

# Ingénieur.e en apprentissage automatique Data scientist

 CDD 3 ans

 Début : dès que possible

 Villejuif

 Télétravail partiel

 Bac +5 ou plus

L'Inserm est le seul organisme public français entièrement dédié à la recherche biologique, médicale et en santé des populations. Il dispose de laboratoires de recherche sur l'ensemble du territoire, regroupés en 12 Délégations Régionales. Notre institut réunit 15 000 chercheurs, ingénieurs, techniciens et personnels administratifs, avec un objectif commun : améliorer la santé de tous par le progrès des connaissances sur le vivant et sur les maladies, l'innovation dans les traitements et la recherche en santé publique.

Rejoindre l'Inserm, c'est intégrer un institut engagé pour la parité et l'égalité professionnelle, la diversité et l'accompagnement de ses agents en situation de handicap, dès le recrutement et tout au long de la carrière. Afin de préserver le bien-être au travail, l'Inserm mène une politique active en matière de conditions de travail, reposant notamment sur un juste équilibre entre vie personnelle et vie professionnelle.

L'Inserm a reçu en 2016 le label européen HR Excellence in Research et s'est engagé à faire évoluer ses pratiques de recrutement et d'évaluation des chercheurs.

## Emploi

**Poste ouvert aux candidats**

- Agents fonctionnaires de l'Inserm par voie de mobilité interne
- Agents fonctionnaires non Inserm par voie de détachement
- CDD agents contractuels

**Catégorie** A

**Corps** IR

**Emploi-Type** Expert en information statistique  
E1D44

## Structure d'accueil

**Unité**

CépiDc-US10

**A propos de la Structure**

Le CépiDc, unité de service de l'Inserm, a pour mission de produire la base de données statistique sur les causes médicales de décès en France, de la diffuser et de réaliser des analyses sur cette base de données. Cette base de données statistique repose sur la collecte et le traitement des volets médicaux des certificats de décès. Ses finalités d'usage sont multiples : la connaissance de l'état de santé de la France et de ses territoires et sa comparaison au niveau international, en vue d'aider au pilotage des politiques de santé publique ; la recherche et les études, les données alimentant le système national des données de santé ; la veille et l'alerte sanitaire, par la production de la donnée la plus pertinente possible dans des délais de quelques jours. Les principaux traitements réalisés au CépiDc concernent d'abord l'accueil, le contrôle et l'intégration des données collectées, pour une mise à disposition immédiate à des

fins de veille de sanitaire. Puis, le CépiDc assure le dédoublonnage et la correction de non-réponse totale *via* une mise en cohérence avec les décès déclarés à l'état civil et gérés par l'Insee, pour finir par la construction des variables statistiques, avec en particulier le codage des causes médicales de décès.

Concernant ce dernier aspect, il s'agit d'analyser et de coder les textes rédigés par les professionnels de santé lors du constat de décès dans la classification internationale des maladies (CIM). La stratégie de codage combine de façon optimale trois modes : utilisation d'un système-expert de règles en batch et en interactif et utilisation de modèles d'intelligence artificielle prédictifs, construits et entraînés *in-house*. Capitalisant sur les millions d'observations analysées par des experts suivant des standards internationaux et dans le contexte d'une profonde rénovation du processus de production des causes de décès, le CépiDc intègre ces méthodes dans sa chaîne de production pour gagner en temps (afin de respecter les délais réglementaires de diffusion des données) et en qualité, tout en adoptant une démarche statistique rigoureuse et novatrice. L'élaboration de la base de données sur les causes de décès suit les recommandations de l'OMS, et doit satisfaire les normes de qualité d'une statistique officielle et du code des bonnes pratiques en matière de statistique européenne.

Le CépiDc est composé d'une vingtaine d'agents, répartis en deux pôles : pôle production des données et pôle exploitation-diffusion.

<b>Directrice</b>	Hélène Chaput
<b>Adresse</b>	Bâtiment Laplace, Hôpital Paul Brousse, 12 Av. Paul Vaillant Couturier, 94800 Villejuif
<b>Délégation Régionale</b>	Paris Sud
<b>Description du poste</b>	
<b>Mission principale</b>	<p>Vous participez à la mise en œuvre en production courante des outils d'intelligence artificielle pour le codage des causes de décès. Ces outils fondés sur de l'apprentissage profond et du traitement automatique des langues améliorent la qualité et la rapidité de codage automatique, de façon à satisfaire les délais réglementaires de diffusion de la base. Vous êtes responsable de leur adaptation pour tenir compte du prochain changement de nomenclature (passage de la CIM 10 à la CIM 11) et vous êtes partie prenante de l'évolution du système d'information du CépiDc qui en découle. Vous bénéficiez d'un accès à des ressources de calcul (GPU) permettant de concevoir, entraîner et tester des modèles et de réaliser des prédictions.</p> <p>Au sein du pôle production des données du CépiDc, vous travaillez dans l'équipe automatisation, sous la responsabilité de la cheffe d'équipe, et en étroite collaboration avec le <i>data scientist senior</i>. Vous collaborez également avec le reste de l'équipe pluridisciplinaire (codeurs, nosologistes, responsables de production, statisticiens) et vous êtes partie prenante de l'écosystème formé avec les partenaires de recherche et développement (médecins spécialisés en informatique médicale et data scientists, de l'AP-HP, LISN-CNRS, Insee, Santé publique France, Inserm).</p>
<b>Activités principales</b>	<ul style="list-style-type: none"> <li>• Mettre en production, maintenir, monitorer et valider une chaîne de traitements de données textuelles comprenant des prédicteurs de type réseaux de neurones (<i>transformers</i>) pour aider/automatiser le codage du texte libre des certificats de décès dans la CIM (annotation, training/fine-tuning, monitoring).</li> <li>• Mettre en production le ciblage des certificats à allouer aux différentes modalités de codage (IA, manuel), évaluer l'amélioration continue du codage automatique (en taux de codage et en qualité) en vue d'une boucle d'apprentissage continue (on line) à partir de la validation/correction des codeurs des propositions de l'algorithme.</li> <li>• Adapter l'architecture du modèle et le <i>feature engineering</i> en vue d'améliorer la classification des causes, en adéquation avec la finalité statistique du traitement et les bonnes pratiques.</li> <li>• Participer à l'internationalisation de ces méthodes en lien avec les instances représentatives</li> </ul>

<p>françaises à l'OMS et au sein de l'Europe.</p> <ul style="list-style-type: none"> <li>Assurer une veille scientifique sur les modèles et les algorithmes à l'état de l'art dans le domaine.</li> <li>Participer activement à des groupes d'échanges de bonnes pratiques existants ou à construire regroupant data-scientists, statisticiens et chercheurs en épidémiologie et informatique (Insee, DREES, Inserm, Inria,...) autour de l'usage de l'IA/TAL sur ces thématiques.</li> </ul>
<p><b>Spécificité(s) et environnement du poste</b></p> <ul style="list-style-type: none"> <li>Confidentialité des données</li> <li>Contraintes de production.</li> </ul>
<p><b>Connaissances</b></p> <ul style="list-style-type: none"> <li>Apprentissage automatique, traitement automatique des langues, <i>deep learning</i>, sciences des données</li> <li>Maîtrise de l'ensemble des étapes allant du développement à la mise en production</li> <li>Maîtrise des environnements de production</li> <li>De bonnes bases statistiques</li> <li>Des connaissances en biostatistique et un intérêt pour l'épidémiologie sont des plus</li> </ul>
<p><b>Savoir-faire</b></p> <ul style="list-style-type: none"> <li>Très bonne maîtrise de Python et des librairies de <i>deep learning</i> (Tensorflow, Pytorch) en particulier celles appliquées au traitement automatique des langues.</li> <li>Entraînement et monitoring d'algorithmes de <i>deep learning</i></li> <li>Mise en production d'algorithmes de <i>machine learning</i>, MLops</li> <li>Git, outil de versioning</li> <li>Design et maintien de pipeline de <i>machine learning</i>, ces expériences sont des plus, de même que l'utilisation de Docker, MLFlow, et de technologies cloud</li> </ul>
<p><b>Aptitudes</b></p> <ul style="list-style-type: none"> <li>Proactivité, force de proposition</li> <li>Aisance relationnelle, sens de la communication et de la pédagogie</li> <li>Capacités d'organisation, de planification et de rigueur</li> <li>Discrétion et confidentialité</li> <li>Savoir s'insérer et interagir avec des équipes multidisciplinaires : pôle de production, experts métiers chargés de production, statisticiens, stagiaires, chercheurs</li> <li>Savoir se maintenir à l'état de l'art des connaissances</li> </ul>
<p><b>Expérience(s) souhaité(s)</b></p> <ul style="list-style-type: none"> <li>Ce poste convient à un sortant d'école motivé, formé à l'usage de Python et des librairies d'apprentissage profond.</li> </ul>
<p><b>Niveau de diplôme et formation(s)</b></p> <ul style="list-style-type: none"> <li>Diplôme d'ingénieur de grandes écoles, Master en data science ou équivalence professionnelle</li> </ul>

### Informations Générales

<p><b>Date de prise de fonction</b></p> <p>Dès que possible</p>
<p><b>Durée (CDD et détachements)</b></p> <p>36 mois</p> <p>Renouvelable : <input checked="" type="checkbox"/> OUI <input type="checkbox"/> NON</p>
<p><b>Temps de travail</b></p> <ul style="list-style-type: none"> <li>Temps plein</li> <li>Nombre d'heures hebdomadaires : 38h30</li> <li>45 jours congés annuels et RTT</li> </ul>
<p><b>Activités télétravaillables</b></p> <p><input checked="" type="checkbox"/> OUI, en partie <input type="checkbox"/> NON</p> <p>* Préciser les modalités de télétravail possible.</p>
<p><b>Rémunération</b></p> <ul style="list-style-type: none"> <li><b>Selon l'expérience et le profil de candidature</b></li> </ul>

- Prise en charge d'une partie de la mutuelle.

### Modalités de candidature

**Date limite de candidature**

15/02/2026

**Contact**[recrutement.cepидc@inserm.fr](mailto:recrutement.cepидc@inserm.fr) ; [aude.robert@inserm.fr](mailto:aude.robert@inserm.fr)**Contractuels**

- Envoyer CV et lettre de motivation à [recrutement.cepидc@inserm.fr](mailto:recrutement.cepидc@inserm.fr)
- Précisez vos préférences salariales.

**Pour en savoir +**

- Sur l'Inserm : <https://www.inserm.fr/> ; site RH : <https://rh.inserm.fr/Pages/default.aspx>
- Sur la politique handicap de l'Inserm et sur la mise en place d'aménagements de poste de travail, contactez la Mission Handicap : [emploi.handicap@inserm.fr](mailto:emploi.handicap@inserm.fr)