

# # \*\* سند ۴,۹: مستندات مجموعه داده (Data Sheet + Dictionary) — ShahnamehMap\*\*

نسخه: \*\* ۱,۰ \*\*\*

تاریخ: \*\* ۵/۰۹/۱۴۰۳ \*\*\*

\* مدیر فنی (CCO) تهیه‌کننده: \* مدیر ارشد محتوا \*

(محصول ML برای تیم) وضعیت: \* فعال \*

---

## ## \*\* ۱. مقدمه: نقش داده و ShahnamehMap\*\*

به صورت (ML) از \*\*یادگیری ماشین ShahnamehMap در حال حاضر، محدود و متمرکز استفاده می‌کند. تمرکز اصلی بر \*\*ساخت پایگاه داده و استفاده از \*\*قوانین (Knowledge Graph) دانش ساختاریافته است. با این حال، ما از روز (Rule-based Systems) مبتنی بر دانش مانند ) \*ML\* اول زیرساخت جمع‌آوری داده را برای \* کاربردهای آینده توصیه‌گر محتوا، تشخیص محتوای نامناسب، بالانس خودکار بازی) طراحی کردہ‌ایم.

این سند سه مجموعه داده اصلی را مستند می‌کند: ۱) داده هسته ساختاریافته شاهنامه، ۲) داده رفتاری کاربران، و ۳) داده محتوای کاربرساخته (UGC)\*\*.

---

## ## \*\*۲. Data Sheet: (Shahnameh Core Knowledge Graph)\*\*

### ### \*\*۲,۱. (Motivation)\*\*

\* معتبر، ایجاد یک منبع حقیقت (Single Source of Truth)\*\* ساختاریافته و قابل ماشین خوانی برای تمام موجودیت‌ها و روابط در شاهنامه فردوسی.

### ### \*\*۲,۲. (Composition)\*\*

\* ذخیره می‌شوند Neo4j در Graph) داده‌ها به صورت گراف.

\* \* \* (Nodes/Entities): نودها، موجودات شخصیت‌ها، مکان‌ها، افسانه‌ای، اشیاء، رویدادها.

\* \* \* (Edges/Relationships): روابط بین نودها (مثلاً: یال‌ها). مตولد\_شده\_در، فرزند\_است\_از، کشته\_شد\_توسط، سفر\_کرد\_به.

\* \* \* (Properties): ویژگی‌های هر نود (مانند نام، نام پدر، لقب، جنس، وابستگی قومی) ویژگی‌ها.

### ### \* \* ۲,۳. (Collection Process) جمع‌آوری فرآیند

\* منبع اولیه: متن تصحیح شده شاهنامه (چاپ مسکو و خالقی مطلق) به عنوان منبع مرجع.

\* منابع ثانویه: کتب و مقالات دانشگاهی معتبر برای تکمیل و تصحیح روابط جغرافیایی و شجره‌نامه‌ای (مانند "جغرافیای شاهنامه" از احمدی).

\* استخراج دستی و نیمه‌خودکار توسط تیم متخصص محتوا (با تحصیلات ادبیات فارسی) با استفاده از ابزارهای annotating و سپس تبدیل به گراف.

### ### \* \* ۲,۴. (Preprocessing/Cleaning) پیش‌پردازش/تمیزکاری

- \* های یک نام (مثلاً "rstem"، variant یکسان‌سازی نام‌ها: همه می‌شوند canonical "rstem زال") به یک نود.
- \* برای نام‌های مشترک (مثلاً چند Disambiguation): حل ابهام و منابع برای ایجاد نودهای جداگانه Context ("گردآفرید")، از زمینه استفاده شده است.
- \* اعتبارسنجی: هر رکورد توسط حداقل دو عضو تیم و نهایتاً توسط یک استاد مشاور خارجی بررسی و تأیید شده است.

### ### \*\*۲,۵. (Distribution) توزیع \*\*\*

- \* (Read-only) دسترسی عمومی: یک نسخه محدود و فقط-خواندنی API از گراف (حدود ۲۰۰ موجودیت کلیدی) از طریق در آینده قابل دسترس خواهد بود.
- \* دسترسی داخلی: گراف کامل به عنوان سرویس داخلی برای محصول ما استفاده می‌شود.

### ### \*\*۲,۶. (Maintenance) نگهداشت \*\*\*

- \* بروزرسانی: دستی و بر اساس پروژه. فعلاً برنامه‌ای برای بروزرسانی خودکار نیست.

\* \*\* نسخه‌بندی: هر تغییر اساسی، یک نسخه جدید از مجموعه داده ایجاد می‌کند (مثلاً `shahnameh-core-v1.2`).

---

### ## \*\*۳. Data Sheet: (User Behavioral Event Data)\*\*

#### ### \*\*۱. انجیزه\*\*

\* درک تعامل کاربران با محصول برای بهینه‌سازی تجربه کاربری (UX) \*\*و تحلیل رشد (Personalization)\*\*، (Growth Analytics)\*\*.

#### ### \*\*۲. ترکیب داده\*\*

\* ذخیره داده‌ها به صورت رویداد (Event)\*\* در ClickHouse\*\* می‌شوند. هر رکورد یک رویداد است.

\* \*\*`user\_events`\*\*: جدول

### \*\*\*فرآیند جمع‌آوری \*\*\* #### ۳,۳.

منبع:\*\*\* تعاملات کاربران با فرانت‌اند و بک‌اند پلتفرم \* \* .  
با \*\*\* کدهای (Event Tracking) روش:\*\*\* ردیابی رویداد \* \*  
در فرانت‌اند (با رعایت حریم خصوصی) و \* \* Instrumentation \* \*  
لاغ‌های سرور\* \* در بک‌اند \* \* .

### \*\*\*پیش‌پردازش/تمیزکاری \*\*\* #### ۳,۴.

\* \* \* واقعی به یک `user\_id` (Anonymization):\*\*\* ناشناس‌سازی  
پس از ۲۴ IP یک‌طرفه تبدیل می‌شود. \* \* آدرس `pseudo\_user\_id` \* \* . ساعت حذف می‌شود

\* \* \* Schema:\*\*\* اعتبارسنجی \* \* Schema  
اعتبارسنجی می‌شوند تا از \* \* Schema (Apache Avro یا Protobuf) با) ثبت شده  
انسجام داده اطمینان حاصل شود

### \*\*\*توزیع . ۳,۵ \*\*\* #### ۳,۵.

\* \* فقط دسترسی داخلی. \* \* توسط تیم‌های محصول، تحلیل داده و \* \*  
آینده) با کنترل‌های دسترسی \* \* نیاز-به-دانستن (Need-to-Know) \* \* مهندسی  
استفاده می‌شود

### ### \*\*نگهداشت .۳,۶\*\*

\* حذف خودکار: داده‌های رویداد پس از ۱۴ ماه به طور \*  
\* مطابق با خطمشی حریم خصوصی و) خودکار حذف می‌شوند (GDPR).

\* اسکریپت‌های مانیتورینگ، تغییرات ناگهانی در Drift: \*\* نظارت بر \*  
\* مثلاً اگر ناگهان رویداد) حجم یا توزیع رویدادها را گزارش می‌کند  
\* (کاهش یابد ۵۰٪ `character\_created`).

---

### ## \*\*۴. Data Sheet:\*\* (UGC)\*\*

#### ### \*\*۴,۱. انگیزه\*\*

\* برای کمپین‌های Recommender) توسعه سیستم توصیه‌گر \*  
\* و Content Moderation) کاربری، \* تشخیص محتوای نامناسب  
\* \*\*\* تحلیل خلاقیت جامعه.

#### ### \*\*۴,۲. ترکیب داده\*\*

فایل‌های ) \*\*\*MinIO\*\*\* ساختار کمپیون) و ) \*\*\*Neo4j\*\*\* داده‌ها در آپلودشده) ذخیره می‌شوند.

\* جدول/ساختار:\*\*\* شامل متادیتای کمپیون (عنوان، توضیح، سازنده، بروچسب‌ها) و گراف صحنه‌ها

\*\*\*فرآیند جمع‌آوری ### \* ۴,۳۰.

\* استفاده Campaign Builder منبع:\*\*\* کاربران خالق که از ابزار می‌کنند.

\* \*\*\*TOS) رضایت و مالکیت:\*\*\* کاربران با پذیرش \*\*\*قوانين خدمات و \*\*\*خط مشی محتوا\*\*\*، صراحتاً

تأیید می‌کنند که \*\*\*مالک محتوای تولیدشده\*\*\* هستند . ۱.

۲. مجوز غیرانحصاری، جهانی و رایگان \*\*\* به ShahnamehMap \*\*\* برای \*\*\*ذخیره‌سازی، نمایش، توزیع و ایجاد آثار مشتق (برای بهبود الگوریتم‌ها)\*\*\* می‌دهند.

موافقت می‌کنند که محتوای آن‌ها ممکن است برای \*\*\*تمرین ۳. داخلی\*\*\* مورد استفاده قرار گیرد ML مدل‌های

\*\*\*پیش‌پردازش/تمیزکاری ### \* ۴,۴۰.

\* XSS سانیتايز کردن متن: \* ورودی های متنی برای جلوگیری از \*\* سانیتايز می شوند.

\* بررسی اولیه محتوای نامناسب: \* یک \* فیلتر کلمه کلیدی ساده \* برای مسدود کردن آشکارترین محتواهای (Keyword Filter) \*\* توهین آمیز در لحظه آپلود اجرا می شود.

#### \*\*\*\* توزیع . ٤,٥

\* عمومی: \* کمپین های منتشر شده توسط کاربران برای همه کاربران \*\* قابل مشاهده است.

\* شامل کمپین های ( UGC مجموعه داده کامل \* ML برای ) داخلی \*\* منتشرن شده) ممکن است برای \* تمرین مدل های داخلی تشخیص محتوا یا سیستم توصیه گر \* استفاده شود.

#### \*\*\*\* نگهداری . ٤,٦

\* حذف به درخواست کاربر: \* کاربران می توانند محتوای خود را حذف \*\* حذف \* خواهد \* ML کنند. در این صورت، داده از مجموعه های آموزشی شد.

\* حذف به دلیل تخلف: \* محتوای حذف شده توسط مادها، در یک برای 'violation' مجموعه داده جداگانه " منتشر نشده " با برچسب آموزش مدل تشخیص تخلف \* آینده نگهداری می شود.

—

## \*\*<sup>۵</sup>\*\* ویژگی‌های کلیدی - (Data Dictionary) فرهنگ داده.

\*: از مجموعه داده هسته شاهنامه ۵,۱\*\*\*


| `description` | Text | پهلوان بزرگ | توضیح متنی از منبع. | سیستان، فرزند زال و روتابه... | ~۱۰٪ (برای موجودات کم‌همیت) | منبع اصلی (شماره بیت). | "شاهنامه، جلد ۲، بیت ۱۲۳۴" | ۵٪ (برای اطلاعات استنتاج شده) | `canonical\_source` | String | |

| `gender` | Enum ('MALE', 'FEMALE', 'UNKNOWN') | جنسیت | (برای موجودات افسانه‌ای) | ~۳۰٪ | `MALE` |

\*\*\* #: از مجموعه داده رفتاری کاربران ۰,۲۵\*\*\*

| \*Properties\* | \*\*\* \*\*| نام رویداد\*\* | \*\*\* \*\*| ویژگی‌های کلیدی | توضیح | --- | --- | --- |

| `game\_session\_started` | `session\_id`, `campaign\_id`, `pseudo\_user\_id`, `character\_id` | شروع یک جلسه بازی |

| `player\_action` | `session\_id`, `action\_type` ('ATTACK', 'DIALOG\_CHOICE', 'MOVE'), `target\_id`, `outcome` | هر عمل بازیکن در بازی |

| `campaign\_published` | `campaign\_id`,  
`pseudo\_user\_id`, `category`, `word\_count` | انتشار یک  
| کمپین کاربرساخته  
| `payment\_succeeded` | `pseudo\_user\_id`, `amount`,  
`plan\_type` | پرداخت موفق کاربر. \*\*(حساس) |

### UGC:\*\*\* ۳،۵\*\*\* از مجموعه داده \*\*\*

| نام ویژگی\*\*\* | \*\*\*نوع\*\*\* | \*\*\*توضیح\*\*\* | \*\*\*مسائل کیفیت\*\*\* |  
| Quality Concerns |  
| --- | --- | --- | --- |  
| متن اصلی داستان/دیالوگ | Text |  
| بازتاب دیدگاه فرهنگی/سیاسی کاربر | \*\*Bias:\*\* | نوشته شده توسط کاربر  
| داستان های ناتمام |  
| امبدینگ از | `creator\_style\_embedding` | Vector (Float[]) |  
| سبک نوشتاری کاربر (آینده). | نیاز به مدل برای تولید  
| میانگین امتیاز کاربران دیگر | `user\_rating\_avg` | Float |  
| ممکن است تحت تأثیر | \*\*Bias:\*\* | \*\*Missingness:\*\* | بالا\*\* در ابتدا  
| دوستان سازنده باشد |

---

## ## \*\*(Drift) و انحراف (Bias) ملاحظات اخلاقی، سوگیری .۶\*\*

### #### \*\*۱.۶ سوگیری‌های شناخته‌شده (Known Biases):\*\*

- \* \*\*سوگیری جنسیتی در داده هسته:\*\* شاهنامه جهان‌سازی مردانه دارد. تعداد شخصیت‌های مرد بسیار بیشتر از زن است. هر مدلی که روی این داده آموزش ببیند، این سوگیری را به ارت می‌برد.
- \* \*\*اقدام کاهش:\*\* در کاربردهای آموزشی، این سوگیری به وضوح محتوای تولیدشده توسط کاربران اولیه UGC: سوگیری در داده ممکن است تم‌ها و سبک خاصی داشته باشد که نماینده RPG گیمرهای کل جامعه ایران نباشد.
- \* \*\*اقدام کاهش:\*\* تشویق کاربران متنوع (معلمان، دانشآموزان) به خلق محتوا.

## ### \*\*۶،۲. (Drift) انحراف احتمالی\*\*

\* در داده رفتاری: اگر ویژگی‌های جدیدی به بازی اضافه شود، الگوهای رفتار کاربران تغییر می‌کند.

\* مانیتورینگ: نظارت بر توزیع رویدادهای کلیدی برای شناسایی تغییرات ناگهانی (``action_type``)

## \*\*۶،۳. رضایت و اخلاقی ملاحظات\*\*

\* همانطور که در بخش ۴،۳ آمده، مجوز استفاده از ML: رضایت برای UGC گنجانده شده است. ما از "انتخاب عدم TOS" در ML برای برای آموزش مدل با محتوای کاربر استفاده "Opt-out" مشارکت نمی‌کنیم، زیرا این امر اثربخشی مدل‌های جمعی را مختل می‌کند. در عوض، کاربر می‌تواند محتوای خود را به طور کامل حذف کند.

\* حريم خصوصی: همه داده‌های رفتاری ناشناس شده هستند. هیچ کاربران ندارد (Re-identify) مدلی سعی در شناسایی مجدد.

---

## ## \*\*۷. (Ownership, Licensing, Accountability)\*\*

مجموعه داده	مالک اصلی	مجوز استفاده داخلی				
\*ShahnamehMap\*\*	\*مجوز انتشار عمومی	---	---	---	---	---
اما (Public Domain)\*\* هسته شاهنامه	متون کهن					
ساختار و گردآوری ما	متعلق به شرکت است.	مجوز کامل برای				
با ) CC BY-NC-SA 4.0\*\* استفاده در محصول.	بخشی تحت مجوز					
ذکر منبع) منتشر می‌شود						
داده رفتاری	کاربران (به عنوان موضوع داده)، اما حقوق					
جمع‌آوری و تحلیل متعلق به شرکت است (طبق خطمشی حریم						
خصوصی).	فقط برای تحلیل و بهبود محصول.	هرگز منتشر نمی‌شود.				
ممکن است به صورت تجمعی و ناشناس در گزارش‌ها ذکر شود						
کاربر خالق (مالکیت معنوی). کاربر به شرکت UGC\*\* داده						
مجوز استفاده داده است.	مجوز برای نمایش، ذخیره‌سازی و استفاده					
داخلی.	فقط همان‌طور که کاربر انتخاب کرده است (عموماً ML در					
منتشر شده).						

مسئول نهایی کیفیت، \*\*(CCO) مسئول: \*\*\* مدیر ارشد محتوا \*  
یکپارچگی و استفاده اخلاقی از همه مجموعه داده‌ها است. \* مدیر فنی  
\*\* (CTO) مسئول امنیت و زیرساخت داده است.

---

نتیجه: داده به عنوان دارایی کنترل شده \*\* ##

یک دارایی ShahnamehMap این مستندات نشان می‌دهد که داده در  
استراتژیک اما کنترل شده است. ما

\* منبع و مالکیت هر داده را می‌دانیم \*  
آن را مستند (Bias, Missingness) کیفیت و محدودیت \*  
کرده‌ایم.

\* حقوق کاربران را در قلب فرآیند داده قرار داده‌ایم \*  
ایجاد کرده‌ایم \* ML چارچوبی برای استفاده مسئولانه و آینده‌نگر از \*

این رویکرد از تبدیل داده به یک بدھی حقوقی، اخلاقی یا فنی در آینده  
جلوگیری می‌کند و آن را به موتور قابل اعتماد نوآوری محصول تبدیل می‌کند.

داده‌های شفاف و خوش‌تعریف، سنتگ بنای هر سیستم هوش مصنوعی  
ارزشمندی هستند.