

Abstract - Español

Los procesos Gaussianos (Gaussian Processes, GPs) son modelos Bayesianos no paramétricos que representan el actual estado del arte en tareas de aprendizaje supervisado tales como regresión y clasificación. Por este motivo, son uno de los bloques básicos usados en la construcción de otros algoritmos de aprendizaje máquina más sofisticados. Asimismo, los GPs tienen una variedad de propiedades muy deseables: Son prácticamente inmunes al sobreajuste, disponen de mecanismos sensatos y cómodos para la selección de modelo y proporcionan las llamadas "barras de error", es decir, son capaces de estimar la incertidumbre de sus propias predicciones.

Desafortunadamente, los GPs completos no pueden aplicarse directamente a bases de datos de gran tamaño, cada vez más frecuentes en la actualidad. Para n muestras, el tiempo de cómputo necesario para entrenar un GP escala como $O(n^3)$, lo que hace que un ordenador doméstico actual sea incapaz de manejar conjuntos de datos con más de unos pocos miles de muestras. Para solventar este problema se han propuesto recientemente varias aproximaciones "dispersas", que escalan linealmente con el número de muestras. De entre éstas, el método conocido como "procesos Gaussianos dispersos usando pseudo-entradas" (Sparse Pseudo-inputs GP, SPGP), representa el actual estado del arte. Aunque este tipo de aproximaciones dispersas permiten tratar bases de datos mucho mayores, obviamente no alcanzan el rendimiento de los GPs completos

En esta tesis se introducen varios modelos de GP disperso que presentan un rendimiento mayor que el del SPGP, tanto en cuanto a capacidad predictiva como a calidad de las barras de error. Los modelos propuestos convergen al GP completo que aproximan bajo determinadas condiciones, pero el objetivo de esta tesis no es tanto aproximar fielmente el GP completo original como proporcionar modelos prácticos de alta capacidad predictiva. Tanto es así que, en ocasiones, los nuevos modelos llegan a batir al GP completo que los inspira.

Se proporcionan dos clases generales de modelos: Redes marginalizadas (Marginalized Networks, MNs) y GPs inter-dominio (Inter-Domain GPs, IDGPs). Las MNs pueden verse como modelos que se encuentran a mitad de camino entre las redes neuronales clásicas (Neural Networks, NNs) y los GPs completos, intentando combinar las ventajas de ambos. Aunque la fase de entrenamiento de una MN es diferente, cuando se utiliza para predicción mantiene la estructura de una NN clásica, de manera que las MNs pueden ser interpretadas como una manera novedosa de entrenar NNs clásicas, al tiempo que se añaden beneficios adicionales, como resistencia al sobreajuste y "barras de error" dependientes de la entrada. Los IDGPs generalizan el SPGP, permitiendo a las "pseudo-entradas" residir en un dominio diferente del de entrada, incrementado así la flexibilidad y el rendimiento. Además, proporcionan un marco probabilístico adecuado para entender modelos dispersos anteriores

Así pues, todos los algoritmos propuestos son puestos a prueba y comparados con el SPGP sobre varios conjuntos de datos estándar de diferentes propiedades y de gran tamaño. Se intentan identificar además las fortalezas y debilidades de cada uno de los métodos, de manera que sea más sencillo elegir el mejor candidato para cada aplicación potencial.