# Probabilistic models of subjective judgments

by

Zhiwei Li

A dissertation submitted in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
Center for Neural Science
New York University
2021

Advisor:

Todd Gureckis

Zhiwei Li

zhiwei.li@nyu.edu

# TABLE OF CONTENTS

CHAPTER

# LIST OF FIGURES

# LIST OF TABLES

# ABSTRACT

Probabilistic models have been very influential in many cognitive science topics such as language, concepts learning, decision-making etc. The topic of subjective feelings and judgment, however, has been less studied in this computational framework. My thesis focuses on the subjective judgments or emotions that do not obviously relate to instrumental values. My work explores different ways to extract latent variables from the underlying probabilistic model aiming to explain three different subjective judgments: the subjective feeling of suspense, satisfaction of explanation, as well as preference for food.

The first chapter applies a probabilistic model to explain the dynamically changing feeling of suspense preceding the arrival of new information. The central idea of this project is to test and evaluate a model taken from the economic literature (Ely et al., 2015) using a novel empirical paradigm. Succinctly, the model quantifies suspense as the expected belief update in the nearest future, with belief update being quantified as changes in posterior probability. Other heuristics proposed by the larger literature regarding suspense in story telling or movie watching are formalized and compared. Evidence from a variety of stimuli and carefully contrasted conditions indicates that the "future belief update" model best captures the subjective report of suspense.

The second chapter focuses on the feeling of satisfaction of explanation, which is similarly an emotion closely related with the arrival and interpretation of new information. In previous literature, the "simplicity preference" of an explanation has been argued to be a major consideration in how people prefer some explanations over others. I designed a new experimental paradigm that more clearly shows how the prior and causal strength of a causal system can affect people's overall preference for simple or complex explanations. I found that instead of being a universal preference, a simplicity preference for explanation is only present when the prior of each cause existing is low and the causal strength is high. Moreover, a standard Bayesian estimation of the posterior of some explanation being true is not an accurate account of people's preference; rather, people heavily overweight the importance of causal strength than the prior when comparing candidate explanations.

The third and last chapter is a theoretical model regarding how the information collected through active attention relates to how people make value-based decisions. This work is based on the empirical findings from Krajbich et al. that when one snack item has been fixated on more than the other one, it's more likely to be chosen. To explain this effect, a novel model was devel-

oped where the utility of an item is a weighted sum of the posterior mean and the negative posterior standard deviation, with the latter accounting for risk aversion. This model explains the data better than the original attentional drift-diffusion model proposed by Krajbich et al. but worse than a variant with a collapsing bound.

In summary, by constructing different latent variables based upon probabilitic models and testing with new quantitative experiments, my work advances a formal, predictive account of what are previously through to be highly subjective, individualized judgment and emotions.