувати обидві метрики.

#### 4. ROC-крива та AUC:

- Відображає компроміс між чутливістю та специфічністю.
- AUC (площа під кривою) — узагальнена метрика продуктивності.

#### 5. Log Loss / Cross-Entropy:

- Оцінює впевненість моделі в передбаченнях.
- Штрафує за упевнені неправильні прогнози.

### Метрики для задач регресії

### 1. Середня абсолютна помилка (MAE):

- Середнє абсолютне відхилення прогнозів від фактичних значень.
- Менш чутлива до викидів порівняно з MSE.

### 2. Середньоквадратична помилка (MSE):

- Середнє квадратичне відхилення прогнозів від фактичних значень.
- Сильно штрафує за великі помилки.

### 3. Корінь середньоквадратичної помилки (RMSE):

- Корінь з MSE, має ті ж одиниці виміру, що й цільова змінна.
- Часто використовується для інтерпретабельності.

### 4. Коефіцієнт детермінації ( $R^2$ ):

- Показує частку дисперсії цільової змінної, яка пояснюється моделлю.
- Варіюється від 0 до 1, де 1 — ідеальна відповідність.

## Бізнес-орієнтовані метрики

### 1. Метрики вартості помилок:

- Враховують різну вартість хибно-позитивних та хибно-негативних результатів.
- Наприклад, вартість невиявлення шахрайства vs. вартість помилкового позначення законної транзакції як шахрайської.

### 2. Користувацькі метрики:

- Спеціалізовані метрики, що відображають конкретні бізнес-цілі.
- Наприклад, збільшення конверсії, зменшення відтоку клієнтів тощо.

## Підходи до валідації в реальних умовах

### 1. Бекдор-тестування (Backtesting)

- Імітація роботи моделі на історичних даних так, ніби вона використовувалася в той момент часу.
- Особливо важливо для фінансових моделей та моделей часових рядів.
- Дозволяє оцінити продуктивність моделі в різних економічних умовах та ринкових циклах.

### 2. A/B-тестування

- Порівняння продуктивності нової моделі з поточною на частині реального трафіку.
- Дозволяє оцінити вплив моделі на реальних користувачів та бізнес-метрики.
- Вимагає ретельного планування та моніторингу для мінімізації ризиків.

### 3. Тіньове (shadow) розгортання

- Розгортання нової моделі паралельно з існуючою, але без використання її результатів для прийняття рішень.

- Дозволяє моніторити продуктивність нової моделі без ризику для користувачів.
- Може виявити проблеми, які не були очевидні під час офлайн-тестування.

#### **4. Каналова (Canary) релізація**

- Поступове збільшення частки трафіку, яка обробляється новою моделлю.
- Дозволяє обмежити потенційний негативний вплив проблемної моделі.
- Забезпечує плавний перехід від старої моделі до нової.

### **Тестування стійкості моделей**

#### **1. Тестування на різних сегментах даних**

- Перевірка продуктивності моделі на різних підгрупах користувачів/об'єктів.
- Виявлення сегментів, де модель працює гірше, та потенційних упереджень.
- Забезпечення справедливого обслуговування всіх груп користувачів.

#### **2. Стрес-тестування**

- Оцінка продуктивності моделі в екстремальних умовах.
- Тестування на штучно створених "важких" прикладах.
- Виявлення потенційних вразливостей та граничних випадків.

#### **3. Адверсаріальне тестування**

- Перевірка стійкості моделі до навмисно створених прикладів, призначених для обману.
- Особливо важливо для систем безпеки та модерації контенту.
- Допомогає виявити та усунути вразливості моделі.

### **Інтерпретація результатів та значущість**

#### **1. Статистична значущість:**

- Оцінка достовірності спостережуваних відмінностей у продуктивності.
- Використання довірчих інтервалів та статистичних тестів.

#### **2. Аналіз помилок:**

- Детальне вивчення випадків, де модель припускається помилок.
- Виявлення патернів та систематичних проблем.

#### **3. Інтерпретація та пояснення:**

- Використання методів інтерпретації моделей (SHAP, LIME, тощо).
- Пояснення причин конкретних прогнозів для підвищення довіри.
- Виявлення ознак, які найбільше впливають на рішення моделі.

### **Документування процесу та результатів**

Належне документування процесу тестування та валідації є критичним для відтворюваності та аудиту:

### 1. Модельні картки (Model Cards):

- Стандартизована документація, що описує призначення моделі, обмеження, продуктивність та етичні міркування.
- Допомогає користувачам розуміти можливості та обмеження моделі.

### 2. Звіти про валідацію:

- Детальний опис процесу валідації, використаних метрик та результатів.
- Включає аналіз продуктивності на різних сегментах даних та потенційних упереджень.

### 3. Журнали експериментів:

- Відстеження всіх проведених експериментів, їхніх параметрів та результатів.
- Використання спеціалізованих інструментів для відстеження (MLflow, Neptune, Weights & Biases).

## Розгортання моделей у виробництво

---

### Стратегії розгортання ML-моделей

Розгортання моделей машинного навчання — це процес інтеграції навченої моделі в виробниче середовище, де вона може приймати рішення на основі реальних даних. Існує кілька стратегій розгортання, кожна з яких має свої переваги та недоліки:

#### 1. Вбудовані моделі (Embedded Models)

- **Опис:** Модель інтегрується безпосередньо в прикладний код та розгортається разом з ним.
- **Переваги:** Низька латентність, відсутність мережевих залежностей, простота розгортання.
- **Недоліки:** Складність оновлення, обмеження у виборі технологій, потенційно підвищені вимоги до ресурсів клієнтських пристроїв.
- **Приклади застосування:** Мобільні додатки, вбудовані системи, системи безпеки.

#### 2. REST API / Мікросервіси

- **Опис:** Модель розгортається як окремий веб-сервіс, що надає API для взаємодії.
- **Переваги:** Незалежність від клієнтських технологій, простота оновлення, централізоване управління.
- **Недоліки:** Вища латентність через мережеві запити, необхідність управління інфраструктурою.
- **Приклади застосування:** Веб-додатки, бізнес-системи, стандартизовані сервіси ML.

#### 3. Пакетна обробка (Batch Processing)

- **Опис:** Модель періодично обробляє великі обсяги даних та зберігає результати для подальшого використання.

- **Переваги:** Ефективність обробки великих обсягів даних, менші вимоги до інфраструктури реального часу.
- **Недоліки:** Затримка між оновленнями, не підходить для задач, що вимагають миттєвих рішень.
- **Приклади застосування:** Рекомендаційні системи, аналіз клієнтської бази, періодичні звіти.

#### 4. Поточкова обробка (Stream Processing)

- **Опис:** Модель обробляє безперервний потік даних у реальному або близькому до реального часу.
- **Переваги:** Близьке до реального часу прийняття рішень, здатність обробляти великі обсяги динамічних даних.
- **Недоліки:** Складність інфраструктури, вищі вимоги до ресурсів, складність моніторингу.
- **Приклади застосування:** Виявлення шахрайства, моніторинг IoT, торгові системи.

#### 5. Гібридні підходи

- **Опис:** Комбінація різних стратегій розгортання для оптимального балансу між латентністю, масштабованістю та ефективністю.
- **Приклад:** Використання пакетної обробки для створення базових рекомендацій та API для персоналізації в реальному часі.

### Технічні аспекти розгортання моделей

#### 1. Серіалізація та формати моделей

Для розгортання моделей потрібно зберегти у форматі, який можна легко завантажити та використовувати у виробничому середовищі:

- **Стандартні формати:**
  - Pickle (Python): Простий, але має обмеження безпеки та сумісності.
  - ONNX (Open Neural Network Exchange): Міжплатформений формат для обміну моделями між різними фреймворками.
  - PMML (Predictive Model Markup Language): XML-формат для опису предиктивних моделей.
  - TensorFlow SavedModel: Містить графи TensorFlow та їхні ваги.
  - TorchScript: Оптимізована версія PyTorch-моделей.
- **Вимоги до серіалізації:**
  - Версіонування: Можливість чіткого визначення версії моделі.
  - Метадані: Включення інформації про вхідні/вихідні дані, гіперпараметри.
  - Препроцесинг: Включення необхідних трансформацій даних.

#### 2. Контейнеризація та оркестрація

Контейнеризація забезпечує узгоджене середовище виконання та спрощує розгортання:

- **Docker:** Створення ізольованих контейнерів з усіма залежностями.

- **Kubernetes:** Оркестрація контейнеризованих сервісів для автоматичного масштабування, розподілу навантаження та відновлення після збоїв.
- **Оркестрація ML-конвеєрів:** Використання спеціалізованих інструментів (Kubeflow, Airflow) для управління конвеєрами ML від підготовки даних до розгортання.

### 3. Хмарні сервіси для розгортання ML-моделей

Сучасні хмарні платформи пропонують спеціалізовані сервіси для розгортання ML-моделей:

- **AWS SageMaker:** Повний цикл створення, навчання та розгортання моделей.
- **Google Vertex AI:** Платформа для навчання та розгортання моделей з автоматичним масштабуванням.
- **Azure Machine Learning:** Інтегрована платформа для повного життєвого циклу ML.
- **Databricks:** Платформа для аналітики даних та ML на основі Apache Spark.

## Інтеграція з операційними системами

### 1. Конвеєри даних для виробництва

Для ефективної роботи моделей у виробництві необхідно забезпечити надійне постачання та обробку даних:

- **Узгодження форматів:** Забезпечення сумісності форматів даних між різними системами.
- **Обробка в реальному часі:** Системи для обробки поточкових даних (Apache Kafka, Apache Flink).
- **Зберігання та доступ:** Оптимізація доступу до даних для мінімізації латентності (кешування, індексування).
- **Моніторинг даних:** Відстеження якості та узгодженості вхідних даних.

### 2. Інтеграція з бізнес-процесами

Ефективне використання ML-моделей вимагає їхньої інтеграції з бізнес-процесами:

- **API-інтеграція:** Стандартизовані інтерфейси для взаємодії з іншими системами.
- **Вебхуки та події:** Механізми для асинхронної взаємодії між системами.
- **Бізнес-правила:** Комбінування прогнозів моделей з бізнес-логікою.
- **Автоматизація прийняття рішень:** Визначення, які рішення можуть бути повністю автоматизовані, а які вимагають людського втручання.

## Стратегії забезпечення надійності

### 1. Обробка помилок та відмовостійкість

- **Таймаути та повторні спроби:** Механізми для обробки тимчасових збоїв.
- **Деградація функціональності:** Стратегії відступу при недоступності компонентів (fallback strategies).
- **Розподіл навантаження:** Балансування запитів між кількома екземплярами моделі.
- **Асинхронна обробка:** Використання черг повідомлень для буферизації та забезпечення надійної доставки.



## 2. Моніторинг продуктивності

- **Метрики часу відповіді:** Відстеження латентності обробки запитів.
- **Утилізація ресурсів:** Моніторинг використання CPU, GPU, пам'яті та мережі.
- **Логуювання:** Детальне логуювання для аналізу та налагодження.
- **Алерти:** Налаштування повідомлень про проблеми продуктивності.

## 3. Масштабування

- **Горизонтальне масштабування:** Додавання нових екземплярів для обробки збільшення навантаження.
- **Вертикальне масштабування:** Збільшення ресурсів, доступних для окремих екземплярів.
- **Автоматичне масштабування:** Налаштування правил для автоматичного масштабування на основі навантаження.
- **Локальне кешування:** Кешування результатів для зменшення обчислювального навантаження.

## Практичні аспекти розгортання

### 1. Процес розгортання в виробництво

#### 1. Підготовка до розгортання:

- Версіонування моделі та артефактів.
- Підготовка документації та інструкцій для операційних команд.
- Перевірка відповідності безпековим та регуляторним вимогам.

#### 2. Стратегії розгортання:

- Синє-зелене розгортання (Blue-Green Deployment): Паралельне існування двох версій з швидким переключенням.
- Канаркове розгортання (Canary Deployment): Поступове збільшення трафіку на нову версію.
- Розгортання за ознакою (Feature Flag Deployment): Активація нових функцій для певних користувачів або умов.

#### 3. Автоматизація розгортання:

- CI/CD для моделей машинного навчання.
- Автоматизовані тести перед розгортанням.
- Інструменти для автоматичного відкату у випадку проблем.

### 2. Оптимізація моделей для виробництва

- **Квантизація:** Зменшення точності обчислень для покращення продуктивності.
- **Дистиляція моделей:** Створення менших моделей на основі більших для зменшення обчислювальних вимог.
- **Прунінг:** Видалення менш важливих нейронів або з'єднань для зменшення розміру моделі.
- **Компіляція моделей:** Оптимізація моделей для конкретного обладнання (TensorRT, ONNX Runtime).

### 3. Безпека розгортання ML-моделей

- **Захист даних:** Шифрування чутливих даних у русі та у спокої.
- **Автентифікація та авторизація:** Контроль доступу до API моделей.
- **Захист від адверсаріальних атак:** Впровадження методів для виявлення та запобігання навмисним спробам обману моделі.
- **Аудит моделей:** Регулярна перевірка на вразливості та відповідність вимогам.

## Моніторинг та обслуговування ML-систем

### Необхідність безперервного моніторингу

На відміну від традиційного програмного забезпечення, ML-системи взаємодіють з динамічним середовищем, де дані та умови постійно змінюються. Це створює унікальні виклики:

- **Аудит моделей:** Регулярна перевірка на вразливості та відповідність вимогам.

## Деградація моделей та оновлення

### Причини деградації моделей з часом

Моделі машинного навчання, розгорнуті у виробництві, часто з часом втрачають свою ефективність. Це відбувається з різних причин:

#### 1. Дрейф даних (Data Drift)

Зміна розподілу вхідних даних, коли нові дані відрізняються від тих, на яких навчалася модель:

- **Приклади дрейфу даних:**
  - Зміна демографічних характеристик користувачів.
  - Сезонні зміни у поведінці споживачів.
  - Зміни в інтерфейсі, що впливають на поведінку користувачів.
  - Нові продукти або категорії в рекомендаційних системах.
- **Виявлення дрейфу даних:**
  - Статистичні тести (Колмогорова-Смирнова,  $\chi^2$ -квадрат).
  - Моніторинг відстаней між розподілами (KL-дивергенція, JS-дивергенція).
  - Візуальний аналіз розподілів ознак у часі.

#### 2. Дрейф концепцій (Concept Drift)

Зміна взаємозв'язків між ознаками та цільовою змінною:

- **Приклади дрейфу концепцій:**
  - Зміна переваг користувачів (те, що раніше подобалося, тепер не цікаво).
  - Зміна економічних умов, що впливають на фінансові рішення.

- Нові стратегії шахраїв у системах виявлення шахрайства.

- **Виявлення дрейфу концепцій:**

- Моніторинг продуктивності моделі на нових розмічених даних.
- Аналіз помилок прогнозування за підгрупами та часовими періодами.
- Оцінка стабільності важливості ознак.

### **3. Зміни у взаємодії з моделлю**

- **Поведінкові зміни з боку користувачів:**

- Адаптація користувачів до рекомендацій системи.
- Активне або пасивне ігнорування порад алгоритму.

- **Зворотний зв'язок системи:**

- Рекомендаційні системи створюють "фільтрові бульбашки".
- Моделі ціноутворення впливають на поведінку покупців.
- Системи модерації змінюють поведінку спільноти.

### **4. Технічні причини**

- **Зміни в інфраструктурі:**

- Оновлення бібліотек та залежностей.
- Зміни у форматі або схемі вхідних даних.
- Проблеми з конвеєрами препроцесингу.

- **Накопичення технічного боргу:**

- Тимчасові виправлення, що стають постійними.
- Відсутність належного документування та управління кодом.
- Неоптимальні рішення в періоди пікового навантаження.

## **Методи виявлення та кількісної оцінки деградації**

### **1. Метрики продуктивності в часі**

- **Підходи до моніторингу:**

- Відстеження тренду метрик у часі (з контрольними межами).
- Порівняння з базовим періодом або попередньою версією моделі.
- Аналіз сезонних змін та циклічних патернів.

- **Візуалізація деградації:**

- Часові ряди ключових метрик.
- Графіки розподілу помилок у різні періоди.
- Теплові карти продуктивності за сегментами та часом.

## 2. Моніторинг за сегментами

- **Сегментація за:**

- Демографічними характеристиками користувачів.
- Типами об'єктів або транзакцій.
- Географічними регіонами або каналами взаємодії.

- **Виявлення нерівномірної деградації:**

- Ідентифікація сегментів, де продуктивність моделі падає швидше.
- Розрахунок парності та дисперсії продуктивності між сегментами.
- Визначення "сліпих плям" моделі.

## 3. Проактивний моніторинг та симуляції

- **Контрфактуальне тестування:**

- Оцінка продуктивності на нових даних з використанням попередніх версій моделі.
- Порівняння результатів різних версій моделі на однакових даних.

- **Тестування на синтетичних даних:**

- Генерація штучних сценаріїв для оцінки стійкості моделі.
- Симуляція екстремальних випадків або граничних умов.

## Стратегії оновлення моделей

### 1. Регулярні перетренування

- **Фіксований графік оновлень:**

- Щоденне, щотижневе або щомісячне перенавчання.
- Автоматизовані конвеєри для оновлення моделей за розкладом.

- **Переваги та недоліки:**

- Переваги: передбачуваність процесу, простота планування ресурсів.
- Недоліки: можливе зайве або запізнile оновлення, неврахування реальних змін у даних.

### Приклад реалізації:

```
# Спрощений приклад планового перенавчання
from airflow import DAG
from airflow.operators.python_operator import PythonOperator
from datetime import datetime, timedelta

default_args = {
    'owner': 'data_science_team',
    'depends_on_past': False,
    'start_date': datetime(2025, 1, 1),
    'email_on_failure': True,
```

```

    'retries': 1,
    'retry_delay': timedelta(minutes=5),
}

dag = DAG(
    'model_retraining_weekly',
    default_args=default_args,
    description='Weekly model retraining pipeline',
    schedule_interval='0 0 * * 0', # Every Sunday at midnight
)

def retrain_model(**kwargs):
    # Логіка перетренування моделі
    pass

retraining_task = PythonOperator(
    task_id='retrain_model',
    python_callable=retrain_model,
    dag=dag,
)

```

## 2. Оновлення на основі тригерів

- Тригери для оновлення:

- Падіння продуктивності нижче порогового значення.
- Виявлення значного дрейфу даних або концепцій.
- Накопичення достатньої кількості нових розмічених даних.

- Переваги та недоліки:

- Переваги: своєчасне реагування на зміни, ефективне використання ресурсів.
- Недоліки: складність визначення оптимальних порогів, потенційна нестабільність процесу.

### Приклад тригерної системи:

```

# Концептуальний приклад тригерного оновлення
def check_model_performance():
    current_performance = evaluate_model_on_recent_data()
    baseline_performance = get_baseline_performance()

    # Визначення порогу для перенавчання
    threshold = baseline_performance * 0.95 # 5% падіння продуктивності

    if current_performance < threshold:
        trigger_model_retraining()
        send_notification("Model retraining triggered due to performance drop")
    else:
        log_performance_check("Performance check passed")

```

## 3. Інкрементальне навчання

- Підходи до інкрементального навчання:

- Онлайн-навчання, де модель оновлюється при кожному новому спостереженні.
- Міні-пакетне навчання, де оновлення відбувається після накопичення групи нових прикладів.
- Вікно ковзання, де модель тренується на нещодавніх даних із забуванням старих.

- **Підходящі алгоритми:**

- Стохастичний градієнтний спуск (SGD).
- Деякі типи ансамблевих методів (підсилення, баггінг).
- Онлайн-версії лінійних моделей та деяких нейронних мереж.

- **Переваги та недоліки:**

- Переваги: швидке адаптування до нових даних, економія ресурсів.
- Недоліки: складність реалізації, не всі алгоритми підтримують цей підхід, ризик "забування" важливих патернів.

#### 4. Навчання з передачею (Transfer Learning)

- **Застосування:**

- Використання попередньо навченої моделі як основи.
- Додаткове навчання лише на нових даних або модифікація окремих компонентів.
- Особливо ефективно для глибоких нейронних мереж.

- **Переваги та недоліки:**

- Переваги: швидше навчання, менший обсяг необхідних даних.
- Недоліки: можливість збереження застарілих шаблонів, складність налаштування.

### A/B тестування та розгортання оновлень

#### 1. A/B тестування моделей

- **Методологія:**

- Паралельне розгортання поточної (A) та нової (B) версій моделей.
- Розподіл трафіку між версіями для статистично значущого порівняння.
- Оцінка як технічних метрик, так і бізнес-показників.

- **Ключові аспекти:**

- Визначення статистично значущих відмінностей.
- Врахування затримки між прогнозами та спостережуваними результатами.
- Моніторинг непередбачених побічних ефектів.

#### Приклад дизайну A/B тесту для моделі:

Тривалість: 14 днів

Розподіл трафіку: 50% A (поточна модель), 50% B (нова модель)

Первинні метрики: Конверсія, середній дохід на користувача

Вторинні метрики: Латентність відповіді, частота помилок, задоволеність користувачів  
Умови дострокового завершення: >15% падіння конверсії, >500мс збільшення латентності

## 2. Поступове розгортання

- **Стратегії розгортання:**

- Канаркове розгортання: спочатку нова модель обслуговує малий відсоток трафіку.
- Розгортання за сегментами: послідовне впровадження для різних груп користувачів або регіонів.
- Розгортання зі перемикачами функцій (feature flags): можливість швидкого вимкнення нової моделі.

- **Моніторинг під час розгортання:**

- Підвищена частота перевірок метрик.
- Додаткові алерти та поріг для швидкого реагування.
- Детальне логування для діагностики проблем.

## 3. Планування відкату

- **Підготовка до можливих проблем:**

- Документовані процедури відкату.
- Автоматизовані механізми швидкого перемикання на попередню версію.
- Резервне копіювання стану системи перед оновленням.

- **Критерії для відкату:**

- Визначені порогові значення для ключових метрик.
- Часові вікна для прийняття рішень.
- Ланцюжок відповідальності та повноважень.

## Балансування стабільності та інновацій

Управління оновленнями моделей вимагає балансу між стабільністю системи та інноваціями:

### 1. Багатошарова архітектура моделей

- **Різні цикли оновлення:**

- Базові компоненти з повільними оновленнями для стабільності.
- Динамічні компоненти з частішими оновленнями для адаптивності.

- **Композитні моделі:**

- Ансамблі різних моделей з динамічним зважуванням.
- Спеціалізовані моделі для різних сегментів або сценаріїв.

### 2. Експериментальне середовище

- **Постійні експерименти:**
  - Паралельне тестування потенційних покращень.
  - Виділення частини трафіку для експериментів.
  - Культура постійного вдосконалення та навчання.
- **Структурована програма експериментів:**
  - Пріоритезація гіпотез на основі потенційного впливу та зусиль.
  - Документування результатів, навіть негативних.
  - Перенесення успішних експериментів у виробництво.

### 3. Оцінка ризиків та планування оновлень

- **Класифікація оновлень за рівнем ризику:**
  - Низький ризик: незначні налаштування, оновлення даних.
  - Середній ризик: зміна гіперпараметрів, додавання ознак.
  - Високий ризик: нова архітектура, зміна алгоритму, суттєві зміни у даних.
- **Стратегії відповідно до рівня ризику:**
  - Низький ризик: автоматизоване розгортання з мінімальним контролем.
  - Середній ризик: обов'язкове A/B тестування, поступове розгортання.
  - Високий ризик: розширене тестування, обмежене розгортання, розширений моніторинг.

## MLOps: принципи та практики

### Концепція MLOps

MLOps (Machine Learning Operations) — це набір практик, що поєднує розробку моделей машинного навчання (ML) та операційні процеси (Ops) з метою стандартизації та спрощення розгортання ML-систем у виробництво. Це розширення принципів DevOps, адаптоване для унікальних викликів ML-систем.

#### Відмінності між MLOps та DevOps

Аспект	DevOps	MLOps
Фокус	Код та конфігурації	Код, дані та моделі
Тестування	Функціональне, інтеграційне	Додатково: валідація даних, оцінка моделі
Версіонування	Код	Код, дані, моделі, гіперпараметри
Моніторинг	Доступність, продуктивність	Додатково: деградація моделі, дрейф даних
Цикл оновлення	Керується розробниками	Керується даними та продуктивністю моделі

### Ключові компоненти MLOps



## 1. Конвеєри ML (ML Pipelines)

Автоматизовані конвеєри, що забезпечують відтворюваність процесу створення та розгортання моделей:

- **Підготовка даних:** Збір, очищення, валідація та трансформація даних.
- **Навчання моделі:** Вибір алгоритму, оптимізація гіперпараметрів, навчання.
- **Оцінка:** Тестування моделі на різних наборах даних, оцінка метрик.
- **Розгортання:** Підготовка моделі до використання у виробництві.
- **Моніторинг:** Відстеження продуктивності та поведінки моделі.

**Приклад використання Apache Airflow для створення ML-конвеєра:**

```
from airflow import DAG
from airflow.operators.python import PythonOperator
from datetime import datetime

default_args = {
    'owner': 'ml_team',
    'start_date': datetime(2025, 1, 1),
    'retries': 1
}

dag = DAG(
    'ml_pipeline',
    default_args=default_args,
    schedule_interval='@daily'
)

def extract_data(**kwargs):
    # Витягнення даних з різних джерел
    return {'data_path': '/path/to/extracted_data.csv'}

def validate_data(**kwargs):
    # Перевірка якості даних
    ti = kwargs['ti']
    data_info = ti.xcom_pull(task_ids='extract_data')
    # Валідація даних...
    return {'validation_result': 'pass'}

def train_model(**kwargs):
    # Навчання та оцінка моделі
    # ...
    return {'model_path': '/path/to/model.pkl', 'accuracy': 0.92}

def deploy_model(**kwargs):
    # Розгортання моделі у виробництво
    # ...
    pass

extract_task = PythonOperator(
    task_id='extract_data',
    python_callable=extract_data,
    dag=dag
)
```

```

validate_task = PythonOperator(
    task_id='validate_data',
    python_callable=validate_data,
    dag=dag
)

train_task = PythonOperator(
    task_id='train_model',
    python_callable=train_model,
    dag=dag
)

deploy_task = PythonOperator(
    task_id='deploy_model',
    python_callable=deploy_model,
    dag=dag
)

extract_task >> validate_task >> train_task >> deploy_task

```

## 2. Версіонування

Відстеження всіх компонентів ML-системи для забезпечення відтворюваності:

- **Версіонування коду:** Використання Git та інших систем контролю версій.
- **Версіонування даних:** Зберігання та відстеження наборів даних (DVC, Pachyderm).
- **Версіонування моделей:** Реєстрація моделей з метаданими (MLflow, Neptune).
- **Версіонування експериментів:** Фіксація параметрів, результатів та артефактів експериментів.

**Приклад використання DVC для версіонування даних:**

```

# Ініціалізація DVC у проєкті
dvc init

# Додавання даних під контроль DVC
dvc add data/training_data.csv

# Додавання файлів DVC під контроль Git
git add data/training_data.csv.dvc .dvc/config

# Фіксація змін у Git
git commit -m "Add training data"

# Налаштування віддаленого сховища для даних
dvc remote add -d storage s3://mybucket/dvcstore

# Завантаження даних у віддалене сховище
dvc push

```

## 3. Безперервна інтеграція та розгортання (CI/CD)

Автоматизація процесу від розробки до розгортання:

- **CI для ML:** Автоматичне тестування коду, моделей та конвеєрів даних.
- **CD для ML:** Автоматизоване розгортання моделей у різні середовища.
- **Збірка артефактів:** Створення пакетів моделей для розгортання.
- **Стратегії розгортання:** Blue-Green, Canary, Shadow deployments.

#### 4. Моніторинг та зворотний зв'язок

Комплексне відстеження поведінки ML-системи у виробництві:

- **Моніторинг моделі:** Відстеження метрик точності, часу відгуку, розподілу прогнозів.
- **Моніторинг даних:** Виявлення дрейфу даних, аномалій, зміни розподілів.
- **Інфраструктурний моніторинг:** Використання ресурсів, доступність, стабільність.
- **Зворотний зв'язок:** Використання результатів моніторингу для покращення моделей та процесів.

### Рівні зрілості MLOps

#### Рівень 0: Ручний процес

- **Характеристики:** Ручне експериментування, немає стандартизації, ML-інженер виконує всі ролі.
- **Обмеження:** Низька відтворюваність, складність масштабування, висока залежність від індивідуальних знань.

#### Рівень 1: Автоматизація розгортання ML

- **Характеристики:** Автоматизоване розгортання моделей, CI/CD, базове версіонування.
- **Переваги:** Швидше і надійніше розгортання, краща відтворюваність.

#### Рівень 2: Автоматизація CI/CD та тренування

- **Характеристики:** Автоматизовані конвеєри даних та тренування, тригери для перенавчання, комплексний моніторинг.
- **Переваги:** Швидша ітерація моделей, рання ідентифікація проблем, краща адаптація до змін у даних.

### Інструменти MLOps

#### 1. Керування даними

- **DVC (Data Version Control):** Версіонування наборів даних.
- **Pachyderm:** Управління даними та конвеєрами у контейнерах.
- **Delta Lake:** Відкрита платформа для озера даних.

#### 2. Експериментальне відстеження

- **MLflow:** Управління життєвим циклом ML, включаючи експерименти, відтворюваність та розгортання.
- **Weights & Biases:** Платформа для відстеження експериментів, візуалізації та колаборації.
- **Neptune.ai:** Управління метаданими експериментів ML.

### 3. Оркестрація конвеєрів

- **Kubeflow:** Платформа для розгортання робочих потоків ML на Kubernetes.
- **Apache Airflow:** Інструмент для оркестрації і планування робочих потоків.
- **Metaflow:** Фреймворк для науки про дані, що спрощує розробку, оркестрацію та розгортання.

### 4. Моделювання та обслуговування

- **TensorFlow Serving:** Система для обслуговування моделей TensorFlow.
- **TorchServe:** Обслуговування моделей PyTorch.
- **Seldon Core:** Платформа для розгортання ML-моделей на Kubernetes.
- **BentoML:** Фреймворк для пакування та обслуговування моделей.

### 5. Моніторинг

- **Prometheus:** Моніторинг системних метрик.
- **Grafana:** Візуалізація та аналітика метрик.
- **Evidently:** Бібліотека для оцінки та моніторингу ML-моделей.
- **WhyLabs:** Моніторинг ML і якості даних.

## Організація команд та процесів

### 1. Структура команд

- **Спеціалізовані ролі:**
  - **Дослідники даних:** Фокус на дослідженні та розробці моделей.
  - **ML-інженери:** Реалізація та розгортання моделей.
  - **Інженери даних:** Конвеєри та інфраструктура даних.
  - **DevOps-інженери:** Інфраструктура та операційні аспекти.
- **Інтегровані команди:**
  - Міжфункціональні команди, що поєднують різні експертизи.
  - Спільна відповідальність за повний життєвий цикл моделі.

### 2. Робочі процеси та колаборація

- **Стандартизовані процеси:**
  - Узгоджені етапи від ідеї до розгортання.
  - Визначені "воротаря" (gates) для переходу між етапами.
  - Процеси перегляду моделей (Model Review) аналогічно до перегляду коду.
- **Документація та обмін знаннями:**
  - Стандартизовані шаблони для документування моделей.
  - Внутрішні каталоги моделей та компонентів.
  - Регулярні сесії обміну знаннями та ретроспективи.

# Приклади впровадження MLOps

## Перехід від ручного процесу до MLOps

### Вихідний стан (Рівень 0):

- Дослідник даних розробляє модель на своєму комп'ютері
- Ручне збирання та підготовка даних
- Результати моделі експортуються та передаються інженерам для інтеграції
- Оновлення моделі потребує повторення всього процесу

### Цільовий стан (Рівень 2):

- Централізоване сховище для даних та коду
- Автоматизовані конвеєри для збору та підготовки даних
- CI/CD для тестування та розгортання моделей
- Автоматичне перенавчання при виявленні дрейфу даних
- Повний моніторинг продуктивності моделі

### Кроки впровадження:

1. Стандартизація розробки (єдине середовище, стилі коду, шаблони проєктів)
2. Впровадження контролю версій для коду та даних
3. Створення автоматизованих конвеєрів підготовки даних
4. Розробка інфраструктури для тренування та оцінки моделей
5. Впровадження CI/CD для ML-компонентів
6. Налаштування моніторингу та алертів
7. Автоматизація циклу зворотного зв'язку