

تحليل البيانات المتدفقة باستخدام التعلم الآلى على الأنظمة الكبيرة

مقدمة من:

فوزية رمضان سيد حسان

مدرس مساعد، قسم علوم الحاسب،
كلية الحاسبات والمعلومات، جامعة الفيوم

رسالة مقدمة الى كلية الحاسبات والمعلومات، جامعة الفيوم
كجزء من متطلبات الحصول على درجة دكتوراة الفلسفة في الحاسبات و
المعلومات

أ.د/ عبدالمجيد أمين على

د/ مسعود إسماعيل مسعود شاهين

قسم علوم الحاسب
كلية الحاسبات والمعلومات
جامعة الفيوم

قسم علوم الحاسب
كلية الحاسبات والمعلومات
جامعة المنيا

تخصص علوم الحاسب

إشراف:

كلية الحاسبات و المعلومات
جامعة الفيوم

ملخص الرسالة

مع الزيادة المستمرة في تدفق البيانات وسرعتها الهائلة في الأنظمة الكبيرة ومع الاحتياجات المتزايدة لتحليل البيانات وخصوصا في مجال الرعاية الصحية. لقد أصبح تحليل تلك البيانات والتنبؤ بالأمراض المزمنة في الوقت الحالي هو أمر ضروري للغاية ، وذلك لأن من الصعب استيعاب ومعالجة وتحليل مثل هذه البيانات الضخمة لأتخاذ إجراء في الوقت الحالي في حالات الطوارئ باستخدام الطرق التقليدية. لذلك ، فإن العمل في هذه الأطروحة يتعلق بكيفية بناء نظام يمكنه تحليل ومعالجة تدفق البيانات في الوقت الحالي باستخدام بيانات وسائل التواصل الاجتماعي المستندة إلى الصحة أو أجهزة الاستشعار الطبية التي يمكن ارتداؤها والتنبؤ بالحالة الحالية لصحة المريض. لذلك فإن العمل في هذه الأطروحة هدفه تحسين تحليل البيانات المتدفقة في الوقت الحالي، وقد تم ذلك من خلال مساهمتين مختلفتين .

المساهمة الأولى وهي تقديم نظام للتنبؤ لمرض السكري عند وصول البيانات من وسائل التواصل الاجتماعي (تويتر). وقد تم تطوير هذا النظام لأستقبال البيانات المتدفقة المستندة إلى الصحة وذلك للتنبؤ بالحالة الصحية للمريض يهدف النظام المقترح إلى إيجاد نموذج التعلم الآلي الأكثر دقة التي لديها أعلى دقة من التنبؤ بمرض السكري. وقد حددت النتائج التجريبية أن نموذج Random Forest حقق أعلى دقة بين النماذج الأخرى عند 84.11%. للتنبؤ عبر الإنترنت من خلال وسائل التواصل الاجتماعي ، لقد تم تنفيذ نظام مقترح للتعامل مع بيانات تويتر المتدفقة الخاصة بصحة المرضى. فلقد تم دمج تدفق Kafka و Spark في الواجهة الخلفية للنظام المقترح. ومن ثم، يتم استخدام Random Forest للتنبؤ بالحالة الصحية الحالية للمريض في الوقت الحالي.

المساهمة الثانية وهي عبارة عن نظام Online Prediction System الذي قد تم تطويره في هذه الأطروحة لحل مشكلة تحليل البيانات المتدفقة في الوقت الحالي من خلال تطبيق ال Streaming machine learning على البيانات الصحية المتدفقة التي يتم تطبيقها على Spark Streaming من خلال Kafka. ولقد تم إجراء النتائج التجريبية على مجموعات مختلفة للبيانات الطبية التاريخية وتحويلها لبيانات متدفقة وايضا بيانات أجهزة الاستشعار الطبية القابلة للارتداء. وقد أثبتت النتائج التجريبية أن Online Prediction System قادر على التعلم وتحديث النموذج وفقا لوصول البيانات الجديدة وحجم النافذة.

