



Machine Learning Engineer - P5

Note technique :

Catégorisation automatique des questions du site openclassroom

Créé le : 10/07/2024

Par Pégeot Guillaume

Table des matières

Introduction	4
1 Préparation des données	5
1.1 Données d'entraînement	5
1.1.1 Récupération des données pour l'entraînement	5
1.1.2 Point d'attention	6
1.1.3 Traitement des données en vue de l'entraînement	6
2 Modèle BERT	10
2.1 Limitations matériel	10
2.2 Paramétrage	10
2.3 Performances	13
3 Mise en place de l'api	15
3.1 API	15
3.2 Tests unitaire	17
3.3 Déploiement	18
Conclusion	20

Table des figures

1.1	Répartition des questions dans le temps pour le jeu d'entraînement	6
1.2	Tags les plus représenté lors de l'entraînement	7
1.3	stopword & lemmation	8
2.1	Définition du modèle et du tokenizer	11
2.2	Définition des paramètres du trainer	12
2.3	Définition du trainer avec les paramètres précédement définis .	12
2.4	Performance observé avec ml flow	13
2.5	Accuracy observé avec ml flow	13
2.6	Recall observé avec ml flow	14
2.7	F1 score observé avec ml flow	14
3.1	Définition du modèle et du tokenizer	15
3.2	Prédire les tags avec l'API	16
3.3	call web pour avoir les tag via navigation	16
3.4	appel API sans interface	17
3.5	Tests unitaire via github action	18
3.6	Statuts des data upload via github	18
3.7	Commande terminal pour build et commit l'image docker . . .	19
3.8	Récupérer le password pour upload docker image	19

Introduction

Stack overflow est le site du réseau Stack Exchange Network pour toute question relative aux problèmes de programmations. Le site sert de plateforme pour que les utilisateur puissent poser des questions et aussi répondre à celle des autres utilisateurs de manière collaborative.

Le site est une véritable référence pour toute personne travaillant dans le monde l'informatique avec ses 23 millions d'utilisateurs enregistrés, 24 millions de question et 35 millions de réponses (chiffres de mars 2024)

Il faut savoir que les questions sont classées en fonction de leurs tags qui sont définis par les utilisateurs lors de la création de leur question. C'est ici que nous intervenons afin de proposer des tags en fonction de l'énoncé de la question

1 Préparation des données

1.1 Données d'entraînement

1.1.1 Récupération des données pour l'entraînement

Pour entraîner notre modèle nous avons fait appel au site data.stackexchange.com afin d'effectuer des requêtes sql, qui nous ont permis de récupérer :

- Id.
- Titre.
- Contenu de la question (BODY).
- Score.
- Tags.
- Date de création.
- Nombre de vues.
- Nombre de réponses.
- Nombre de commentaires.
- Nombre de favoris.

Il est à noter que les requêtes retournent au maximum 50 000 entrées.

Nous avons pris la décision d'entraîner notre modèle sur les individus ayant le meilleur score et au moins 4 tags sur la période allant du 01/01/2021 au 01/01/2024.

1.1.2 Point d'attention

Parmis les données récupérés, seul le Titre et les Tags nous intéressent pour notre modèle actuel. Il est cependant intéressant de regarder la distribution des questions dans le temps.

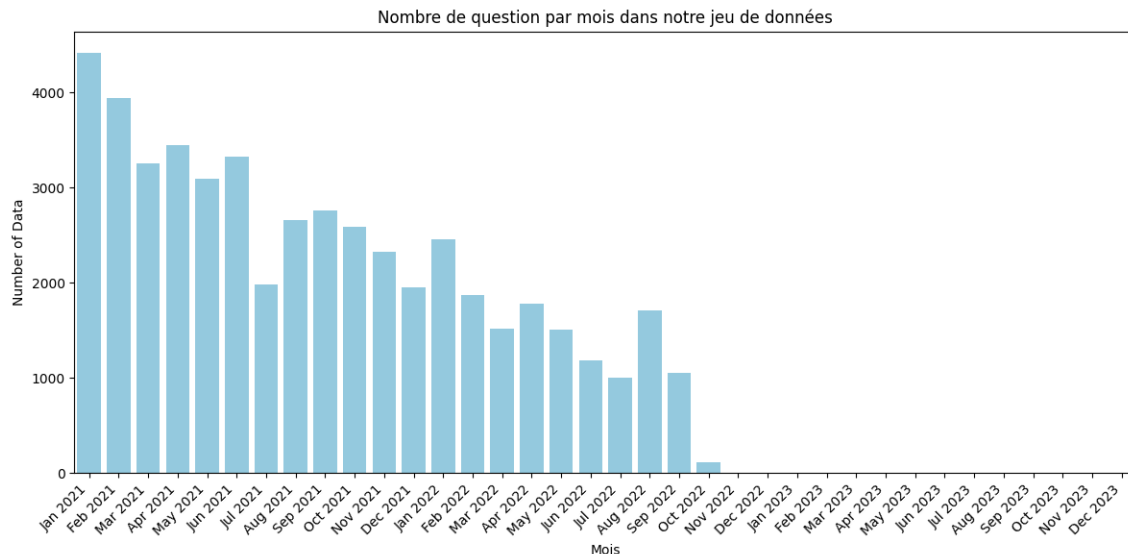


FIGURE 1.1 – Répartition des questions dans le temps pour le jeu d'entrainement

On remarque une grande diminution du nombre de question dans le temps, avec un très faible nombre de question sur l'année 2023. Cela s'explique par le choix de trier les question en fonction de leur score, en effet le score va dépendre du nombre de votes positifs que celle-ci a reçu, favorisant alors les questions les plus anciennes. Ce biais à des chances de pénaliser notre modèle, car il pénalise les documents récent et donc risque de passer à coter des problèmes et des tags émergents.

1.1.3 Traitement des données en vue de l'entrainement

Pour l'entrainement seul les titres et les tags seront conservés.

1.1.3.1 Traitement des tags

Les valeurs obtenu dans le champs tags sont de la forme suivante :

1.1.3.2 Traitement des titres

Les titres se présentent sous forme de string, nous allons donc leur appliquer :

1. Split la liste sur le caractère espace.
2. Passer en lower l'ensemble des individus issue du split
3. Retirer les individus dans les stopwords de la librairie nltk
4. Appliquer un stemmer sur les individus restant
5. Joindre les individus

```
1 import nltk
2 from nltk.corpus import stopwords
3 from nltk.stem import PorterStemmer
4
5 # get the stopwords corpus
6 nltk.download('stopwords')
7 stop_words = set(stopwords.words('english'))
8
9 # get the stemmer used
10 stemmer = PorterStemmer()
11
12 def preprocess_text(text):
13     """
14     Preprocess text to use it for our trained model
15     => tokenize => to_lower => apply stopwords => apply stemmer
16     => join
17
18     Parameters:
19     text (str): text to process.
20
21     Returns:
22     str: Text that can be send to the model.
23     """
24     tokens = text.split()
25     tokens = [stemmer.stem(token) for token in tokens if token.
26               lower() not in stop_words]
27     return ' '.join(tokens)
```

FIGURE 1.3 – stopword & lemmation

1.1.3.3 Cas du *Body/Post* (non utile pour le moment)

Lors de notre analyse l'utilisation du body pour déterminer les tags a été envisagé. Il faudrait alors retiré les reliquats de balises html présente dans le champ *Body* puis le rajouter

au *Titre* dans une nouvelle variable que l'on appellera *Post*. Il faudra alors appliquer au *Post* le même traitement que celui effectué au *Titre*.

L'utilisation du *Body* n'a pas été plus approfondie, suite à des soucis d'ordre logistique, la machine à disposition mettant plus de 2 heures uniquement pour process la partie de préparation des données.

2 Modèle BERT

Avant toute chose il faut savoir pourquoi l'utilisation de BERT a été choisi.

2.1 Limitations matériel

Pour des raisons logistique, et pour des raisons de délai l'entraînement s'est vu restreindre sur deux point :

1. Limiter le nombre de tags à estimer aux 10 tags les plus fréquent.
2. Réduire la quantité de titres utilisé pour l'entraînement.

2.2 Paramétrage

Pour faire fonctionner notre modèle Bert nous aurons besoin de deux dictionnaire :

1. tag2id pour traduire les tag en id
2. id2tag pour décoder les id sous forme de tag

On utilisera un *BertTokenizer*¹ que l'on applique à nos intitulés de questions. La valeur de `len(unique_tags)` étant ici de 10^2 .

1. De base '*bert-base-uncased*' a été utilisé. Il est possible de prendre d'autres paramètres initiaux pour réduire le temps nécessaire à l'entraînement

2. cf décision prise en 2.1 Limitations matériel

```

1 tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
2
3 model = AutoModelForSequenceClassification.from_pretrained(
4     model_path,
5     num_labels=len(unique_tags),
6     id2label=id2tag,
7     label2id=tag2id,
8     problem_type="multi_label_classification"
9 )
10
11 encoded_texts = tokenizer(texts, truncation=True, padding=True,
12                             max_length=512, return_tensors='pt')

```

FIGURE 2.1 – Définition du modèle et du tokenizer

La fonction *compute_metrics* est utilisé pour calculer les différentes métriques d'évaluation. On a choisi une valeur seuil de 0.5, tout logits ayant une valeur supérieur sera traduite par une évaluation a 1 du tag.

```

1 def compute_metrics(pred):
2     labels = pred.label_ids
3     logits = pred.predictions
4
5     threshold = 0.5
6     preds = (logits > threshold).astype(float)
7
8     precision, recall, f1, _ = precision_recall_fscore_support(
9         labels, preds, average='micro', zero_division=0)
10    acc = accuracy_score(labels, preds)
11
12    return {
13        'accuracy': acc,
14        'f1': f1,
15        'precision': precision,
16        'recall': recall,
17    }

```

Le choix de *BCEWithLogitsLoss* est couramment utilisée pour des tâches de classification binaire où chaque classe est traitée de manière indépendante. On choisira de travailler sur 10 epochs. A chaque fin d'époque on évaluera leur performance et on le sauvegardera toutes les epochs³ dans le dossier renseigné dans *output_dir*. L'époque avec la meilleur accuracy sera celle chargé à la fin de l'entraînement.

3. Attention, la sauvegarde des modèles prend de la place!

On utilisera une *DataCollatorWithPadding*⁴ pour assurer que nos inputs auront la bonne taille. En cas d'input avec une taille inférieure, notre data collector injectera des valeurs pour correspondre à nos attentes en vue du training.

```
1     loss_fn = torch.nn.BCEWithLogitsLoss()
2
3     training_args = TrainingArguments(
4         output_dir="./model_epochs",
5         learning_rate=2e-5,
6         per_device_train_batch_size=16,
7         per_device_eval_batch_size=16,
8         num_train_epochs=10,
9         weight_decay=0.01,
10        evaluation_strategy="epoch",
11        save_strategy="epoch",
12        metric_for_best_model="accuracy",
13        load_best_model_at_end=True,
14    )
15    data_collator = DataCollatorWithPadding(tokenizer)
```

FIGURE 2.2 – Définition des paramètres du trainer

Si l'on désire augmenter le nombre d'époch pour notre entraînement, il faudrait songer à modifier notre *save_strategy*, afin d'effectuer une sauvegarde toutes les X époques afin d'alléger la place prise par notre entraînement ainsi que son temps de calcul.

```
1     trainer = Trainer(
2         model=model,
3         args=training_args,
4         train_dataset=train_dataset,
5         eval_dataset=test_dataset,
6         tokenizer=tokenizer,
7         data_collator=data_collator,
8         compute_metrics=compute_metrics,
9     )
```

FIGURE 2.3 – Définition du trainer avec les paramètres précédemment définis

4. Padding \iff Rembourrage

2.3 Performances

Le modèle a été entraîné sur 10 epochs.

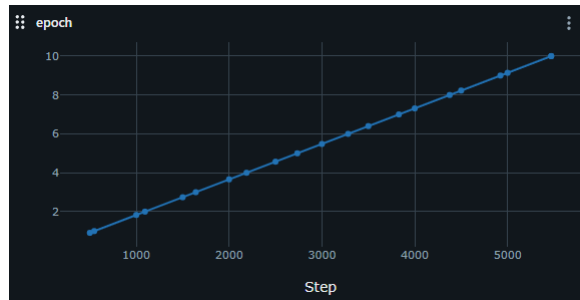


FIGURE 2.4 – Performance observé avec ml flow

Hormis pour *python* & *javascript*, la répartition des tags dans notre jeu de donnée est assez équilibré rendant précision pertinente.

$$Accuracy = \frac{\text{Nombre de predictions correctes}}{\text{Nombre total de predictions}}$$

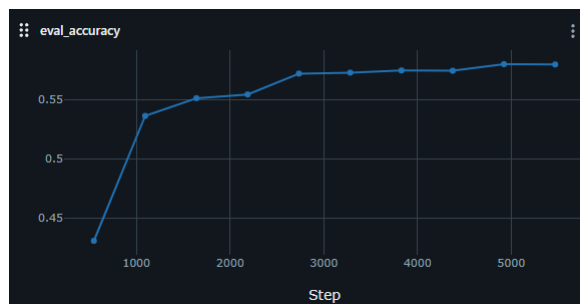


FIGURE 2.5 – Accuracy observé avec ml flow

Le Recall permet d'évaluer la capacité de notre modèle à identifier les instances positives. On arrive donc à prédire 30% des tags, de manière pertinente

$$Recall = \frac{\text{Vrais positifs}}{\text{Vrais positifs} + \text{Faux négatif}}$$

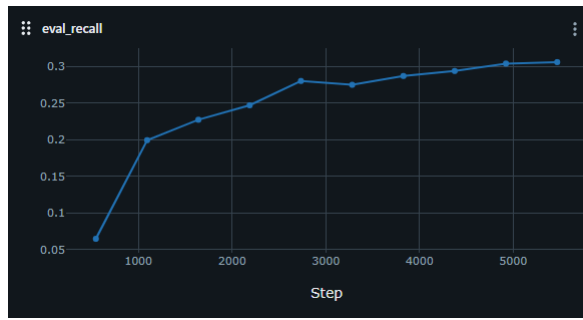


FIGURE 2.6 – Recall observé avec ml flow

Le F1-score est la moyenne harmonique de la précision et du rappel. Ici nous avons une F1-score semblable à notre Recall.

$$F1 - score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

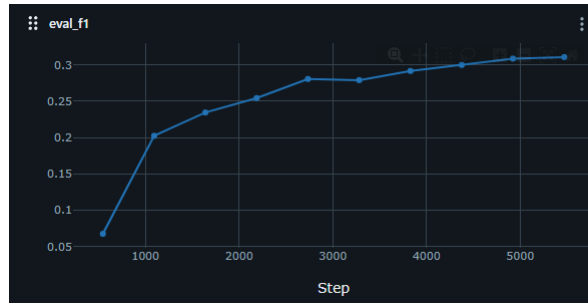


FIGURE 2.7 – F1 score observé avec ml flow

3 Mise en place de l'api

3.1 API

L'API a été développé sur la base /Azure-Samples/msdocs-python-flask-webapp-quickstar, qui est une base par Azure pour avoir une api python utilisant fastapi à déployer sur leur site.

Nous avons légèrement modifié l'API hello world pour qu'elle effectue une prédiction pour un titre donné plutôt que de juste dire bonjour au nom de la personne saisie.

Nous avons donc rajouté le code vu dans la figure 1.3 ainsi que de charger le modèle que l'ont aura au préalable sauvegardé dans le dossier ./model⁵

```
1 model_path = './model'
2
3 model = BertForSequenceClassification.from_pretrained(model_path)
4 tokenizer_config_file = os.path.join(model_path, "
    tokenizer_config.json")
5 with open(tokenizer_config_file, "r") as f:
6     tokenizer_config_data = f.read()
7     print("Content of tokenizer_config.json:")
8     print(tokenizer_config_data)
9
10 tokenizer = BertTokenizer.from_pretrained(model_path)
11 print(tokenizer) # check if the tokenizer is initialized
12
13 id2label = model.config.id2label
```

FIGURE 3.1 – Définition du modèle et du tokenizer

On passera au modèle des titres qui auront donc traité avec stopwords et lemmation et on récupérera une liste de tags. Il sera aussi possible de définir un seuil⁶ qui définira la certitude minimale des tags en sortie. Plus la valeur est proche de 1, plus le modèle est sûr de la

5. Penser à utiliser git lfs pour les fichiers volumineux

présence du tag.

```
1 def predict(texts, threshold=0.2):
2     # Appliquer le pr traitement
3     processed_texts = [preprocess_text(text) for text in texts]
4
5     # Tokenizer les entr es
6     print(processed_texts)
7     inputs = tokenizer(processed_texts, padding=True, truncation=
8     True, return_tensors="pt")
9     with torch.no_grad():
10         outputs = model(**inputs)
11
12     # Appliquer la fonction sigmo de pour obtenir les
13     probabilit s
14     probabilities = torch.sigmoid(outputs.logits).cpu().numpy()
15
16     # Convertir les probabilit s en labels
17     predictions = []
18     for probs in probabilities:
19         labels = [id2label[idx] for idx, prob in enumerate(probs)
20         if prob > threshold]
21         predictions.append(labels)
22     return predictions
```

FIGURE 3.2 – Prédire les tags avec l'API

Comme on s'est basé sur une base possédant un affichage web, on a défini 2 cas d'usage pour la détermination des tags en fonction d'un titre. Soit on passe par le navigateur via

```
1 @app.get("/get_tags", response_class=HTMLResponse)
2 async def get_tags(request: Request, title: str = Query(...)):
3     predictions = predict([title])
4     return templates.TemplateResponse('hello.html', {"request":
5     request, "predictions": predictions[0], 'title': title})
```

FIGURE 3.3 – call web pour avoir les tag via navigation

Ou bien par le call directement auprès de l'API sans passer par aucune interface

6. threshold $\in [0, 1]$


```

1 @app.post("/get_tags_api")
2 async def get_tags_api_endpoint(request: TitleRequest):
3     title = request.title
4     threshold = request.threshold
5     tags = get_tags_api(title, threshold)
6     return {"tags": tags}

```

FIGURE 3.4 – appel API sans interface

3.2 Tests unitaire

Des tests unitaire ont été mis en place. Ils sont effectué à chaque commit sur la branche principale via github action.

	Test	effectue
1	index	la racine répond ?
2	preprocess_text	test de la fonction qui effectue stopwords + lemmation
3	get_tags_api_endpoint	appel réussi via API
4	get_tags_api_endpoint_invalid	appel raté car titre vide via
5	get_tags_api_function	test de la fonction de prédiction appelé par l'api
6	get_tags_api_endpoint_valid_threshold	appel api réussi avec un threshold
7	get_tags_api_endpoint_invalid_threshold_below	appel api threshold<0 est interdit
8	get_tags_api_endpoint_invalid_threshold_above	appel api threshold>1 est interdit
9	hello	hello.html est le point d'entré web
10	get_tags	/get_tags?title répond avec un title
11	get_tags_empty_title	/get_tags?title répond sans un title

```

Run tests
1 ▶ Run source venv/bin/activate
11 ===== test session starts =====
12 platform linux -- Python 3.10.14, pytest-8.2.2, pluggy-1.5.0
13 rootdir: /home/runner/work/Project5-API/Project5-API
14 plugins: anyio-4.4.0
15 collected 11 items
16
17 test_main.py ..... [100%]
18
19 ===== warnings summary =====
20 main.py:93
21 /home/runner/work/Project5-API/Project5-API/main.py:93: PydanticDeprecatedSince20: Pydantic V1 style '@validator' validators are deprecated. You should migrate to Pydantic V2 style '@field_validator' validators,
22 see the migration guide for more details. Deprecated in Pydantic V2.0 to be removed in V3.0. See Pydantic V2 Migration Guide at https://errors.pydantic.dev/2.9/migration/
23 @validator('title')
24
25 main.py:99
26 /home/runner/work/Project5-API/Project5-API/main.py:99: PydanticDeprecatedSince20: Pydantic V1 style '@validator' validators are deprecated. You should migrate to Pydantic V2 style '@field_validator' validators,
27 see the migration guide for more details. Deprecated in Pydantic V2.0 to be removed in V3.0. See Pydantic V2 Migration Guide at https://errors.pydantic.dev/2.9/migration/
28 @validator('threshold')
29
30 test_main.py::test_index
31 test_main.py::test_hello
32 test_main.py::test_get_tags
33 test_main.py::test_get_tags_empty_title
34 /home/runner/work/Project5-API/Project5-API/venv/lib/python3.10/site-packages/starlette/templating.py:178: DeprecationWarning: The 'name' is not the first parameter anymore. The first parameter should be the
35 'Request' instance.
36 Replace 'TemplateResponse(name, ***request*: request*)' by 'TemplateResponse(request, name)'.
37 warnings.warn(
38
39 -- Docs: https://docs.pytest.org/en/stable/how-to/capture-warnings.html
40 ===== 11 passed, 6 warnings in 6.95s =====

```

FIGURE 3.5 – Tests unitaire via github action

On remarquera la présence de warning qu'il faudra, à l'avenir, retirer.

3.3 Déploiement

Un déploiement automatisé via github action sur la plateforme Azure avait été mis en chantier mais suite à des erreurs de manipulation, la quantité de data disponible pour l'upload via git a été saturé.

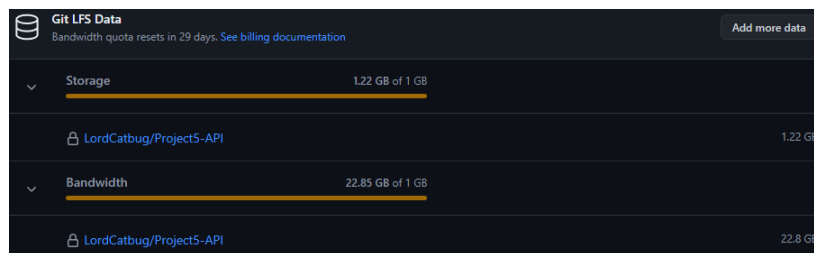


FIGURE 3.6 – Statuts des data upload via github

Le déploiement s'est donc fait via une encapsulation docker ainsi que via un push via le client azure AD.

Le ressource groupe *rg-project5* a été créé contenant :

- Container instances : *registryfastapi*.

- Container registry : *registryfastapi* : .

```
1 #créer l'image pour docker
2 docker build -t myimage .
3
4 #run l'image locallly
5 docker run -d --name mycontainer 80:80 myimage
6
7 # cf figure3.7 pour avoir les informations
8 docker login registryfastapi.azurecr.io -u registryfastapi -p "
  password"
9
10 # changer le X de buildX par le num ro de build/version
11 docker build -t registryfastapi.azurecr.io/rlfastapi:buildX
12 docker push registryfastapi.azurecr.io/rlfastapi:buildX
```

FIGURE 3.7 – Commande terminal pour build et commit l'image docker

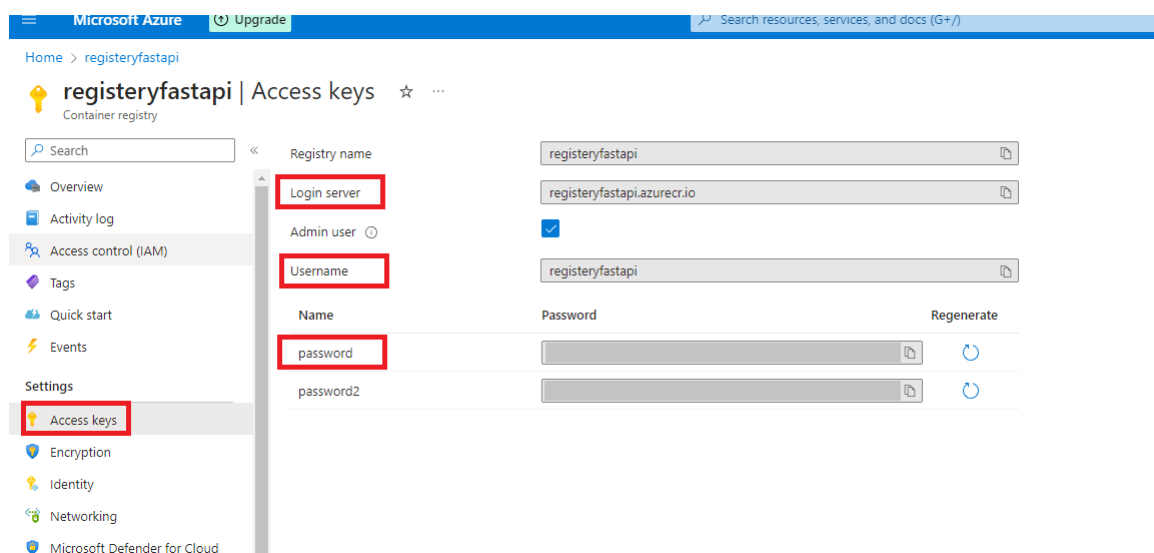


FIGURE 3.8 – Récupérer le password pour upload docker image

Conclusion

Nous avons entraîné un modèle et avons rendu possible l'appel à celui-ci via des appels Api auprès d'Azure.

Les résultats du modèle ne sont pas parfait mais son encourageant. Une des pistes envisageable pour améliorer les performances serait d'intégrer les *Body* afin d'avoir une description plus précise du problème. Cependant, cela induira un temps de pré-traitement et d'entraînement drastiquement plus long.

Nos résultats ne sont compris que dans une liste de 10 tags, il serait bon de réentraîner un modèle avec une liste de tag plus conséquente.

Il est aussi à noter que nous utilisons Azure avec un compte d'essai, ce qui nous limite dans le type de machine dont nous avons accès, se répercutant sur le temps de traitement de nos requêtes