# Feature Engineering for Machine Learning

## Principles and Techniques for Data Scientists

*Alice Zheng and Amanda Casari*

# Table of Contents