

# # \*\* MLOps سند ۱۱،۴: برنامه ShahnamehMap\*\* و مانیتورینگ

نسخه: \*\* ۱،۰ \*\*\*

تاریخ: \*\* ۱۵/۰۹/۱۴۰۳ \*\*\*

\*\* ML مهندس / (CTO) تهیه‌کننده: \*\*\* مدیر فنی

وضعیت: \*\* فعال \*\*\*

---

## \*\* ۱. فلسفه MLOps: مفهوم و مقیاس پذیر \*\* يکپارچه‌سازی ایمن

اتوماسیون، کنترل و نظارت بر \*، ShahnamehMap در MLOps هدف  
چرخه حیات مدل‌های یادگیری ماشین\*\*\* است تا اطمینان حاصل شود که  
مدل‌ها \*\*\* به طور مداوم، قابل اعتماد و با کیفیت\*\*\* در سرویس  
Production نسخه‌بندی، ) عمل می‌کنند. ما از \*\*\* اصول مهندسی نرم‌افزار استفاده می‌کنیم ML برای سیستم‌های CI/CD)\*\* تست،

---

## ## \*\*۲. مراحل کلیدی (ML Lifecycle) چرخه حیات مدل

```mermaid

graph LR

A[آزمایش & توسعه] --> B[نیاز کسبوکار / مشکل];

B --> C[انتشار & استقرار]; کد + داده + مدل -->

C --> D[سرویسدهی & مانیتورینگ];

D -- Drift/اعتبارسنجی --> E[بازآموزی & اعتبارسنجی]; افت کیفیت -->

E --> C;

D --> F[بازنشانی & آرشیو];

...

## ## \*\*۳. MLOps Stack (ZirSaxt و ابزارها)

| \*\*\*لایه\*\*\* | \*\*\*ابزار (انتخاب فعلی)\*\*\* | \*\*\*هدف\*\*\* | \*\*\*دلیل انتخاب\*\*\*

|

| :--- | :--- | :--- | :--- |

| \*\*GitLab\*\* | \*\*DVC (Data Version Control)\*\* | داده/مدل) | ردیابی نسخه‌های مختلف داده، کد و ) سبک و مناسب داده‌های کوچک- DVC . موجود CI/CD مدل. | یکپارچه با متوسط

| \*\*Model Registry)\*\* | \*\*DVC Remote Storage\*\* + \*\*S3/MinIO\*\* با متادیتا. | ذخیره‌سازی مدل‌ها ذخیره، وابسته‌سازی و بازیابی مدل‌های آموزش‌دیده. | سادگی. در صورت نیاز ارتقا \*\*MLflow\*\* به قابلیت‌های پیشرفته (مانند مدل‌های مرحله‌بندی) به خواهیم داد

| \*\*CI/CD\*\* + \*\*GitLab CI/CD\*\* خط لوله\*\* Docker\*\* خودکارسازی آموزش، تست و استقرار مدل. | یکپارچگی با موجود DevOps زیرساخت

| \*\*Runtime\*\* | \*\*Docker\*\* محیط اجرا\*\* | \*\*Kubernetes (K8s)\*\* | ایزوله کردن و مقیاس سرویس مدل. | استفاده از زیرساخت موجود

| \*\*FastAPI (Serving)\*\* | \*\*API REST\*\* | استقرار\*\* در کانتینر. | ارائه پیش‌بینی‌های مدل به سرویس‌های دیگر. | سادگی، سازگاری با معماری میکروسرویس

برای آزمایش) + | \*\*JupyterHub\*\* ( استاندارد Python / برای اجرای دسته‌ای) | اسکریپت‌های Batch. | کنترل و انعطاف | Batch. محیط توسعه و اجرای | متریک ( Prometheus | \*\*مانیتورینگ & لاگینگ | لاگ)، \*\*الگوریتم‌های تشخیص (Loki)، (داشبورد Drift). | نظارت بر عملکرد، سلامت و کیفیت مدل. | استفاده از Observability موجود استک |.

—

## \*\*۴. نقش‌ها و مسئولیت‌ها (RACI Matrix)\*\*

| مهندس DevOps | مهندس ML فعالیت | مهندس داده محصول مدیر \*

| ---: | ---: | ---: | ---: | ---: |

| مدل توسعه آزمایش (انجام) R | زیرساخت (مشاوره) C | (مشاوره داده) C | (تصویب نیاز) A |

| \*\*R\*\* | (تعریف مراحل) CI/CD ایجاد خط لوله \* | \*\*R\*\* |  
| | | | (پیاده‌سازی)

| \*\*R\*\* | آموزش و اعتبارسنجی مدل \* | \*\*R\*\* | | | C | \*\*R\*\* | (تهیه)  
| | | | (داده)

| \*\*R\*\* | (تأیید کیفیت) A | \*\*R\*\* | \*انتشار و استقرار مدل\* |  
| | | | (استقرار)

| (تحلیل، تنظیم آستانه) \*\*R\*\* | \*\*R\*\* | مانیتورینگ و هشدار مدل \* |  
| | | | (تنظیم ابزار) C |

| \*\*R\*\* | (پیشنهاد) /Rollback\*\* | \*\*R\*\* | (تصمیم بازآموزی) C | A  
| | | | (تصویب کسب‌وکار)

---

## \*\*۵. ML Pipeline)\*\* (برای مدل CI/CD خط لوله .

یا `main` یا اسکریپت‌های آموزش در شاخه (`/ml`) هر بار که کد مدل `develop` آپدیت می‌شود، خط لوله زیر به طور خودکار اجرا می‌شود:

```
```yaml
# .gitlab-ci.yml (ML بخش)

stages:
  - test-ml
  - train-model
  - evaluate-model
  - build-model-image
  - deploy-to-staging

train-and-evaluate:
  stage: train-model
  script:
    - python train.py --config configs/model_v1.yaml # آموزش با DVC
      نمایش معیارهای ارزیابی
    - dvc metrics show # اجرای تست‌های واحد مدل
      artifacts:
        paths:
```

- خروجی آموزش # model.pkl

- metrics.json

only:

- main

- develop

- /<sup>^</sup>release-.\*/

deploy-model-staging:

stage: deploy-to-staging

script:

- docker build -t shahnameh-rec-model:\$CI\_COMMIT\_SHA .

- kubectl set image deployment/ml-recommender model=shahnameh-rec-model:\$CI\_COMMIT\_SHA -n staging

only:

- develop به استیجینگ # اتوماتیک

...

## \*\* دروازه‌های کیفیت (Quality Gates):\*\*

۱. \*\*Gate 1 (Unit Tests):\*\* مدل (Tested) کد را عبور از تست‌های واحد (Unit Tests).
۲. \*\*Gate 2 (Data Quality):\*\* بروزرسانی کیفیت داده ورودی آموزش (Data Quality).  
عدم (Tested) داده بررسی کیفیت داده ورودی آموزش (Data Quality).
۳. \*\*Gate 3 (Feature Completeness):\*\* مدل جدید باید از لحاظ (Feature Completeness) معیارهای ارزیابی (Feature Completeness) مطابق باشد.  
حداقل به اندازه  $Precision@5 \geq 95\%$  و  $Recall@10 \geq 95\%$ .  
عملکرد داشته باشد. در غیر این صورت، خط لوله (Production) مدل قبلی متوقف و هشدار داده می‌شود.
۴. \*\*Gate 4 (A/B Test):\*\* Production قبل از Staging در Staging جدید در محیط با  $10\%$  از ترافیک واقعی به مدت حداقل ۴۸ ساعت تست می‌شود. تنها در صورت بهبود یا عدم تغییر معنی‌دار A/B CTR مانند (CTR) معیارهای کسب و کار راه می‌یابد.

---

## ## \*\*Production\*\* مانیتورینگ جامع در ۶\*

:مانیتور می‌کنیم Production ما چهار نوع مشکل اصلی را در مدل

## مانیتورینگ عملکرد سخت افزاری / زیرساختی . ۶،۱. # ## \* \* # مانیتورینگ (Infrastructure Monitoring) \*\*

- \* \*\* CPU/Memory مصرفی کانتینر مدل، زمان API (Latency)، نرخ خطای HTTP (5xx). پاسخگویی ابزار های Prometheus از طریق Exporter K8s).
- \* \*\* هشدار: اگر Latency P95 > 300ms یا نرخ خطا < ۱٪ برای ۵ دقیقه.

## مانیتورینگ دقت مدل (Model Performance Monitoring) ۶،۲. # ## \* \* #

- \* \*\* چالش: در سیستم های توصیه گر، بر چسب حقیقت زمینی بلا فاصله در دسترس نیست (کاربر ممکن است هفته بعد با یک توصیه تعامل کند)

### \* \*\* (Proxy Metrics): راهکار

- \* \*\* Real-time. نرخ کلیک (CTR) بر روی توصیه ها: ردیابی
- \* \*\* Conversion Rate): نرخ تبدیل (نرخ شروع بازی پس از کلیک روی توصیه.

\* میانگین زمان صرف شده در بخش (Engagement):\*\* مشارکت  
توصیه‌ها.

\* روزانه نسبت به میانگین متحرک ۷ روزه \*\*بیش CTR هشدار:\*\*\* اگر از ۲۵٪ افت\*\* کند.

### \*\*\* # مانیتورینگ انحراف داده (Data Drift Monitoring)\*\*

می‌تواند عملکرد (Feature Drift)\*\* انحراف در توزیع \*\*ورودی‌های مدل را خراب کند.

\* \*\*: ویژگی‌های تحت نظرارت `user\_avg\_rating`,  
`campaign\_difficulty\_tag`, `interaction\_type\_weight`.

\* \*\*: مقایسه توزیع \*\*ویژگی‌های هر درخواست ورودی Production\*\* با توزیع \*\*داده‌های آموزشی\*\* (مرجع) با استفاده از KL Divergence یا PSI (Population Stability Index)\*\* به صورت روزانه.

\* \*\*: برای هر ویژگی کلیدی \*\*<۰,۲\*\* شود PSI هشدار:\*\*\* اگر (نشان‌دهنده تغییر معنی‌دار).

## #### \*\*\* ۶.۴. مانیتورینگ انحراف مفهومی (Concept Drift Monitoring)\*\*\*

انحراف در رابطه بین \*\*ورودی و خروجی\*\*. برای ما، یعنی آیا کاربران هنوز به ویژگی‌های قدیمی (مثلًا "نژاد پهلوان") همانند گذشته واکنش نشان می‌دهند؟

\* \*\*روش:\*\* ذخیره \*یک نمونه تصادفی ۱٪ از درخواست‌های پیش‌بینی\* به همراه \*نتیجه تعامل کاربر (مثلًا کلیک/عدم کلیک)\* در روی Logistic Regression)\*\* آینده. هر هفته، یک \*مدل ساده مدلی که روی داده AUC آن را با AUC این داده جدید آموزش می‌دهیم و قدیمی آموزش دیده مقایسه می‌کنیم.

\* \*\*باشد  $AUC > 0.5$ \*\* هشدار: \*اگر \*افت\*

---

## \*\*۷. مدیریت نسخه (Retraining Policy) سیاست بازآموزی\*\*

#### \*\*۷.۱. بازآموزی برنامه‌ریزی شده (Scheduled Retraining):\*\*

\* \*\*\* هفتگی برای مدل توصیه‌گر (با داده فعلی). این فرکانس با رشد داده می‌تواند تغییر کند.

\* \*\*\* (Trigger): زمان‌بند (Cron Job) در GitLab CI. به طور خودکار ('train-model' مرحله CI/CD روند: خط لوله) می‌شود. اگر مدل جدید از معیارهای ارزیابی (Gate 3) عبور کرد، ) A/B مستقر شده و وارد مرحله Staging به طور خودکار در Test می‌شود.

### \*\*\* ۷.۲. بازآموزی مبتنی بر رویداد (Event-Triggered Retraining):\*\*

\* \*\*\* ماشه‌ها:

یا افت  $0.3 > \text{PSI}$ : شدید Data/Concept Drift هشدار .  
 $\text{AUC} > 0.1$ .

۲. افت شدید Proxy Metrics: (CTR  $< 40\%$ ).

۳. تغییر اساسی در محصول: راهاندازی یک ویژگی بزرگ که الگوهای رفتار کاربر را تغییر می‌دهد.

\* \*\*\* را مطلع می‌کند. آنها به صورت دستی خط ML روند: هشدار، تیم لوله بازآموزی را راهاندازی و فرآیند را تسريع می‌کنند.

## #### \*\*۷.۳. مدیریت نسخه‌ها (Versioning):\*\*

- \* \*\*`MODEL\_NAME-vMAJOR.MINOR.PATCH` قالب:\*\*
- \* \*\*MAJOR:\*\* تغییر معماری یا شکست در سازگاری (باعث می‌شود سرویس‌های مصرف‌کننده نیاز به بروزرسانی داشته باشند)
- \* \*\*MINOR:\*\* اضافه شدن ویژگی‌ها یا بهبود عملکرد (سازگار عقب‌رو)
- \* \*\*PATCH:\*\* رفع باگ یا بازآموزی با داده جدید (سازگار عقب‌رو) ذخیره‌سازی:
  - \* هر نسخه از مدل، همراه با کد، داده و در Registry مدل (DVC از طریق) هایپرپارامترهای دقیقاً تکرار پذیر ذخیره می‌شود.

---

## ## \*\*۸. Rollback و بازیابی استراتژی (Rollback & Recovery)\*\*

اصل: \*\*توانایی بازگشت سریع به آخرین نسخه سالم مدل در کمتر از ۱۰ دقیقه\*\*.

\* \*\*Deployment روش:\*\*\* مدل‌های \*  
آن ساده است rollback مدیریت می‌شوند. مکانیزم \*K8s\*\*:

```bash

kubectl rollout undo deployment/ml-recommender

...

این دستور، Deployment را به نسخه قبلی (stable)\*\* برمی‌گرداند.

\* \*\*Rollback م Ashe‌های \*: دستی

هشدارهای بحرانی مانیتورینگ که نشان‌دهنده خرابی کامل مدل است ۱.

استیجینگ A/B Test شناسایی باگ‌های جدی در ۲.

درخواست مدیر محصول به دلیل رفتار غیرمنتظره مدل ۳.

\* \*\*Fallback State:\*\* Rollback در طول ( ) یا خرابی سرویس مدل، ( ) سیستم به الگوریتم هیورستیک ساده مثلًا (Heuristic Fallback) \*\*مثالاً ( ) توصیه محبوب‌ترین‌ها + جدیدترین‌ها) تغییر حالت می‌دهد تا تجربه کاربر کاملاً قطع نشود.

---

## \*\* جمع‌بندی: قابل اعتماد و قابل نگهداری؟ ۹. \*\*##

ما در مسیر قابل اعتماد و قابل ML با استناد به این برنامه، سیستم نگهداری بودن قرار دارد، زیرا

| ریسک/اقدام ShahnamehMap\*\* | \*\*وضعیت\*\* | آینده

| :--- | :--- | :--- |

| عالی: کد، داده و Reproducibility)\*\* | ✓\*\* تکرارپذیری | - | .ورژن می‌شوند DVC مدل با

| برای CI/CD خوب\*\* (Automation)\*\* | ✓\*\* اتوماسیون | Drift آموزش و استقرار پایه وجود دارد. | نیاز به اتوماسیون کامل تشخیص بازآموزی |

| در حال توسعه: مانیتورینگ Monitoring)\*\* | ⚠\*\* مانیتورینگ | Drift برقرار است. مانیتورینگ Proxy Metrics مانیتورینگ زیرساخت و Drift در حال پیاده‌سازی. | تکمیل مانیتورینگ | و یکپارچه‌سازی هشدارها

| \*\*\* خوب: \*\* تست واحد و \*\* (Testing) |  \*\* تست و اعتبارسنجی |  
وجود دارد. | افزودن تست‌های یکپارچگی CI دروازه کیفیت مبتنی بر معیار در  
(Integration) | خودکار A/B Test و |

| \*\*\* خوب: \* مدیریت خطأ و \* Rollback |  \*\* استراتژی Rollback  
تعریف شده است. | مستندسازی دقیق‌تر سناریوهای Fallback ساده و  
Rollback. |

اگر تغییرات Drift.\*\* بزرگ‌ترین ریسک فعلی: تشخیص به موقع\*\*  
را دیر تشخیص دهیم، مدل به (Concept Drift) تدریجی در رفتار کاربران  
قوی (بخش Drift تدریج بی‌اثر می‌شود. \* پیاده‌سازی سیستم مانیتورینگ  
۳.۶ و ۴.۶) اولویت فوری ماست

نتیجه: \*\*\* ما از \*\*\* بمب ساعتی\*\*\* (مدل بدون نظارت) فاصله گرفته‌ایم و \*\*  
یک \*\*\* چرخه حیات کنترل شده و مبتنی بر داده\*\*\* برای مدل‌های خود ایجاد  
کرده‌ایم. این چارچوب به ما امکان می‌دهد با اطمینان، مدل‌های بیشتری را در  
آینده توسعه و مستقر کنیم.