

Versionskontrolle für Machine-Learning-Projekte

Originalartikel

Backup

<html> <p>„Betrieb“<a href=„<https://www.informatik-aktuell.de/betrieb.html>“ target=„_self“>Betrieb“<a href=„<https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz.html>“ target=„_self“>Künstliche Intelligenz</p>Veit Schiele 11. August 2020<header></header><div class=„ce-textpic ce-left ce-intext“><div class=„ce-gallery ce-row ce-column“ data-ce-columns=„1“ data-ce-images=„1“><figure class=„image“><figcaption class=„image-caption“>„Mit DVC können Sie sprachunabhängig ML-Pipelines mit den zugehörigen Trainingsdaten, Konfigurationen, Leistungsmetriken usw. definieren.“</figcaption></figure></div><div class=„ce-bodytext“><p>In diesem Artikel erfahren Sie, wie die Modellentwicklung für maschinelles Lernen (ML) systematisch organisiert werden kann. So kann die Leistung eines Modells verbessert werden, wenn die Parameter feiner abgestimmt oder wenn mehr Trainingsdaten verfügbar werden. Um die Verbesserung messen zu können, sollte nachverfolgt werden können, welche Daten für das Training in welcher Modelldefinition und -konfiguration (Parameter etc.) verwendet und welche Modellleistungen damit erzielt wurden. Dabei sollten sowohl die Daten wie auch der zugehörige Programmcode in einer Version erfasst werden.</p><p>DVC (Data Version Control) wurde entwickelt, um Sie genau bei dieser Aufgabe zu unterstützen <a href=„<https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/versionskontrolle-fuer-mchine-learning-projekte.html#c32318>“>[1]. Durch die Implementierung einer DVC-Pipeline werden alle Daten geladen, vorverarbeitet, trainiert und die Leistung bewertet, wobei der Vorgang vollständig reproduzierbar und automatisierbar ist. Trainingsdaten, Modellkonfiguration, das Modell und Leistungsmetriken sind so versioniert, dass Sie bequem zu einer bestimmten Version zurückkehren und alle zugehörigen Konfigurationen und Daten überprüfen können. Außerdem bietet DVC einen Überblick über Metriken für alle Versionen Ihrer Pipeline, mit deren Hilfe Sie die beste Version ermitteln können. Zudem können Sie die Trainingsdaten, Modelle, Leistungsmetriken usw. mit anderen teilen und eine effiziente Zusammenarbeit ermöglichen.</p><h2>Warum DVC?</h2>Git-annex speichert wie DVC große Dateien nicht im Git-Repository selbst, sondern in einem lokalen Schüssel-Wert-Speicher und verwendet Hardlinks oder Symlinks, anstatt Dateien zu duplizieren <a href=„<https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/versionskontrolle-fuer-mchine-learning-projekte.html#c32318>“>[2].Git-LFS verwendet Reflinks oder Hardlinks, um Kopiervorgänge zu vermeiden und so große Dateien effizienter verarbeiten zu können. DVC ist jedoch kompatibel zu deutlich mehr Remote-Speichern (S3, Google Cloud, Azure, SSH usw.) <a href=„<https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/versionskontrolle-fuer-mchine-learning-projekte.html#c32318>“>[3].Andere Workflow-Management-Systeme wie MLflow sind meist sehr allgemein und nicht speziell für die Verwaltung von

Daten in ML-Projekten entwickelt worden <a href=„<https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/versionskontrolle-fuer-machine-learning-projekte.html#c32318>“>[4]. <p>DAGsHub ist ein DVC-Äquivalent, jedoch nur für GitHub <a href=„<https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/versionskontrolle-fuer-machine-learning-projekte.html#c32318>“>[5]. </p><h2>Ein Beispielprojekt</h2><p>Dieser Artikel führt Sie durch ein Beispielprojekt mit folgenden Phasen:</p>Repositories erstellenDatenpipelines definierenReproduzierenPipeline visualisierenDaten teilen</div></div><header><h3>1. Repositories erstellen</h3></header><p>1. Zunächst wird ein Repository mit einer Versionsverwaltung Ihrer Wahl vorbereitet. In unserem Beispiel ist dies Git, DVC kann jedoch auch mit jeder anderen Versionsverwaltung zusammenarbeiten.</p><pre>\$ git init</pre><p>2. In diesem Repository DVC initialisieren:</p><pre>\$ dvc initYou can now commit the changes to git....</pre><p>3. Initiales Git-Repository einchecken:</p><pre>\$ git status neue Datei: .dvc/.gitignore neue Datei: .dvc/config\$ git add .dvc \$ git commit -m „Initial repo“</pre><p>4. Daten mit DVC verwalten:</p><pre>\$ mkdir data\$ dvc get <https://github.com/iterative/dataset-registry> get-started/data.xml \ -o data/data.xml\$ dvc add data/data.xml</pre><p>5. Datenänderungen mit Git synchronisieren:</p><pre>\$ git add data/data.xml.dvc data/.gitignore\$ git commit -m „Add raw data“</pre><p>6. Entfernen Datenspeicher konfigurieren:</p><p>Sie können DVC-Daten und -Modelle mit dvc push auch außerhalb Ihres lokalen Caches an einem entfernten Ort speichern, damit sie später auch in anderen Umgebungen abgerufen werden können. Außerdem sind dies entfernte Cloud-Services (S3, Azure Blob Storage, Google Cloud Storage), aber auch SSH, HDFS, und HTTP sind möglich. Um das Beispiel möglichst einfach nachvollziehbar zu halten, wählen wir hier einen lokalen Speicherort:</p><pre>\$ sudo mkdir -p /var/dvc-storage\$ dvc remote add -d local /var/dvc-storageSetting 'local' as a default remote.</pre><p>7. Speichern und teilen:</p><p>Mit dvc push kopieren Sie Dateien aus Ihrem lokalen DVC-Cache in den zuvor konfigurierten Remote-Storage.</p><pre>\$ dvc push</pre><p>8. Überprüfen:</p><p>Sie können dies überprüfen, z. B. mit:</p><pre>\$ ls -R /var/dvc-storagea3 f1/var/dvc-storage/a3:04afb96060aad90176268345e10355/var/dvc-storage/f1:5a7474cd26c014ce0cf7a8a3d50516.dir</pre><p>Beachten Sie, dass beide Versionen der Daten gespeichert sind und bestimmen sollten mit dvc cache. </p><div class=„ce-image ce-left ce-above ce-gallery ce-row ce-column“><figure class=„image“><figcaption class=„image-caption“>Abb. 1: Repositories erstellen. Cusy GmbH, 2020</figcaption></figure></div><header><h3>2. Datenpipelines definieren</h3></header><p>Die Versionierung großer Daten für Data Science ist ein Schritt in die richtige Richtung, aber noch nicht ausreichend, wenn Daten gefiltert, transformiert oder zum Trainieren von ML-Modellen verwendet werden sollen. Daher erfasst DVC auch die Abfolge der Prozesse bei der Manipulation der Daten. Damit können die Ergebnisse später genau so reproduziert werden, wie sie entstanden sind. Unser Beispiel soll die Pipeline aus den folgenden fünf Phasen bestehen:</p>VorbereitenAufteilen von Trainings-

und

TestdatenMerkmalsextraktionTrainierenEvaluieren<p>Das Ergebnis dieser Pipeline sind die Leistungsmetriken des trainierten Modells. Das gesamte Schema sieht so aus:</p><p>Rohdaten → aufteilen → Trainings-/Testdaten → extrahieren → Features → trainieren → Modell → evaluieren → Metriken</p><p>2.1. Vorbereiten: Als Vorbereitung auf die Datenpipeline unseres Beispiels benötigen wir zunächst eine virtuelle Python-Umgebung:</p><pre>\$ python3 -m venv venv\$ source venv/bin/activate\$ pip install -r requirements.txt</pre><p>In unserem Beispiel enthält die requirements.txt-Datei die vier Python-Pakete pandas, sklearn, scikit-learn und scipy. Beachten Sie jedoch bitte, dass DVC sprachunabhängig und nicht an Python gebunden ist. Niemand kann Sie davon abhalten, die Phasen in Bash, C oder in einer anderen Lieblingssprache oder einem anderen Framework wie Spark, PyTorch usw. zu implementieren.</p><p>2.2. Aufteilen von Trainings- und Testdaten: Mit dvc run können Sie einzelne Verarbeitungsphasen erstellen, wobei jede Phase durch eine mit Git verwaltete Quellcode-Datei sowie die Abhängigkeiten und Ausgabedaten beschrieben wird. Alle Phasen zusammen bilden dann die DVC-Pipeline. Unsere erste Phase soll die Daten in Trainings- und Testdaten aufteilen:</p><pre>\$ dvc run -n split -d src/split.py -d data/data.xml -o data/splitted \python src/split.py data/data.xml</pre>-n split gibt den Namen mit der Beschreibung der Verarbeitungsphase an.-d src/split.py und -d data/data.xml geben die Abhängigkeiten (dependencies) an. Wenn sich später eine dieser Daten ändert, erkennt DVC, dass die Ergebnisse neu berechnet werden müssen.-o data/splitted gibt das Verzeichnis an, in das die Ergebnisse geschrieben werden sollen. In unserem Fall sollte sich der Arbeitsbereich geändert haben in: ├── data │ ├── data.xml │ ├── data.xml.dvc │ └── splitted+│ ├── test.tsv+│ └── train.tsv+│+├── dvc.lock+├── dvc.yaml ├── requirements.txt └── src └── split.py</pre>python src/split.py data/data.xml ist der Befehl, der in dieser Verarbeitungsphase ausgeführt wird. Die resultierende dvc.yaml-Datei sieht dann so aus:stages: split: cmd: python src/split.py data/data.xml deps: - data/data.xml - src/split.py outs: - data/splitted</pre>In dvc.lock werden hingegen die MD5-Hashwerte gespeichert, anhand derer DVC erkennen kann, ob Änderungen an den Dateien vorgenommen wurden:</p><pre> split: cmd: python src/split.py data/data.xml deps: - path: data/data.xml md5: a304afb96060aad90176268345e10355 - path: src/split.py md5: ffa32f4104c363040f27d2bd22db127d outs: - path: data/splitted md5: 1ce9051bf386e57c03fe779d476d93e7.dir</pre><p>Da die Daten im Ausgabeverzeichnis nie mit Git versioniert werden sollten, hat dvc run dies auch bereits in die data/.gitignore-Datei geschrieben:</p><pre> /data.xml+ /splitted</pre><p>2.3. Merkmalsextraktion:</p><p>Die nächsten Sie nun erstellen, indem die Ausgabe der vorhergehenden als Abhängigkeit definiert wird, in unserem Beispiel mit:</p><pre>dvc run -n featurize -d src/featurization.py -d data/splitted \ -o data/features \ python src/featurization.py data/splitted data/features</pre><p>Sie können diese Verarbeitungsphase jedoch auch parametrisieren. Hierfür erstellen wir in unserem Beispiel die Datei params.yaml mit folgendem Inhalt:</p><pre> max_features: 6000 ngram_range: lo: 1 hi: 2</pre><p>Der Aufruf führt dann den obigen Befehl noch -p <filename>: <params_list> hinzu:</p><pre>\$ dvc run -n featurize -d src/featurization.py -d data/splitted \ -p params.yaml: max_features, ngram_range.lo, ngram_range.hi -o data/features \ python</pre>

+ *| train
| +---
+
*
 **
 +---+ | evaluate |
 +---+</p><header><h3>5. Daten teilen</h3></header><p>Sie können
nun Ihren Code und Ihre Trainingsdaten einfach mit anderen Teammitgliedern teilen. Sofern die
Teammitglieder ebenfalls auf unser lokales DVC-Repository in /var/dvc-storage zugreifen
können, können Sie die Ergebnisse unseres Beispiels reproduzieren mit:</p><pre> \$ git
clone <https://github.com/veit/dvc-example.git> \$ cd dvc-example \$ dvc pull -TR A data/data.xml 1 file
added \$ ls data/ data.xml data.xml.dvc</pre><header><h2>Fazit</h2></header><p>Mit DVC
können Sie sprachunabhängig reproduzierbare ML-Pipelines definieren und zusammen
mit den zugehörigen Trainingsdaten, Konfigurationen, Leistungsmetriken usw. versioniert
speichern. Dabei arbeitet DVC mit allen modernen Versionsverwaltungen zusammen und
unterstützt viele verschiedene Speicherarten wie S3, Google Cloud, Azure, SSH usw. Damit
strukturiert DVC nicht nur die Datenhaltung, sondern durch einzelne, atomare Phasen der DVC-
Pipeline bleiben Änderungen in den Daten auch transparent und nachvollziehbar. Insgesamt
erleichtert und effektiviert dies die Arbeit an ML-Projekten erheblich.</p><div id=„c32318“
class=„frame quellen frame-type-text frame-
layout-0“><header><p>Quellen</p></header><a href=„<https://dvc.org/>“ title=„DVC“
target=„_blank“ rel=„noreferrer“>DVC<a href=„<https://git-annex.branchable.com/>“
title=„git annex“ target=„_blank“ rel=„noreferrer“>git-annex<a
href=„<https://git-lfs.github.com/>“ title=„LFS“ target=„_blank“ rel=„noreferrer“>Git Large File Storage
(LFS) <a href=„<https://mlflow.org/>“ title=„MLflow“ target=„_blank“
rel=„noreferrer“>MLflow<a href=„<https://dagshub.com/>“ title=„DAGsHub“
target=„_blank“ rel=„noreferrer“>DAGsHub<p>Github: <a
href=„<https://github.com/iterative/dvc>“ title=„DVC“ target=„_blank“
rel=„noreferrer“>DVC Beispielprojekt: Github: <a href=„<https://github.com/veit/dvc-example>“
title=„dvc-example“ target=„_blank“ rel=„noreferrer“>dvc-example</p><p><a
href=„[https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/versionskontrolle-fuer-machin
e-learning-projekte.html#top“><img
src=„<https://www.informatik-aktuell.de/fileadmin/templates/wr/images/LinkToTop.png>“ width=„50“
height=„50“ alt=““ referrerpolicy=„no-referrer“
/></p></div><header><p>Autor</p></header><div class=„tx-pwteaser-pi1 main odd“><a
href=„<https://www.informatik-aktuell.de/autoren-cv/veit-schiele.html>“><img
src=„\[https://www.informatik-aktuell.de/fileadmin/_processed/_4/0/csm_200-veit-schiele_f5dbff384.jp
g“ width=„150“ height=„150“ alt=““ referrerpolicy=„no-referrer“ /><h3><a
href=„<https://www.informatik-aktuell.de/autoren-cv/veit-schiele.html>“>Veit Schiele</h3>Veit
Schiele ist erfahrener Trainer für Analysten, Wissenschaftler und Ingenieure bei der Entwicklung
von Analyse- und Forschungssoftware sowie Autor des Jupyter-Tutorial und des PyViz-Tutorial.<a
href=„\\[https://www.informatik-aktuell.de/autoren-cv/veit-
schiele.html\\]\\(https://www.informatik-aktuell.de/autoren-cv/veit-schiele.html\\)“>>> Weiterlesen</div><header><p>Das könnte Sie auch
interessieren</p></header><p>Kommentare \\(0\\)</p> </html>\]\(https://www.informatik-aktuell.de/fileadmin/_processed/_4/0/csm_200-veit-schiele_f5dbff384.jp\)](https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/versionskontrolle-fuer-machin)

From:
<https://schnipsl.qgelm.de/> - **Qgelm**

Permanent link:
<https://schnipsl.qgelm.de/doku.php?id=wallabag:wb2versionskontrolle-fur-machine-learning-projekte>

Last update: **2025/06/27 11:17**

