# Appunti integrati di Machine Learning and Data Analitics (2017/2018)

A cura di

L. B. [1]

---

[1]Lo scopo è quello di dare un ordine gerarchico agli argomenti presenti nel materiale del prof. Eric Medvet

# Indice

## 5 Section 5: Support Vector Machines 70