

5 - REFERENTIELS

Article L6113-1 [En savoir plus sur cet article...](#) Créé par [LOI n°2018-771 du 5 septembre 2018 - art. 31 \(V\)](#)

« Les certifications professionnelles enregistrées au répertoire national des certifications professionnelles permettent une validation des compétences et des connaissances acquises nécessaires à l'exercice d'activités professionnelles. Elles sont définies notamment par un **référentiel d'activités** qui décrit les situations de travail et les activités exercées, les métiers ou emplois visés, un **référentiel de compétences** qui identifie les compétences et les connaissances, y compris transversales, qui en découlent et un **référentiel d'évaluation** qui définit les critères et les modalités d'évaluation des acquis. »

Chaque bloc de compétences est évalué sous la forme de projet, accompagné et validé par un mentor (professionnel du métier), puis présenté au jury du titre.

REFERENTIEL D'ACTIVITES <i>décrit les situations de travail et les activités exercées, les métiers ou emplois visés</i>	REFERENTIEL DE COMPETENCES <i>identifie les compétences et les connaissances, y compris transversales, qui découlent du référentiel d'activités</i>	REFERENTIEL D'EVALUATION <i>définit les critères et les modalités d'évaluation des acquis</i>	
		MODALITÉS D'ÉVALUATION	CRITÈRES D'ÉVALUATION
BLOC 1 – Prétraiter et analyser des données structurées pour répondre à un problème métier.			
1.1- Cadrage et problématisation d'une analyse de données.	Déterminer les objectifs d'une analyse de données à partir d'un problème métier.	Mise en situation professionnelle (projet) Et Soutenance orale	La demande du client interne/ externe est prise en compte.
			La problématique métier est correctement identifiée.
1.2- Nettoyage de données (data cleaning).	Effectuer des opérations de nettoyage sur des données structurées.		Les objectifs de l'analyse de données à mener sont pertinents et permettent de répondre à la problématique métier.
			Les données manquantes sont correctement identifiées. Des actions correctives sont proposées.

			<p>Les duplicats de variables et d'enregistrement sont traités.</p> <p>Les variables non pertinentes pour la problématique sont éliminées.</p>
<p>1.3- Analyse exploratoire de données structurées.</p>	<p>Effectuer une analyse statistique univariée à partir de données structurées et nettoyées.</p>		<p>Les métriques adaptées sont utilisées.</p> <p>Les valeurs aberrantes sont identifiées. Des actions correctives sont proposées.</p> <p>Les distributions observées sont correctement caractérisées.</p> <p>Les quantiles sont correctement définis.</p>
	<p>Effectuer une analyse statistique multivariée à partir de données structurées et nettoyées.</p>		<p>L'analyse effectuée est pertinente au regard de la problématique.</p> <p>Les résultats de l'analyse sont corrects.</p>
	<p>Représenter des données grâce à des graphiques clairs et pertinents.</p>		<p>Les graphiques produits comprennent au moins un graphique de chaque parmi les suivants : box-plot (boîte à moustaches), barplot (graphique en barres), pie chart (diagramme circulaire),</p>

			<p>histogramme, scatter plot (nuage de points).</p> <p>Les graphiques sont lisibles.</p>
1.4- Formulation de préconisations pour une analyse automatisée des données.	Formuler des préconisations pour un traitement des données permettant au client interne/ externe d'en automatiser l'analyse.		<p>Des propositions sont formulées pour un traitement potentiel des données.</p> <p>Ces propositions sont cohérentes et pertinentes.</p>
BLOC 2 – Entraîner un modèle d'apprentissage automatique supervisé pour réaliser une analyse prédictive.			
2.1- Mise en place d'un modèle d'apprentissage automatique supervisé.	Sélectionner et/ ou transformer les variables pertinentes pour la modélisation (<i>feature engineering</i>).	Mise en situation professionnelle (projet) Et Soutenance orale	<p>Les variables catégorielles sont identifiées et transformées de façon adaptée.</p> <p>Si nécessaire, les variables sont normalisées.</p> <p>Si nécessaires, des transformations mathématiques sont utilisées.</p>
	Sélectionner et mettre en place un modèle d'apprentissage supervisé adapté à une problématique métier.		<p>La (les) variable(s) cibles sont correctement choisies et corrélées.</p> <p>Au moins un algorithme de régression linéaire et un algorithme de régression non linéaire sont testés.</p>

2.2- Évaluation et amélioration d'un modèle d'apprentissage automatique supervisé.	Évaluer les performances d'un modèle d'apprentissage supervisé.		<p>La métrique choisie est adaptée (RMSE, écart moyen).</p> <p>D'autres indicateurs de performance que le score sont utilisés pour analyser le résultat (coefficients des variables en fonction de la pénalisation, visualisation des erreurs en fonction des variables du modèle, temps de calcul...)</p> <p>Les données sont séparées en train/test pour les évaluer de façon pertinente et détecter le surapprentissage (<i>overfitting</i>).</p> <p>Une <i>baseline data</i> est mise en place pour évaluer le pouvoir prédictif du modèle choisi (variable qualitative – <i>dummy regressor</i>).</p> <p>L'ensemble des résultats sont présentés en allant des modèles les plus simples aux plus complexes</p>

	Adapter les paramètres d'un modèle d'apprentissage supervisé afin de l'améliorer.		<p>Une validation croisée est mise en place pour optimiser les hyperparamètres.</p> <p>Les hyperparamètres pertinents sont optimisés dans les différents algorithmes.</p> <p>Le choix final d'algorithme et d'hyperparamètres est justifié.</p>

BLOC 3 – Entraîner un modèle d'apprentissage non supervisé adapté à une problématique de segmentation ou de réduction de données.

3.1- Mise en place d'un modèle d'apprentissage automatique non supervisé.	Sélectionner, transformer et créer les variables pertinentes pour la modélisation (<i>feature engineering</i>).	Mise en situation professionnelle (projet) Et Soutenance orale	<p>Les variables pertinentes sont transformées pour permettre leur exploitation (variables catégorielles en particulier).</p> <p>Une ou plusieurs variables pertinentes pour l'amélioration de la solution sont créées.</p>
--	---	--	---

	<p>Sélectionner et mettre en place un modèle d'apprentissage non supervisé adapté une problématique métier.</p>		<p>Le nombre de segments et la répartition bien choisis.</p> <p>La stratégie d'ajout de nouveaux segments est explicitée.</p> <p>La nature des variables d'entrée est prise en compte dans le choix de l'algorithme et de la distance.</p>
<p>3.2- Évaluation et amélioration d'un modèle d'apprentissage automatique non supervisé.</p>	<p>Évaluer les performances d'un modèle d'apprentissage non supervisé.</p>		<p>La forme des clusters est évaluée.</p> <p>La stabilité des clusters est évaluée.</p> <p>La compatibilité avec des connaissances extérieures ("test set" et/ou a priori) est évaluée.</p>
	<p>Adapter les paramètres d'un modèle d'apprentissage non supervisé afin de l'améliorer.</p>		<p>Les étapes d'évaluation sont automatisées et permettent de pouvoir tester facilement plusieurs combinaisons de paramètres.</p> <p>Le choix des valeurs de paramètres testés est pertinent.</p>

BLOC 4 – Prétraiter et analyser des données non structurées (texte, images).			
4.1- Collecte et prétraitement de données non structurées.	Collecter des données répondant à des critères définis <i>via</i> une interface de programmation (API).	Mise en situation professionnelle (projet) Et Soutenance orale	<p>Les conditions d'utilisation de l'API dans la documentation (limite du nombre d'appels et limitations d'utilisation) sont correctement identifiées.</p> <p>La structure de l'API et les champs contenant les données d'intérêt <i>via</i> la documentation sont correctement identifiés.</p> <p>La requête permettant de récupérer les données d'intérêt <i>via</i> l'API est écrite.</p>
	Prétraiter des données textuelles non structurées pour obtenir un jeu de données exploitable.		<p>Les champs de texte sont nettoyés (retirer la ponctuation, les mots de liaison, mettre tout en minuscule, ...)</p> <p>Une fonction permettant de tokeniser une phrase a été écrite et fonctionne correctement.</p> <p>Une fonction permettant de <i>stemmer</i> (raciniser) une phrase a été écrite et fonctionne</p>

			<p>correctement. Une fonction permettant de lemmatiser une phrase a été écrite et fonctionne correctement.</p> <p>Des variables (<i>features</i>) <i>bag-of-word</i> sont construites (avec étapes de nettoyage supplémentaires : seuil de fréquence des mots, normalisation des mots (racines, utilisation du package NLTK).</p> <p>Une phrase (ou un court texte) d'exemple permet d'illustrer et de tester la bonne réalisation des 5 étapes précédentes.</p>
	<p>Prétraiter des données sous forme d'images non structurées pour obtenir un jeu de données exploitable.</p>		<p>Des bibliothèques spécialisées sont utilisées pour un premier traitement du contraste (ex : openCV).</p> <p>Le bruit est filtré.</p> <p>L'histogramme est égalisé.</p> <p>Une fonction permettant d'extraire des <i>features</i> (via l'un des algorithmes suivants : ORB, SIFT, SURF ou réseaux de</p>

			<p>neurones comme par exemple CNN) a été écrite et fonctionne correctement.</p> <p>Une fonction permettant de construire des <i>features</i> de type "bag-of-image" a été écrite et fonctionne correctement.</p>
<p>4.2- Analyse exploratoire de données de grande dimension.</p>	<p>Réduire la dimension de données de grande dimension afin d'optimiser les temps de calcul.</p>		<p>La nécessité de réduction de dimension (en particulier dans le cas de données image et texte) est justifiée (coût d'acquisition, de stockage, complexité de calcul...)</p> <p>Une méthode de méthode de réduction de dimension (ex : ACP) pour des données image et texte est appliquée.</p> <p>Le choix des valeurs des paramètres dans la méthode de réduction de dimension retenue (ex : le nombre de dimensions conservées pour l'ACP) est justifié.</p>
	<p>Représenter graphiquement des données de grande dimension afin d'en réaliser une analyse exploratoire.</p>		<p>Au moins un graphique représentant les informations contenues dans des données à plus de 2D a été réalisé.</p>

			<p>Le graphique réalisé est lisible et compréhensible.</p> <p>Les différents éléments graphiques sont expliqués (variables représentées sur les axes, par les couleurs, la taille, ...)</p>
BLOC 5 – Présenter et déployer un modèle d'apprentissage automatique auprès de ses utilisateurs finaux.			
5.1- Mise en production d'un modèle d'apprentissage automatique.	Déployer un modèle <i>via</i> une interface de programmation (API) dans le web.	Mise en situation professionnelle (projet) Et Soutenance orale	<p>Un fichier (par exemple format pickle) contenant un modèle de machine learning sérialisé a été créé et le modèle chargé depuis le fichier fonctionne.</p> <p>Le modèle d'apprentissage automatique (machine learning) est déployé sous forme d'API (<i>via</i> Flask par exemple).</p> <p>L'API renvoie effectivement la prédiction correspondante à un client en réponse à un identifiant client.</p> <p>L'API est déployée dans le web <i>via</i> un outil gratuit et disponible plusieurs mois (en vue du jury).</p>

<p>5.2- Formalisation et présentation d'une démarche de modélisation et de ses résultats.</p>	<p>Réaliser un tableau de bord (<i>dashboard</i>) pour présenter son travail de modélisation.</p>		<p>Au moins un parcours utilisateur simple permettant de répondre aux besoins des utilisateurs (les différentes actions/clics sur les différents graphiques permettant de répondre à une question que se pose l'utilisateur) a été décrit et conçu.</p> <p>Au moins deux graphiques interactifs permettant aux utilisateurs d'explorer les données clients (permettant de répondre à des questions de type "quel est le client avec le plus de transactions ?") ont été développés.</p> <p>Les graphiques réalisés sont lisibles, clairs et pertinents.</p> <p>Le <i>dashboard</i> est accessible pour d'autres utilisateurs sur leurs postes de travail (déploiement dans le web).</p>
	<p>Réaliser la présentation orale d'une démarche de modélisation à un client interne/ externe.</p>		<p>La démarche de modélisation est détaillée avec :</p> <ul style="list-style-type: none"> - La méthode d'évaluation de la performance du modèle de

			<p>machine learning.</p> <ul style="list-style-type: none"> - La façon d'interpréter les résultats du modèle. - La façon d'interpréter l'importance des variables du modèle. <p>Les explications sont claires et fournies avec assurance.</p> <p>Les explications sont compréhensibles par un public non technique.</p>
	<p>Rédiger une note méthodologique afin de communiquer sa démarche de modélisation.</p>		<p>La démarche de modélisation est présentée de manière synthétique dans la note (2 pages max).</p> <p>La fonction coût, l'algorithme d'optimisation et la métrique d'évaluation sont explicités (1 page max).</p> <p>L'interprétabilité du modèle est explicitée (la façon d'interpréter l'importance des variables n'est pas la même</p>

			<p>pour une régression logistique que pour un random forest) et les limites éventuelles sont précisées (1 page max).</p> <p>Le limites et les améliorations envisageables pour gagner en performance et en interprétabilité de l'approche de modélisation sont décrites (1 page max).</p>
	<p>Assurer l'intégration du modèle auprès de collaborateurs en utilisant un logiciel de version de code.</p>		<p>Un dossier contenant tous les scripts du projet a été créé dans un logiciel de version de code (ex : Github).</p> <p>L'historique des modifications du projet affiche au moins trois versions distinctes et l'on peut accéder à ces anciennes versions.</p> <p>La liste des packages utilisés ainsi que leur numéro de version est disponible et tenu à jour.</p> <p>Le travail est réutilisable par d'autres personnes et la collaboration est facilitée :</p> <ul style="list-style-type: none">- un fichier introductif permet de comprendre l'objectif du

			<p>projet et le découpage des dossiers</p> <p>- les scripts et les fonctions sont commentés.</p>
<p>BLOC 6 – Déployer un modèle d'apprentissage automatique à l'échelle en utilisant les technologies du <i>Big data</i>.</p>			
<p>6.1- Analyse et modélisation de données dans un environnement <i>Big data</i>.</p>	<p>Sélectionner les outils du <i>Cloud</i> permettant de disposer d'un environnement <i>Big Data</i>.</p>	<p>Mise en situation professionnelle (projet)</p> <p>Et</p> <p>Soutenance orale</p>	<p>Les différentes briques d'architecture nécessaires pour la mise en place d'un environnement <i>Big data</i> sont identifiées.</p> <p>Les outils du cloud permettant de mettre en place l'environnement <i>Big Data</i> sont identifiés.</p>
<p>6.2- Réalisation d'opérations de calcul de grands volumes de données.</p>	<p>Prétraiter, analyser et modéliser des données dans un environnement <i>Big data</i> en utilisant les outils du <i>Cloud</i>.</p>		<p>Les fichiers (de départ et ceux après transformation) sont chargés dans un espace de stockage cloud.</p> <p>Les scripts ont été exécutés en utilisant des machines dans le cloud.</p> <p>Un des scripts permet d'écrire les sorties du programme directement dans l'espace de stockage cloud.</p>

	Réaliser des calculs distribués sur des données massives en utilisant les outils adaptés.		Les traitements critiques lors d'un passage à l'échelle en termes de volume de données sont identifiés. Les scripts sont développés en Pyspark (API).
--	---	--	--