# Machine Learning and Data Mining Lecture Notes

## CSC C11/D11

Department of Computer and Mathematical Sciences
University of Toronto Scarborough

Version: September 28, 2015

    

# Conventions and Notation

Scalars are written with lower-case italics, e.g., $x$. Column-vectors are written in bold, lower-case: $\mathbf{x}$, and matrices are written in bold uppercase: $\mathbf{B}$.

The set of real numbers is represented by $\mathbb{R}$; $N$-dimensional Euclidean space is written $\mathbb{R}^N$.

> *Aside:*
> Text in "aside" boxes provide extra background or information that you are not required to know for this course.

# Acknowledgements

Graham Taylor and James Martens assisted with preparation of these notes.