

Universidade de São Paulo
Instituto de Física

Estudo de estratégias para mudanças coletivas em modelos de opinião.

André Schraider Maizel

Orientador: Prof. Dr. Nestor Felipe Caticha Alfonso

Dissertação de mestrado
apresentada ao Instituto de Física
para a obtenção do título de
Mestre em Ciências

Banca Examinadora:

Prof. Dr. Nestor Felipe Caticha Alfonso (IF-USP)
Prof. Dr. Osame Kinouchi Filho (FFCLRP-USP)
Profa. Dra. Renata Zukanovich Funchal (IF-USP)

São Paulo
2014

-
-
-

FICHA CATALOGRÁFICA
Preparada pelo Serviço de Biblioteca e Informação
do Instituto de Física da Universidade de São Paulo

Maizel, André Schraider

Estudo de estratégias para mudanças coletivas em modelos de opinião. São Paulo, 2014.

Dissertação (Mestrado) – Universidade de São Paulo. Instituto de Física. Depto. Física Geral.

Orientador: Prof. Dr. Nestor Felipe Caticha Alfonso

Área de Concentração: Física

Unitermos: 1. Mecânica estatística; 2. Teoria da informação e comunicação; 3. Sociologia; 4. Moral.

USP/IF/SBI-067/2014

Agradecimentos

Gostaria de agradecer a minha família pelo apoio incondicional dado durante todos os anos de minha formação. Sempre me ajudando quando necessário e me provendo de certezas quando haviam apenas incertezas.

Agradeço também a minha namorada e melhor amiga Elisa, que não só teve a paciência de ouvir sobre meu trabalho durante todos estes anos, como também nunca me deixou baixar a cabeça.

Agradeço a meus amigos por sempre se interessarem pelo meu trabalho, mesmo fugindo tanto de suas áreas de atuação. Agradeço também pelos momentos de descontração, sem estes não seria possível seguir em frente.

Agradeço a meu orientador Nestor Caticha, por ter me dado a oportunidade de trabalhar em uma área tão inusitada, inovadora e empolgante.

Agradeço a Jonatas Eduardo César, por todas as conversas extremamente produtivas.

Agradeço ao curso de Física Biológica, por possibilitar uma formação interdisciplinar, e assim permitir um trânsito entre áreas tão interessantes e distintas.

Agradeço a todos os professores que já passaram por minha formação, pois são a pedra fundamental de qualquer pessoa, seja ela pesquisador ou não.

Agradeço a todos os funcionários do Instituto e todos aqueles que tornaram possível direta ou indiretamente a realização deste trabalho.

Abstract

The study of social systems was always seen as out of scope for the physical sciences. However, in the last years, with the rapid development of statistical mechanics and machine learning, along with recent advances in the field of neuroscience, it became possible to create a wide range of models with the objective to investigate quantitatively aspects of sociology that were mainly considered as qualitative features.

Within the considered problems lies the issue of morality, as well as its consequences to opinion dynamics. More specifically, it is considered relevant to understand how the opinion change dynamics undergoes inside a society, as well as strategies to convince a population to alter its moral direction.

Using an agent based model, in which each agent is represented by a moral vector and has an optimally performing algorithm in the professor/student scenario, we study the influence of two different convincement strategies on the macroscopic behaviour of our model society. In the online learning framework, without any noise, it is known that examples distributed perpendicular to the student achieve an exponential decay in its generalization error. Therefore, we study the effect of this technique as a population convincement strategy, along with its efficiency compared to the standard strategy, in which examples are selected uniformly.

Resumo

O estudo de sistemas sociais sempre foi visto como fora do escopo da física. No entanto, nos últimos anos, com o desenvolvimento da mecânica estatística e da aprendizagem de máquinas, em conjunto com recentes avanços na neurociência, tornou-se possível a criação de diversos modelos no intuito de estudar quantitativamente grandezas antes consideradas majoritariamente qualitativas.

Dentre os problemas considerados está a moralidade, bem como suas consequências para as dinâmicas de opinião. Mais especificamente, considera-se relevante estudar como se dá a mudança de opiniões dentro de uma sociedade, bem como estratégias para convencer uma população a alterar sua direção moral.

Utilizando um modelo baseado em agentes, na qual cada agente é representado por um vetor moral e utiliza uma estratégia de aprendizagem ótima para o cenário professor/aluno, estudamos a influência de duas estratégias de convencimento no comportamento macroscópico de nossa sociedade modelo. Tomando como base a aprendizagem sequencial sem a presença de ruído, e o fato de que seleção de exemplos na borda da dúvida gera um decaimento exponencial do erro de generalização em redes neurais artificiais, estudamos o efeito desta técnica como estratégia de convencimento populacional, assim como a comparação de sua eficácia com a estratégia padrão, na qual os exemplos são selecionados uniformemente.

Sumário

Introdução	1
1 Construção do Modelo	7
1.1 Teoria dos Fundamentos Morais	7
1.2 Bases Neurológicas do Modelo	12
1.3 Aprendizagem de Máquinas	18
2 Modelo e métodos	31
2.1 Modelo	31
2.2 Métodos	35
3 Resultados	39
3.1 Dependência com ρ e β	41
3.2 Dependência com o Número de Vizinhos	48
3.3 Comparação entre as estratégias	50
3.4 Zeitgeist livre	55
Conclusão	59
A Informação e Maximização da Entropia	67

Introdução

“The totality of beliefs and sentiments common to the average members of a society forms a determinate system with a life of its own. It can be termed the collective or creative consciousness.”

Émile Durkheim

O estudo de comportamentos sociais e da moralidade sempre foi de enorme interesse, tanto por parte de cientistas como de políticos e governantes. Tradicionalmente, as ciências sociais possuíam caráter mais qualitativo, com maior proximidade com disciplinas como a filosofia e a linguística. No entanto, nas últimas décadas, a sociologia vem ganhando uma nova abordagem. Utilizando-se de conceitos de estatística e teoria de informação, é possível realizar uma análise quantitativa de características que antes acreditava-se que fossem dificilmente quantificáveis. Com o avanço das neurociências, da inteligência artificial e da inferência estatística, foi possível dar um passo além e estender técnicas comumente utilizadas em física para estudos de psicologia social/moral e de sociologia. Devido ao enorme conjunto de dados presente em ferramentas como a internet e a grande possibilidade de extração de dados de diversas outras fontes, estudos que utilizem técnicas avançadas para a extração de características e padrões se tornam cada vez mais comuns. A crescente geração de informação por meio da internet vem convergindo com estudos teóricos para a criação de modelos de sistemas sociais, possibilitando assim a realização de estudos nas mais diversas frentes.

Alguns conceitos fundamentais da física estatística podem facilmente ser aplicados à sociologia, como por exemplo a ideia de comportamentos emergentes. A sociedade é constituída de um número considerável de elementos, que quando analisados em conjunto apresentam propriedades que só existem como uma média de cada comportamento individual. Ainda em meados do

século XIX Émile Durkheim já afirmava que a moralidade era um agente agregador, e que os grupos poderiam ser vistos como propriedades emergentes com características próprias [6, 13]. Suas ideias serviram de inspiração para a nova geração de sociólogos e psicólogos, como Jonathan Haidt. Com a descoberta de reações sociais intuitivas e de grande caráter evolutivo, a hipótese de emergência ganhou força. Dado que a socialização e formação de grupos deve originalmente ser uma consequência evolutiva, é de se esperar que suas propriedades organizacionais e de coalescência surjam a partir de uma coleção de comportamentos individuais, e que dada a natureza das interações, não podem ser observadas individualmente. Desta maneira, se torna extremamente sugestiva e apropriada a utilização de técnicas de mecânica estatística, na qual propriedades emergentes são extremamente comuns, como por exemplo no magnetismo. Como veremos, a sociedade pode ser vista de certa forma como um sistema magnético, possibilitando assim paralelos entre as duas teorias, a princípio tão distintas.

A teoria da informação nos fornece a conexão direta com a estatística, abrindo a possibilidade de interpretar os dados obtidos por um novo ponto de vista. Utilizando de princípios como a máxima entropia, podemos obter indicações de que tipo de tratamento devemos fornecer a nossos dados, nos levando diretamente ao ferramental da mecânica estatística de Boltzmann. Utilizando de técnicas mais complexas de análise de dados é possível também extrair características de dados multidimensionais, classificando-os e agrupando-os por diferentes medidas de similaridade.

Trabalhos realizados na área se focaram principalmente em características fundamentais, como a estrutura de rede, ou o tipo de interação entre os indivíduos. Trabalhos como o de Barabási e Albert por exemplo, visam a compreensão de como propriedades intrínsecas da rede podem influenciar seus grafos associados, como por exemplo a conectividade preferencial das redes livre de escala. Este tipo de propriedade pode por exemplo, ser relacionada à concentração de renda e/ou de capital social (como por exemplo o número de eleitores). Outros trabalhos visam entender como as interações neurológicas e algumas características evolutivas determinam a essência de nosso comportamento coletivo. No entanto, para que se obtenha um panorama satisfatório de um problema tão complexo como o comportamento social, é necessário utilizar de todas as ferramentas e abordagens possíveis. Em diversos estudos anteriores de Nestor Caticha e colaboradores, foram abordados os problemas do surgimento de hierarquias, do aparecimento de moedas de troca, e das relações entre características cognitivas e morais, entre outros. Pretendendo portanto atacar o problema por diversas frentes, a fim de obter uma visão mais ampla do comportamento humano.

Um dos problemas de nosso interesse está relacionado com a estrutura

da moralidade e de sua dinâmica como fator social. Nos últimos anos a psicologia moral vem passando por uma mudança de paradigma, no qual cada vez mais acredita-se que o julgamento moral seja primariamente intuitivo. De tal forma, cria-se um grande interesse para a compreensão dos mecanismos envolvidos que regem tal julgamento, bem como de suas consequências na sociedade e no comportamento coletivo. Utilizando resultados empíricos de psicologia social e de cognição moral, podemos avaliar que tipo de características um modelo coerente de moralidade deve conter. Outro fator importante é a compreensão de que tipo de características do modelo podem ser efetivamente comparadas com estes resultados empíricos, de forma a possibilitar uma validação posterior dos resultados e conclusões obtidas.

Dentro do espectro de problemas associados a moralidade podemos encontrar a dinâmica de opiniões. Baseadas principalmente em julgamentos morais, as sociedades tendem a definir que tipos de comportamentos e opiniões são consideradas corretas e incorretas, criando assim uma espécie de direção moral preferencial. Entender como se dá o processo de organização das opiniões pode fornecer indícios de como amenizar conflitos políticos e/ou ideológicos. Consideremos por exemplo as diferenças entre pessoas de diversas religiões, que muitas vezes culminam em conflitos (vide o caso Israel-Palestina) motivados principalmente à uma oposição total de ideias. É portanto de grande interesse a análise de estratégias para o convencimento de sociedades de maneira eficiente, ou ao menos factível. Sabendo que na maioria dos casos a diferença total de opiniões tende a gerar uma grande intransigência, devemos procurar uma forma de superar esta barreira psicológica e informacional.

Em um de seus últimos trabalhos Nestor Caticha, Renato Vicente e colaboradores desenvolveram um modelo baseado no trabalho do psicólogo social Jonathan Haidt e nos fundamentos de inteligência artificial, com suporte em experimentos de cognição moral [1, 2]. Consistindo em uma sociedade de agentes conformistas, que navegam pelo espaço da moralidade utilizando um função ótima de aprendizado sequencial. Partindo de um modelo com poucos elementos, foi possível entender o surgimento de uma transição de fase por meio de aproximações de campo médio, que nos forneceram indícios de quais parâmetros de ordem observar, nos levando a correlações entre atitude política e estilo cognitivo, posteriormente validado por Jonatas César em sua tese de doutoramento [34] por meio de dados empíricos. Este modelo também possibilitou a compreensão de que a pressão social exercida sobre a sociedade tende a alterar sua dispersão de opiniões. Utilizando dados obtidos por Jonathan Haidt em um questionário feito pela internet, Nestor Caticha e Jonatas César demonstraram que assinaturas estatísticas obtidas por meio de simulações podem ser associadas aos resultados reais, extraídos do trabalho

de Haidt. Tal comparação permitiu a associação de agentes com um perfil cognitivo corroborativo com pessoas conservadoras, e agentes com perfil cognitivo mais explorador (*novelty seekers*) puderam ser associados a liberais. Desta forma, foi possível fazer um paralelo entre os resultados e parâmetros do modelo com características sociais existentes na realidade, possibilitando assim uma interpretação das características observadas.

Partindo deste mesmo modelo, propomos um estudo de estratégias para alterar a opinião média da sociedade, utilizando um caso mais geral no qual é possível a existência de uma influência externa na sociedade. Para este estudo, consideraremos a existência de um *oráculo*, que apresentará exemplos à sociedade, na tentativa de alterar sua direção preferencial. Serão analisadas duas situações distintas, a primeira na qual a sociedade possui uma direção fixa, competindo então com o oráculo. A segunda, na qual a direção preferencial é tomada por meio da média das opiniões da sociedade, possibilitando a análise de como esta mudança ocorre. Utilizando características extraídas de algoritmos de aprendizagem de máquina, temos uma indicação de quais estratégias de convencimento utilizar, bem como de qual estratégia deve performar melhor, assim como os resultados anteriores obtidos por Caticha e César nos indicam quais comportamentos do modelo devem ser mantidos para que ainda haja validade dos resultados neste caso mais geral.

Diversos modelos em sociofísica se propuseram a estudar a dinâmica de opiniões, bem como a evolução social e a influência de uma mídia/agente externo no sistema. Um exemplo clássico é o modelo de Axelrod [3] de 1994, que visava o estudo de polarização e formação de culturas distintas em um modelo baseado em agentes. Posteriormente, o modelo de Axelrod foi estendido para casos mais gerais [4, 5], com a inclusão de um campo externo, na intenção de simular o efeito da mídia e de comunicação em massa, bem como estudos sobre a influência da dimensionalidade no modelo de Axelrod. Nosso modelo difere em diversas características, começando por sua formulação, que leva em conta as características neurológicas para a determinação da forma de interação entre os agentes. Outra diferença fundamental se dá no fato que não há formação de domínios e facções, dado que o objetivo principal de nosso trabalho é estudar a dinâmica de mudança de opiniões em uma sociedade homogênea, bem como a dimensionalidade é um fator fundamental de nosso modelo, derivado da natureza da moral. No entanto, existem diversas semelhanças entre nosso modelo e diversos outros já estudados, como por exemplos a homofilia, representada em nosso caso por uma homogeneidade de estilos cognitivos, bem como o caráter orientacional do modelo, permitindo a análise de similaridades geométricas entre os agentes (como por exemplo as medidas de magnetização).

Esta dissertação se encontra separada em três etapas. Primeiramente

serão apresentados os ingredientes empíricos e teóricos que dão origem ao modelo. Após a apresentação de todos os elementos necessários para a construção do trabalho, o modelo utilizado será rapidamente desenvolvido, bem como uma explicação de quais métodos serão utilizados para a obtenção das grandezas de interesse. Por fim, será feita uma análise dos resultados, seguida de uma conclusão e de perspectivas futuras de trabalho.

Capítulo 1

Construção do Modelo

1.1 Teoria dos Fundamentos Morais

Tratar sobre a moralidade humana é sempre uma tarefa complicada, tanto devido à enorme complexidade do problema quanto ao fato de estarmos embebidos em um enorme viés cultural e histórico. Durante os últimos séculos o assunto foi estudado por filósofos, até então com um caráter metafísico e normativo. Para Platão a moral podia ser vista como uma virtude, uma objetivação da racionalidade humana, onde os sentimentos e sensações eram considerados como um obstáculo na obtenção de um comportamento ético. Para Aristóteles (que pode ser visto de certa maneira como seu sucessor) a ética e a moral também eram frutos da racionalidade humana, tida como a maneira de definir regras pelas quais se poderia viver da melhor maneira possível. Séculos mais tarde, Immanuel Kant considerou que o papel da moral era estabelecer regras para um convívio social ético, a partir do *imperativo categórico*, um conjunto mínimo de regras que afirmou serem suficientes para garantir o pleno exercício da moral. Podemos notar que a moralidade era tradicionalmente considerada como uma virtude, reflexo da racionalidade humana. Porém nos últimos anos, com o avanço de áreas como psicologia social e neurociências, foi possível mudar radicalmente o paradigma do problema associado à moralidade.

Atualmente, o psicólogo Jonathan Haidt propõe que a moral humana seja primariamente um conjunto de reações intuitivas, sendo apenas posteriormente avaliadas com argumentos racionais. Esta teoria foi primeiramente proposta por Jonathan Haidt e Craig Joseph em um artigo de 2004 [8] e estudada mais a fundo pelos mesmos autores em um livro de 2006 [11]. Em um experimento de 2005 [9], Haidt apresentou situações com diversos dilemas morais para um grupo de participantes, onde um percentual deste grupo

havia sido submetido à um tratamento de hipnose no qual uma palavra seria utilizada como gatilho para uma sensação de nojo. Haidt notou que os voluntários que haviam sido hipnotizados possuíam uma resposta mais negativa associada às situações moralmente controversas (como incesto, trapaça, etc), assim como um alto índice de reprovação moral em situações na qual não havia nenhum dilema presente. Em outro experimento [10], Haidt utilizou situações na qual havia um dilema moral presente, porém sem explicação racional clara para a rejeição de uma certa situação, como por exemplo, um caso de incesto consensual no qual ambos estavam protegidos contra uma possível gravidez e ambos concordariam em nunca mais repetir a experiência. Ao perguntar para os participantes da pesquisa se estes eram contra ou à favor, a maioria das pessoas se declaravam extremamente contrárias, porém apresentavam grande dificuldade em elaborar uma explicação racional para sustentar sua opinião, indicando que o processo de racionalização da moralidade é posterior ao julgamento moral intuitivo. Esta característica será parte essencial do modelo.

Outra característica fundamental estudada por Haidt envolve a estrutura da moralidade. Para Haidt a moral humana é composta por um conjunto limitado de dimensões inatas com uma grande carga evolutiva associada (i.e características adquiridas ao longo de nossa evolução) [11, 12]. Primeiramente, os pesquisadores de psicologia moral acreditavam que a fundação moral seria constituída apenas de *justiça/trapaça* e *cuidado/violência*. Porém, ao longo de diversos estudos transculturais, Haidt constatou que havia grande responsabilidade de um viés cultural nas conclusões obtidas pelos pesquisadores, quase sempre liberais provenientes de uma cultura ocidental, causando assim um estreitamento inapropriado da moralidade à estas duas dimensões. Ao analisar textos morais como a *Bíblia*, o *Corão* e o *Código de Hammurabi* ficou evidente que ao menos mais três dimensões deveriam ser adicionadas, sendo elas: *lealdade ao grupo/traição*, *respeito à autoridade/subversão* e *pureza/degradação*. Abaixo estão listadas as cinco dimensões inatas da moralidade¹ e suas prováveis origens biológicas, assim como suas consequências no convívio social.

- **Justiça/Trapaça:** Consequência evolutiva da necessidade de se eliminar aproveitadores (*free riders*) para que atos de cooperação possam

¹Os nomes das dimensões morais são uma tradução livre de:

Fairness/Cheating;

Care/Harm;

(in-group) Loyalty/Betrayal;

Authority/Subversion;

Santity (purity)/Degradation.

elevar a aptidão (*fitness*) média da sociedade. Está presente em todas as culturas, sendo um grande fator de agregação para altruísmo mútuo e cooperação entre indivíduos sem parentesco. Desrespeitar os conceitos de justiça usualmente desencadeiam sentimentos de raiva e culpa, enquanto que o cumprimento das normas sociais relacionadas costumam gerar sentimentos de gratidão.

- **Cuidado/Violência:** Surge a partir da maturidade tardia dos filhotes em primatas e que implicam em uma necessidade de cuidado para com indivíduos frágeis (crianças, idosos, doentes) a fim de garantir a sobrevivência/perpetuação da espécie. Em muitas espécies de primatas, em particular o ser humano, a sensibilidade quanto aos maus-tratos e sentimentos de compaixão se estendem além da própria prole. A aprovação de indivíduos que previnem ações violentas estão relacionadas com bondade e compaixão, enquanto que indivíduos violentos são normalmente associados a vícios como crueldade e agressão.
- **Lealdade ao grupo/Traição:** Ferramenta evolutiva essencial para a manutenção de coalizões e da estrutura de grupos sociais. Fornece uma intuição de que indivíduos podem ser considerados "traidores", gerando uma sensação de repulsa, revolta, ou até mesmo ódio com relação aos indivíduos desleais ou que colocam o grupo em risco. Está também relacionado com o sentimento de desconfiança com indivíduos pertencentes a outros grupos. Diversas espécies de primatas apresentam formas de ostracismo e abandono à animais que ferem o princípio lealdade ao grupo.
- **Respeito à autoridade/Subversão:** Consequência evolutiva do surgimento de hierarquias em sociedades/comunidades. Inicialmente as estruturas hierárquicas surgiram por meio de dominância à partir de força física, onde os machos/fêmeas dominantes devia oferecer proteção ao grupo em troca do privilégio na escolha de parceiros sexuais. Nos seres humanos e em primatas com capacidade cognitiva mais alta esta hierarquização é mais sutil, estando relacionada com prestígio e deferência voluntária. As pessoas normalmente sentem respeito e admiração pelos líderes, onde boa capacidade de liderança é visto também como uma virtude. Diversos mamíferos apresentam estruturas hierarquizadas, a falta de sensibilidade quanto a estrutura de liderança pode levar à resultados desastrosos para o indivíduo, como a morte ou a impossibilidade de acasalamento.
- **Santidade (ou Pureza)/Degradação:** Ferramenta evolutiva para

minimizar as chances de que infecções/doenças causadas por ingestão de alimentos ou substâncias tóxicas se espalhem na comunidade. Serve também como ferramenta de proteção contra objetos/situações potencialmente perigosos(as). Indivíduos que não se enquadrarem no padrão alimentar/comportamental são isolados ou punidos. Está relacionado com o surgimento de religiões e regras comportamentais estritas. Está altamente relacionado com o sentimento de nojo, podendo ser ativado não apenas por hábitos alimentares como também por padrões estéticos (de vestimentas e até mesmo biológicos).

Outro ponto extremamente importante da Teoria dos Fundamentos Morais é a forma com a qual as pessoas atribuem valores a cada dimensão moral. Em um estudo experimental [13], Haidt aplicou dois questionários que visavam avaliar a importância dada a cada uma das componentes morais. No primeiro questionário (Figura 1.1), o participante devia fornecer um valor monetário para realizar uma certa atitude moralmente questionável, onde cada ação representaria uma categoria moral.

How much money would it take to get you to...			
	Column A	Column B	Moral category
1)	Stick a pin into your palm. \$ ____	Stick a pin into the palm of a child you don't know. \$ ____	Harm/care
2)	Accept a plasma screen television that a friend of yours wants to give you. You know that your friend got the television a year ago when the company that made it sent it, by mistake and at no charge, to your friend. \$ ____	Accept a plasma screen television that a friend of yours wants to give you. You know that your friend bought the TV a year ago from a thief who had stolen it from a wealthy family. \$ ____	Fairness/reciprocity
3)	Say something slightly bad about your nation (which you don't believe to be true) while calling in, anonymously, to a talk-radio show in your nation. \$ ____	Say something slightly bad about your nation (which you don't believe to be true) while calling in, anonymously, to a talk-radio show in a foreign nation. \$ ____	Ingroup/loyalty
4)	Slap a friend in the face (with his/her permission) as part of a comedy skit. \$ ____	Slap your father in the face (with his permission) as part of a comedy skit. \$ ____	Authority/respect
5)	Attend a performance art piece in which the actors act like idiots for 30 min, including failing to solve simple problems and falling down repeatedly on stage. \$ ____	Attend a performance art piece in which the actors act like animals for 30 min, including crawling around naked and urinating on stage. \$ ____	Purity/sanctity
	Total for column A: \$ ____	Total for column B: \$ ____	

Figura 1.1: Tabela utilizada em [13] para avaliar a utilização das componentes morais.

No segundo questionário, o participante deveria declarar sua afiliação política, em uma escala de *Muito Liberal* para *Muito Conservador*, e então

responder uma série de 15 questões, onde se deveria atribuir um valor de 1-6 para o grau de relevância da questão apresentada (sendo de nunca relevante à sempre relevante, respectivamente) visando obter uma relação entre a afiliação política e as dimensões morais. Os resultados obtidos (Figura 1.2) mostraram que pessoas auto declaradas como liberais tendem a atribuir um valor maior às dimensões de *cuidado/violência* e *justiça/trapaça*, enquanto que os conservadores atribuem em média a mesma importância à todas as 5 dimensões morais. Este questionário ainda se encontra ativo na internet no endereço www.yourmorals.org, podendo ser respondido por qualquer pessoa.

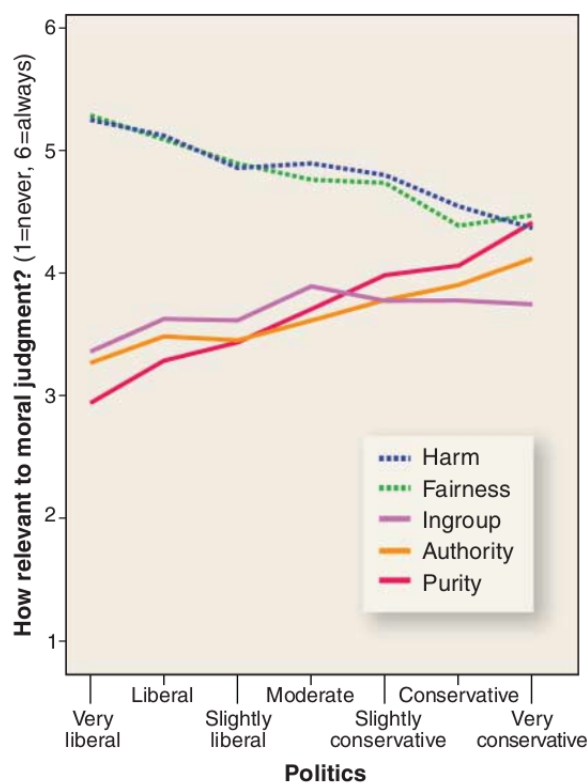


Figura 1.2: Resultado obtido em [13] para a relação entre afiliação política e a importância relativa de cada fundação moral.

Após estes resultados fica claro o viés ao qual as pesquisas na área de psicologia moral estavam sujeitas. Por se tratarem de liberais, os pesquisadores possuíam a falsa impressão de que a moralidade é baseada apenas nas dimensões que os liberais valorizam. Isto acaba por evidenciar também a importância de estudos com um caráter mais objetivo como o realizado por Haidt. Tendo a Teoria dos Fundamentos Morais como ponto de partida

para o desenvolvimento do modelo podemos então avançar para as bases neurológicas de nossa modelagem.

1.2 Bases Neurológicas do Modelo

Tomando como princípio o fato de que o comportamento humano é necessariamente regido por seu sistema nervoso, entender os mecanismos do cérebro associados as características fundamentais do modelo é de extrema importância. Os dois aspectos necessários para a execução deste trabalho são o aprendizado por reforço e a tendência de socialização.

Pressão social, exclusão social e conformidade

É fato incontestável a nossa enorme capacidade de socialização, assim como nossa necessidade por aprovação/aceitação do grupo ao qual fazemos parte. Esta urgência por nos sentirmos enquadrados em um coletivo é tão forte que pode até mesmo mudar nossa visão sobre fatos e verdades, como ficará evidente mais a frente nesta seção. Porém, qual será o fator biológico que nos move e nos faz alterar nossa percepção de forma tão intensa? A exclusão social sempre foi associada com o sentimento de dor, porém de maneira metafórica. Entretanto, nos últimos anos, com o auxílio de estudos de ressonância magnética funcional (fMRI), neurocientistas foram capazes de desvendar as estruturas do cérebro associadas com a sensação de rejeição. Em um estudo de 2003 [15] Naomi Eisenberger, Matthew Lieberman e Kipling Williams submeteram voluntários a um experimento controlado, no qual suas reações a situações de socialização e rejeição eram medidas por fMRI. Desta maneira, os autores puderam avaliar quais regiões do cérebro estão envolvidas no processamento de informações sociais relacionadas a aceitação. O experimento era composto de três situações distintas. A primeira situação consistia em colocar os voluntários no scanner de ressonância, apresentar-lhes um vídeo com um jogo de *CyberBall*² e contar uma falsa história, afirmando que a conexão com aquele scanner ainda não havia sido finalizada e portanto não seria possível jogar a partida, apenas assisti-la; criando desta forma uma sensação de exclusão. A segunda situação era semelhante, porém nesta os participantes poderiam jogar, gerando assim o sentimento de inclusão. E no último teste, os participantes poderiam jogar o jogo, porém após sete arremessos estes seriam impedidos de jogar e poderiam apenas observar o resto do jogo, possibilitando assim uma análise relativa entre exclusão/inclusão. Terminado o experimento, os voluntários deveriam responder um questionário

²Jogo para o console NES[©], de 1988.

inferindo sobre o quão excluídos e o nível de angústia que estes experimentaram ao longo do exame. Os resultados obtidos pelos pesquisadores foram surpreendentes, a região do cérebro ativada pela sensação de rejeição social era exatamente a mesma região ativada pela dor física (Figura 1.3), o córtex anterior cingulado (ACC), elevando assim a expressão *a dor da exclusão* de metáfora à realidade.

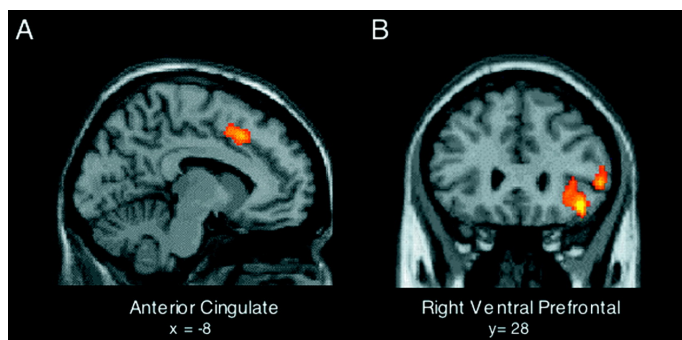


Figura 1.3: Imagem de fMRI obtida no experimento realizado em [15], exibindo as regiões ativas no cérebro durante o experimento de exclusão

Sabemos portanto, que a exclusão social leva à dor, e como um reflexo natural; os seres humanos tenderão a concordar a fim de minimizar esta sensação (o que mais a frente nos possibilitará utilizar uma função custo/energia). Esta tendência foi evidenciada em dois estudos clássicos de psicologia social.

No primeiro experimento, realizado em 1937 por Muzafer Sherif [17], participantes eram colocados em uma sala escura em duas situações distintas: sozinhos e em grupos. Dentro da sala os participantes deveriam dizer se um ponto de luz projetado em uma das paredes era estático ou se estava se movimentando, sem saber que o ponto estava sempre fixo. Quando colocados em grupos, os participantes possuíam uma tendência a entrar consenso, independente da validade da resposta.

No segundo experimento, realizado em 1951 por Solomon Asch, os participantes deveriam comparar o tamanho de uma linha vertical de referência com o de três linhas verticais de tamanhos diferentes, indicando qual destas linhas possuía o mesmo comprimento da referência (Figura 1.4). No entanto, esta escolha deveria ser feita em um grupo contendo pessoas ditas *confederadas*, ou seja, indivíduos que sabiam os objetivos do estudo. Estes *confederados* deveriam indicar uma resposta claramente errada, influenciando a escolha do indivíduo estudado, que na maioria dos casos acabava por optar pela mesma resposta do grupo. No grupo de controle (composto apenas por um participante “real”) a taxa de erros era inferior a 1% das tentativas, enquanto que

nos grupos contendo atores a taxa de erro chegou a incríveis 33%, onde 75% dos participantes cometeriam o erro pelo menos uma vez.

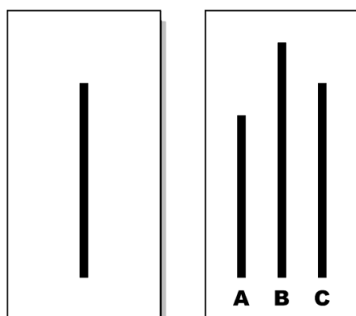


Figura 1.4: Exemplificação das cartas utilizadas no experimento de Asch. Figura retirada de [18].

O experimento de Asch costuma ser interpretado como uma forte evidência para a conformidade e para a *influência social normativa* [20–22], onde a influência normativa é tida como a disposição a se conformar publicamente e obter recompensas sociais a fim de evitar punições e penalidades sociais [23]. É possível assumir que a pressão social possa influenciar não apenas padrões comportamentais simples, mas também fenômenos complexos tais como a atitude política de uma sociedade. Em um estudo de 2009, Nail e McGregor [16] compararam o resultado de pesquisas de opinião que visavam avaliar o grau de conservadorismo da sociedade estadunidense, realizadas em duas épocas distintas: pré 11 de Setembro e pós 11 de Setembro. É de conhecimento geral que após os atentados o nível de tolerância social no Estados Unidos diminuiu drasticamente. Como consequência, podemos assumir que a pressão social se encontrava em níveis bem mais elevados. O estudo constatou que, dentre os respondentes, todo o espectro de afiliações políticas (liberal-conservador) sofria uma deslocamento em direção à uma atitude política mais conservadora, como é possível observar na Figura 1.5. É portanto de se esperar que nosso modelo seja capaz de reproduzir esta característica, ou ao menos que esta característica possa emergir de maneira natural de nossos resultados.

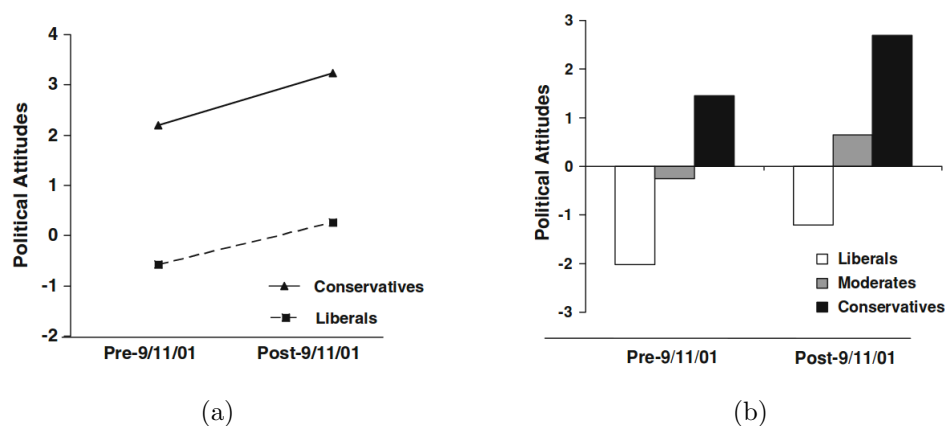


Figura 1.5: Gráficos retirados de [16] comparando as atitudes políticas médias em função da data. Maior atitude política corresponde à um alto nível de conservadorismo. (a) Comparação considerando apenas duas categorias: conservadores e liberais. (b) Comparação levando em conta três classificações: liberais, moderados e conservadores.

Sabendo qual o comportamento humano quanto a pressão social e nossa tendência à conformidade, partiremos agora para um dos mecanismos de aprendizado que permite-nos navegar pelo universo da moral durante nossa vida.

Aprendizado por reforço

Existem muitas maneiras pelas quais um ser humano aprende em sua vida, todas elas podem ser divididas em apenas duas categorias: *aprendizado por associação* e *aprendizado não associativo*.

O aprendizado não associativo está relacionado com uma mudança gradual e permanente de resposta com relação a um estímulo devido a exposição prolongada ao mesmo, e está mais relacionada com respostas químicas como a causada pelo uso de drogas (como a de sensibilização e ou habituação a uma substância).

Já o aprendizado associativo está relacionado com a associação de dois estímulos, ou de um comportamento/resposta e um estímulo. Podendo então ser atribuído boa parte das tarefas de aprendizado cujas quais estamos sujeitos. Considerando o caráter intuitivo de nosso modelo de moral e a grande carga associativa que o aprendizado social apresenta, devemos entender qual os diferentes tipos de estilos cognitivos e saber qual a influência que estes apresentarão em nosso estudo. O mecanismo cerebral envolvido na fixação

de uma informação aprendida está fortemente relacionado com o reforço das conexões sinápticas envolvendo os diferentes estímulos/respostas. Diferentes pessoas terão diferentes estilos cognitivos, ou seja, reagirão de maneiras distintas com relação a fixação/importância dada a cada informação recebida.

Cientistas sociais vêm percebendo que existem claras diferenças cognitivas e motivacionais entre pessoas consideradas *liberais* e pessoas consideradas *conservadoras*³. Estas diferenças cognitivas relacionadas à afiliação política já mostraram grandes evidências quanto a sua origem em caráter genético, com uma grande influência da infância e também relativamente estável durante todo o período de vida [24, 25]. Os pesquisadores envolvidos com o comportamento humano sugerem que as diferenças psico-cognitivas entre liberais e conservadores sejam devidas ao sistema auto regulatório de monitoramento de conflitos [26]. O mecanismo de monitoramento de conflitos consiste em identificar erros em respostas habituais dado a incompatibilidade destas com a situação atual. Este tipo de controle está associado a atividades neurológicas no córtex anterior cingulado (ACC) [27], estrutura cerebral responsável por diversas funções cognitivas, tais como tomada de decisão e antecipação de recompensa (como no célebre experimento de Pavlov), e que curiosamente é também a região cerebral responsável pelo registro da dor (vide os resultados apresentados na seção anterior). Estudos comportamentais mostraram que conservadores costumam possuir um julgamento mais estruturado e persistente em problemas de tomada de decisão, enquanto que liberais se mostraram mais tolerantes à ambigüidade e mais abertos a novas experiências/informações em termos de medidas psicológicas. [28]

Em um experimento publicado em 2007, David Amodio *et al.* se propuseram a estudar a relação entre as diferenças nas respostas do córtex anterior cingulado de liberais e de conservadores, na tentativa de confirmar as diferenças cognitivas estabelecidas por Jost em 2003. Para atingir tal objetivo, Amodio e seus colaboradores submeteram seus voluntários a um estudo do tipo *Go/No-Go*. Os estudos *Go/No-Go* consistem na apresentação de uma sequência de formas geométricas simples (quadrados, círculos etc.), onde um mesmo estímulo é apresentado diversas vezes (estímulo *Go*) até que o participante se torne habituado. Entretanto, em uma pequena parcela das tentativas serão apresentados estímulos inesperados (*No-Go*), de forma a romper com a expectativa gerada pela sequência anterior, possibilitando assim uma medida da atividade relativa das regiões cerebrais envolvidas. Estas medidas estão relacionadas com uma maior atividade no ACC e são obtidas utilizando-se uma técnica de análise eletroencefalográfica conhecida como *Potenciais de*

³de acordo com os padrões estadunidenses de liberal/conservador

*Relacionados a Eventos*⁴ (ERP). A medida relativa da atividade no ACC durante os estímulos Go e os estímulos No-Go é a chamada *Negatividade Relacionada ao Erro*⁵ (ERN), que é a diferença entre os picos de atividades cerebrais detectados após os estímulos (ambos Go e No-Go). Para estabelecer a correlação entre a afiliação política e os resultados obtidos no ERP os participantes deveriam declarar sua afiliação política em uma escala de -5 (muito liberal) até +5 (muito conservador). As amplitudes de ERN foram então comparadas às afiliações políticas (Figura 1.5a). O conjunto de dados apresentou uma correlação de 59% entre a ERN e a escala de afiliações com um valor P inferior a 0,001.

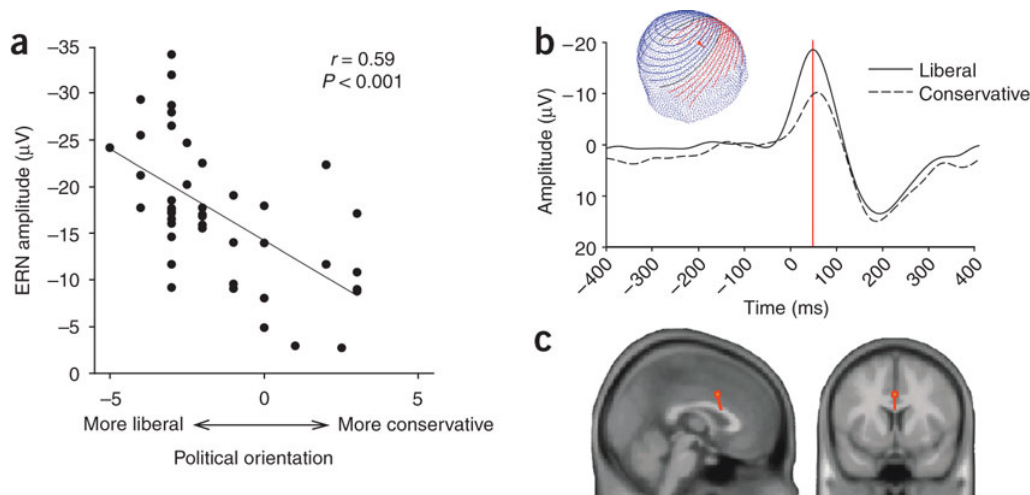


Figura 1.6: (a) Gráfico das Amplitudes de ERN contra a afiliação política auto-declarada. (b) Representação esquemática das formas de ondas ERP correspondendo aos estímulos No-Go com as ondas correspondentes aos estímulos Go já devidamente subtraídas. (c) Indicação da região do Córtex Anterior Cingulado (ACC) como fonte da atividade detectada.

Os resultados obtidos por Amodio e seus colaboradores indicam claramente que há uma diferença cognitiva entre liberais e conservadores. Esta distinção será essencial em nosso modelo, como será mostrado mais à frente nesta dissertação.

Outro fator importante relacionado ao aprendizado por reforço está relacionado com a capacidade de mudança de opiniões. Em um estudo de 1973 Charles Blake Keasey testou a capacidade de formação e mudança de opinião em um grupo de jovens [19], onde foram notadas diferenças signifi-

⁴Tradução livre de *Event-Related Potential*

⁵Tradução livre de *Error-Related Negativity*

tivas nos processos de racionalização e mudança de opinião. Utilizando como inspiração os experimentos de Asch e Scheriff (apresentados na seção 2.2), Keasey apresentou jovens em diferentes estágios do desenvolvimento moral a uma variedade de dilemas, onde um grupo de voluntários deveria fornecer a sua opinião logo após a apresentação dos dilemas, e o outro grupo deveria apresentar sua opinião apenas após a argumentação de dois atores contratados. Após duas semanas, o experimento foi refeito, de forma a examinar possíveis mudanças de opinião, como resultado, foi observado que pessoas com maior interação social foram sujeitas a uma maior variação de opinião. Keasey pode inferir que um maior número de interações sociais está correlacionado com uma maior habilidade de mudar sua opinião, pois ao ser apresentado a diferentes pontos de vista o sujeito é capaz de alterar sua percepção a partir dos mecanismos de aprendizagem por reforço. Como veremos na próxima seção, este resultado possui grande compatibilidade com nosso modelo e com os ferramentais teóricos da aprendizagem de máquinas, onde veremos que agentes apresentados a um maior número de exemplos possuem uma maior facilidade em sua mudança de opiniões, bem como uma modulação do aprendizado que valoriza a novidade.

Estamos agora em posse de todos os ingredientes empíricos para a formulação de nosso modelo. Podemos portanto nos focar agora nos fundamentos matemáticos necessários para o desenvolvimento da modelagem. Necessitaremos primariamente do arcabouço lógico da área de Aprendizagem de Máquina (*Machine Learning* ou ML). Para este trabalho, estaremos interessados apenas no caso especial chamado aprendizado *online* ou *sequencial*. Ambos os tópicos serão discutidos na seção a seguir.

1.3 Aprendizagem de Máquinas

A compreensão da origem de comportamentos inteligentes é um tópico que ocupa há muitos séculos a mente humana. Porém, por muito tempo se acreditou que o cérebro não possuía função biológica e que toda a consciência estaria no coração. Esta ideia foi contestada primeiramente por Hipócrates, afirmando de forma pioneira que o cérebro era a fonte das sensações e a casa da razão. Já no século XVIII, Luigi Galvani descobriu a origem elétrica dos movimentos motores, utilizando correntes elétricas para mover músculos de animais mortos. A partir da criação do microscópio as neurociências experienciaram um enorme salto, passando pelo descobrimento das estruturas neuronais por Santiago Ramón y Cajal e posteriormente por Camilo Golgi. Ficou estabelecido então que os responsáveis pelo processamento de informação e tomada de decisão eram o sistema nervoso e seus componentes.

Em 1943 Warren McCulloch e Walter Pitts desenvolveram o primeiro modelo formal de um neurônio [29], conhecido por *threshold logic unit*. O modelo de *McCulloch-Pitts* consiste em uma função bi-estável que pode ser considerada *ativa* ou *passiva*. Sendo portanto representado por uma variável binária $S = \pm 1$, onde $+1$ equivale a um estado ativo e -1 a um estado passivo. O estado de cada neurônio S_i será definido ao longo do tempo por sinais recebidos por outros neurônios ou por uma entrada externa S_j , ponderados por um acoplamento sináptico J_{ij} . A média ponderada das entradas S_j é chamada de *Potencial Pós-Sináptico*. Podemos então representar o estado do neurônio S_i em um dado instante por:

$$S_i(t+1) = \text{sgn} \left(\sum_j J_{ij} S_j(t) - \theta_i \right) \quad (1.1)$$

Onde θ_i é um limiar de ativação e $\text{sgn}(x)$ é a função sinal, definida por $\text{sgn}(x) = +1$ se $x > 0$ e $\text{sgn}(x) = -1$ se $x < 0$. Entretanto, o comportamento de um único neurônio está longe de sistemas biológicos, que apresentam milhões de neurônios. Devemos então considerar a conexão de um ou mais neurônios, em uma chamada *rede neural artificial*. Podemos notar que devido a natureza do neurônio, o estado da rede é altamente dependente de sua arquitetura, ou seja, da estrutura do grafo representando a conectividade entre seus elementos. Para arquiteturas mais simples, onde apenas um neurônio se conecta a todos os outros, é possível fazer uma análise matemática mais detalhada. Porém, a grande maioria dos casos interessantes e realísticos não podem ser tratados de maneira analítica. Podemos no entanto nos basear em técnicas de mecânica estatística para a compreensão das grandezas e parâmetros de ordem interessantes para a caracterização das redes. De tal forma que podemos considerar a hipótese central da mecânica estatística do aprendizado é que a realização da aprendizagem é uma propriedade emergente, e portanto pode ser estudada por meio de grandes redes de neurônios do tipo McCulloch e Pitts [30].

No caso de arquiteturas do tipo alimentação direta (*feed-forward*), a rede consiste em camadas $l = 1, \dots, L$ de neurônios, na qual as conexões sinápticas dos neurônios da camada l se dão sempre com as camadas posteriores ($l+1$). A primeira camada é dita *entrada* ou *input* e a última camada é dita *saída* ou *output*. Devido a simplicidade das estruturas do tipo *feed-forward* é possível mapear a relação entre *input* e *output* utilizando a dinâmica (1.1).

O caso geral a ser considerado é o de uma rede do tipo *feed-forward* com N neurônios de entrada denotados por S_i com $i = 1, \dots, N$ e apenas um neurônio de saída representado por $\sigma = \text{sgn}(\vec{J} \cdot \vec{S})$. As conexões sinápticas entre o *output* e o *input* são ponderadas por um vetor $\vec{J} = \{J_1, \dots, J_N\}$. O

valor de cada componente do vetor J deve ser ajustado ao longo da dinâmica de aprendizado, de forma a atingir uma relação input-ouput alvo. Este ajuste dos pesos sinápticos é conhecido por *regra*. É extremamente conveniente considerar a relação alvo como uma rede *professor* $\vec{T} = \{T_1, \dots, T_N\}$, nas quais as componentes T_i são fixas. As únicas informações acessíveis ao vetor estudante \vec{J} são as entradas e suas saídas correspondentes dadas pelo professor, contidas em um conjunto de treinamento composto de p exemplos, onde cada *exemplo* correspondente a um par $(\vec{\xi}^\mu, \sigma_T^\mu)$, no qual $\vec{\xi} = \xi_1^\mu, \dots, \xi_N^\mu$ e $\sigma_T^\mu = \pm 1$. O índice μ representa a ordem dos exemplos a serem apresentados, sendo igual a $\mu = 1, \dots, p$. Os exemplos devem ser sorteados aleatoriamente de uma distribuição de probabilidades $P(\mathbf{S})$, onde a escolha desta distribuição pode ser essencial para a rapidez do aprendizado, como será mostrado mais à frente. Devido a natureza booleana de nossa saída, podemos afirmar que a rede neural proposta é um *classificador*, ou seja, é capaz de classificar um conjunto de pontos *linearmente separável* em duas categorias distintas.

Existem dois casos possíveis para o problema geral apresentado anteriormente: aprendizado por lote (*batch learning*) e aprendizado sequencial (*online learning*). No batch learning, os exemplos são fornecidos todos de uma só vez e estão disponíveis a todo momento para o vetor aluno. Já para o aprendizado online, os exemplos são apresentados sequencialmente e apenas uma vez. Nosso modelo corresponde a um algoritmo de aprendizado sequencial, e portanto trataremos apenas deste caso neste trabalho.

Para entendermos um pouco melhor a tarefa de classificação proposta, podemos fazer uma pequena interpretação geométrica de seu significado. Podemos definir formalmente esta separação como sendo $\sigma(\vec{S} \cdot \vec{J})$, onde a classificação ± 1 será dada de acordo com ângulo relativo entre os vetores de entrada e aluno, sendo divididos pelo hiperplano definido pelo ângulo $\pi/2$. Este hiperplano no qual $\sigma(\vec{S} \cdot \vec{J})$ muda de sinal é conhecido como fronteira/borda da decisão/dúvida. Como o tamanho dos vetores não influenciam na tarefa de classificação, é conveniente normalizá-los. Para facilitar alguns desenvolvimentos, a normalização será tal que:

$$\vec{J}^2 = \sum_i J_i^2 = N \quad (1.2)$$

$$\vec{S}^2 = \sum_i S_i^2 = N \quad (1.3)$$

Portanto, todos os vetores considerados para este aprendizado estão contidos em uma hipersfera N-dimensional de raio \sqrt{N} . Para compararmos a classificação feita pelo aluno e pelo professor, introduzimos uma variável chamada

sobreposição (*overlap*), definida por:

$$\rho = \frac{\vec{J} \cdot \vec{T}}{N} \quad (1.4)$$

Devido a normalização dos vetores T e J , a grandeza ρ é apenas o cosseno do ângulo θ entre estes dois vetores. Podemos portanto afirmar que o erro de generalização realizado pelo perceptron será simplesmente o arco cujo cosseno é ρ , ou mais formalmente:

$$\varepsilon_g = \frac{1}{\pi} \arccos \rho \quad (1.5)$$

Portanto, entender a dinâmica de ρ ao longo do aprendizado é essencial para compreender a eficiência da estratégia de treinamento utilizada.

A regra de aprendizagem online mais simples é a chamada *regra de Hebb* [31]. O algoritmo de Hebb consiste em alterar o valor das dimensões do acoplamento sináptico por erro/acerto. Se a dimensão ξ_i do input possuir o mesmo sinal que σ_T^μ então J_i será incrementado de $+1$, caso contrário será reduzido de -1 . Ou mais formalmente:

$$J_i^{\mu+1} = J_i^\mu + \xi_i^\mu \sigma_T^\mu \quad (1.6)$$

Após p exemplos, o vetor J poderá ser escrito como:

$$\vec{J} = \frac{1}{\sqrt{N}} \sum_{\mu=1}^p \vec{\xi}^\mu \sigma_T^\mu \quad (1.7)$$

Existem diversas regras de aprendizagem sequencial, e em sua grande maioria podem ser vistas como casos especiais do algoritmo de Hebb. Em 1958, Frank Rosenblatt desenvolveu o que seria o algoritmo mais famoso da história dos neurônios artificiais, o *perceptron* [32], uma simples alteração da regra de Hebb, porém com desempenho razoavelmente mais elevado. O perceptron consiste em alterar o vetor sináptico apenas em caso de erro, e mantê-lo inalterado em caso de acerto. O que faz muito sentido, já que após um número grande de exemplos a probabilidade de acertos é muito maior do que a de uma classificação incorreta, aprimorando assim a convergência da rede para um grande conjunto de exemplos. Formalmente a regra do perceptron consiste em:

$$\vec{J}^{\mu+1} = \begin{cases} \vec{J}^\mu + \frac{1}{\sqrt{N}} \vec{\xi}^\mu \sigma_T^\mu & \text{se } \xi^\mu \sigma_T^\mu < 0 \\ \vec{J}^\mu & \text{caso contrário} \end{cases} \quad (1.8)$$

É interessante notar que ambos os casos podem ser escritos na forma:

$$\vec{J}^{\mu+1} = \vec{J}^{\mu} + \frac{1}{\sqrt{N}} F^{\mu} \vec{\xi}^{\mu} \quad (1.9)$$

Onde a função F^{μ} é chamada de *amplitude de aprendizagem* ou *função de modulação* e tem um papel importantíssimo na dinâmica de aprendizado. Multiplicando os dois lados da eq. 1.9 por $\vec{\xi}^{\mu}$, temos:

$$\begin{aligned} \vec{J}^{\mu+1} \cdot \vec{\xi}^{\mu} &= \vec{J}^{\mu} \cdot \vec{\xi}^{\mu} + \frac{N}{\sqrt{N}} F^{\mu} \\ F^{\mu} &= \lambda^{\mu} - h^{\mu} \end{aligned} \quad (1.10)$$

Onde

$$\lambda^{\mu} = \frac{\vec{J}^{\mu+1} \cdot \vec{\xi}^{\mu}}{\sqrt{N}} \quad \text{e} \quad h^{\mu} = \frac{\vec{J}^{\mu} \cdot \vec{\xi}^{\mu}}{\sqrt{N}} \quad (1.11)$$

O que nos fornece uma interpretação clara de F^{μ} como sendo a amplitude de realinhamento do vetor sináptico, ou seja, a importância dada ao exemplo $\vec{\xi}^{\mu}$. Dada a natureza aleatória dos exemplos, a dinâmica que rege o vetor \vec{J} deve consequentemente ser de natureza estocástica. É portanto essencial definirmos parâmetros de ordem, na intenção de acompanhar a evolução do erro de generalização ε . Estes parâmetros de ordem são:

$$R^{\mu} = \frac{\vec{J}^{\mu} \cdot \vec{T}}{N} \quad \text{e} \quad Q = \frac{\vec{J}^{\mu} \cdot \vec{J}^{\mu}}{N} \quad (1.12)$$

Podemos agora derivar as dinâmicas para ρ e Q utilizando a 1.9 como ponto de partida. Multiplicando 1.9 por \vec{T} , temos:

$$\begin{aligned} \vec{J}^{\mu+1} \cdot \vec{T} - \vec{J}^{\mu} \cdot \vec{T} &= F^{\mu} \frac{\vec{T} \cdot \vec{\xi}^{\mu}}{\sqrt{N}} \\ N (R^{\mu+1} - R^{\mu}) &= F^{\mu} \frac{\vec{T} \cdot \vec{\xi}^{\mu}}{\sqrt{N}} \end{aligned}$$

Iterando 1.9 l vezes e repetindo o mesmo procedimento obtemos:

$$N \frac{(R^{\mu+l} - R^{\mu})}{l} = \frac{1}{l} \sum_{i=0}^{l-1} F^{\mu+i} \frac{\vec{T} \cdot \vec{\xi}^{\mu+i}}{\sqrt{N}} \quad (1.13)$$

Tomando o limite termodinâmico, onde $N \rightarrow \infty$, $l \rightarrow \infty$ porém com $l/N = d\Lambda$ podemos considerar que $\rho^{\mu} = \rho(\Lambda)$ e que durante um incremento de tempo $d\Lambda$ a diferença $\rho^{\mu+1} - \rho^{\mu}$ se torna também um diferencial $d\rho$. Nos levando assim a equação diferencial que rege ρ , dada por:

$$\frac{dR}{d\Lambda} = \langle Fu \rangle \quad (1.14)$$

Onde

$$u = \frac{\vec{T} \cdot \vec{\xi}^\mu}{\sqrt{N}} \quad (1.15)$$

É o campo local u do professor. O símbolo $\langle \dots \rangle$ representa a média tomada com relação a distribuição de exemplos. Precisamos agora da equação diferencial de Q . Para isso, elevamos ao quadrado os dois lados da eq. 1.9. O que nos fornece:

$$\begin{aligned} \vec{J}^{\mu+1} \cdot \vec{J}^{\mu+1} &= \vec{J}^\mu \cdot \vec{J}^\mu + 2F^\mu h^\mu + (F^\mu)^2 \\ N(Q^{\mu+1} - Q^\mu) &= F^\mu (F^\mu + 2h^\mu) \end{aligned}$$

Realizando o mesmo procedimento de iteração utilizado anteriormente:

$$N \frac{(Q^{\mu+l} - Q^\mu)}{l} = \frac{1}{l} \sum_{i=0}^{l-1} F^{\mu+i} (F^{\mu+i} + 2h^{\mu+i}) \quad (1.16)$$

Tomando novamente o limite termodinâmico e considerando agora $Q^\mu = Q(\Lambda)$ obtemos a equação:

$$\frac{dQ}{d\Lambda} = \langle F(F + 2h) \rangle \quad (1.17)$$

Onde h é idêntico ao definido na eq. 1.11, representando o campo local (ou campo de alinhamento) do aluno. Observando as equações 1.14 e 1.17 podemos notar que escolher a regra de aprendizado é equivalente a definir a função de modulação F do algoritmo. Seguindo a normalização estabelecida para esta demonstração, a forma correta da sobreposição R é dada pela relação $\rho = R/\sqrt{Q}$. Ou seja:

$$\begin{aligned} \frac{d\rho}{d\Lambda} &= \frac{1}{\sqrt{Q}} \frac{dR}{d\Lambda} - \rho \frac{1}{2Q^{3/2}} \frac{dQ}{d\Lambda} \\ \frac{d\rho}{d\Lambda} &= \frac{1}{\sqrt{Q}} \langle Fu \rangle - \frac{R}{2Q^{3/2}} \langle F(F + 2h) \rangle \\ \frac{d\rho}{d\Lambda} &= \frac{\langle 2QFu - RF(F + 2h) \rangle}{2Q^{3/2}} \end{aligned} \quad (1.18)$$

Fazendo algumas manipulações:

$$\begin{aligned}
\frac{d\rho}{d\Lambda} &= \left\langle \frac{2QFu - RF^2 - 2RFh}{2Q^{3/2}} \right\rangle \\
\frac{d\rho}{d\Lambda} &= \left\langle F \left(\frac{u}{Q^{1/2}} - \frac{Rh}{Q^{3/2}} \right) - \frac{RF^2}{2Q^{3/2}} \right\rangle \\
\frac{d\rho}{d\Lambda} &= \left\langle \frac{F}{Q^{1/2}} (u - \rho h) - \frac{\rho F^2}{2Q} \right\rangle \tag{1.19}
\end{aligned}$$

Em posse das equações 1.14, 1.17 e 1.19 podemos partir para o estudo da eficiência dos diferentes algoritmos de aprendizado. Por motivos de simplicidade trataremos apenas das regras de Hebb, perceptron e posteriormente iremos inferir sobre a classificação ótima.

Regra de Hebb

A regra de Hebb consiste em uma função de modulação $F = \sigma(u)$. Precisamos agora realizar as médias $\langle Fu \rangle$ e $\langle F(F + 2h) \rangle$. Para poder efetuar o cálculo destas grandezas é necessário perceber que no limite termodinâmico u e h são variáveis gaussianas correlacionadas, tais que:

$$\langle u^2 \rangle = 1, \quad \left\langle \frac{h^2}{Q} \right\rangle = 1, \quad \left\langle u \frac{h}{\sqrt{Q}} \right\rangle = \rho$$

Como os exemplos $\vec{\xi}$ aparecem nas equações de maneira indireta, basta calcular as médias das grandezas u e t .

Temos então:

$$\langle \sigma(u)u \rangle = \int_{-\infty}^{\infty} du |u| P(u) = \frac{1}{\sqrt{2\pi}} 2 \int_0^{\infty} u e^{-\frac{u^2}{2}} du = \sqrt{\frac{2}{\pi}}$$

$$\langle \sigma(u)(\sigma(u) + 2h) \rangle = 1 + \langle 2\sigma(u)h \rangle = 1 + 2 \int_{-\infty}^{\infty} dudh P(u, h) \sigma(u) h$$

$$\int_{-\infty}^{\infty} dudh P(u, h) \sigma(u) h = \int_{-\infty}^{\infty} dudh P(h|u) P(u) \sigma(u) h$$

Utilizando-se do fato que para variáveis gaussianas h e b com média nula e correlação R a esperança condicional é $E[h|u] = Ru$ chegamos em:

$$\int_{-\infty}^{\infty} dudh P(u, h) \sigma(u) h = \rho \int_{-\infty}^{\infty} \sigma(u) u P(u) du$$

Que é exatamente $\langle \sigma(u)u \rangle = \sqrt{\frac{2}{\pi}}$, portanto:

$$\langle \sigma(u)(\sigma(u) + 2h) \rangle = 1 + \langle 2\sigma(u)h \rangle = 1 + 2R\sqrt{\frac{2}{\pi}}$$

O que nos leva a:

$$\frac{dR}{d\Lambda} = \sqrt{\frac{2}{\pi}} \quad (1.20)$$

$$\frac{dQ}{d\Lambda} = 1 + 2R\sqrt{\frac{2}{\pi}} \quad (1.21)$$

Tomando como condição inicial $R(0) = Q(0) = 0$ e integrando as equações 1.20 e 1.21, obtemos:

$$R(\Lambda) = \sqrt{\frac{2}{\pi}}\Lambda \quad \text{e} \quad Q(\Lambda) = \Lambda + \frac{2}{\pi}\Lambda^2$$

Lembrando que $\rho = R/\sqrt{Q}$ chegamos finalmente em $\varepsilon(\Lambda)$ como sendo:

$$\rho(\Lambda) = \frac{\sqrt{\frac{2}{\pi}}\Lambda}{\sqrt{\Lambda + \frac{2}{\pi}\Lambda^2}} = \left(1 + \frac{\pi}{2\Lambda}\right)^{-1/2} \quad (1.22)$$

$$\varepsilon_g(\Lambda) = \frac{1}{\pi} \arccos \rho = \frac{1}{\pi} \arccos \left(1 + \frac{\pi}{2\Lambda}\right)^{-1/2} \quad (1.23)$$

Que no limite $\Lambda \rightarrow \infty$ possui um comportamento do tipo:

$$\varepsilon_g(\Lambda) \sim \frac{1}{\pi} \arccos \left(1 - \frac{1}{\pi\Lambda}\right) \sim \sqrt{\frac{2}{\pi}}\Lambda^{-1/2} \quad (1.24)$$

Perceptron

Para o perceptron, a função de modulação correspondente é $F = \sigma(u)\Theta(-uh)$ onde Θ é a função de Heaviside, definida por:

$$\Theta(x) = \begin{cases} 1 & \text{se } x > 0 \\ 0 & \text{caso contrário} \end{cases} \quad (1.25)$$

Novamente o procedimento é o mesmo, devemos efetuar as médias $\langle Fu \rangle$ e $\langle F(F + 2h) \rangle$, logo:

$$\langle Fu \rangle = \langle \sigma(u)u\Theta(-hu) \rangle = \langle |u|\Theta(-hu) \rangle$$

$$\langle |u|\Theta(-hu) \rangle = \int_{-\infty}^{\infty} dudhP(h)P(u|h)|u|\Theta(-hu)$$

Calculando as médias e fazendo as substituições necessárias, chegamos na equação diferencial para a sobreposição entre aluno e professor, dada por:

$$\frac{d\rho}{d\Lambda} = \frac{1 - \rho^2}{\sqrt{2\pi Q}} - \frac{\rho}{2\pi Q} \arccos \rho \quad (1.26)$$

É interessante notar que esta equação possui um ponto fixo em $\rho = 1$. Fazendo uma análise da aproximação assintótica de ρ a este ponto fixo, podemos encontrar o comportamento para o erro de generalização em $\Lambda \rightarrow \infty$ como sendo:

$$\varepsilon_g(\Lambda) \sim \frac{2^{1/3}}{3} \pi^{-1} \Lambda^{-1/3} \quad (1.27)$$

Que é um resultado inferior ao algoritmo de Hebb. No entanto, o algoritmo de Hebb possui boa capacidade de generalização apenas para exemplos distribuídos uniformemente, enquanto que o Perceptron pode ser expandido para situações na qual os exemplos são extraídos de qualquer distribuição arbitrária. Além disso, também é possível utilizar uma técnica chamada *quenching*, na qual a função de modulação para o perceptron diminui de acordo com ρ . Veremos na próxima seção que uma função de modulação variável de acordo com o estágio do aprendizado é muito mais eficiente do que funções estáticas como a dos casos apresentados até aqui.

Aprendizado ótimo

Após entender dois dos algoritmos básicos para a aprendizagem online, nos resta agora a pergunta: Qual a melhor performance que é possível obter utilizando uma distribuição uniforme de exemplos? O desenvolvimento a seguir foi primeiramente realizado por Osame Kinouchi e Nestor Caticha [33]

Primeiramente, devemos nos lembrar da equação que rege a sobreposição entre aluno e professor:

$$\frac{d\rho}{d\Lambda} = \left\langle \frac{F}{Q^{1/2}} (u - \rho h) - \frac{\rho F^2}{2Q} \right\rangle \quad (1.28)$$

Sabendo que $d\rho/d\Lambda$ representa o ganho informacional a cada exemplo apresentado, devemos então procurar a função de modulação que maximize a média do lado direito da equação 1.28, ou seja, precisamos que

$$\frac{\delta}{\delta F} \frac{d\rho}{d\Lambda} = 0$$

Que nos fornecerá o formato para a função de modulação ótima para o caso online. Lembrando que no limite termodinâmico $Q = J^2$ chegamos que a função F_{opt} será tal que:

$$F_{opt}^\mu = \sigma_T^\mu W_{opt}^\mu = J^\mu (\kappa^\mu - z^\mu) \quad (1.29)$$

Onde

$$\kappa^\mu = \frac{\sigma_T^\mu u^\mu}{\rho^\mu} \quad \text{e} \quad z^\mu = \sigma_T^\mu h^\mu \quad (1.30)$$

No entanto, este formato para a função de modulação necessita que o aluno possua a informação sobre o campo do professor u^μ , e que em qualquer caso realista é inacessível. Podemos então separar as variáveis de W_{opt} em duas categorias: variáveis visíveis $\mathbb{V} = \{h^\mu, \sigma_T^\mu\}$, e variáveis ocultas $\mathbb{H} = \{|b^\mu|\}$. De tal forma, a melhor função de modulação será aquela que satisfizer a relação 1.29 com relação a média das variáveis ocultas. Portanto, temos que a função ótima de modulação será dada por:

$$\bar{W}_{opt} = \langle W_{opt} \rangle_{\mathbb{H}|\mathbb{V}} = \int d\mathbb{H} P(\mathbb{H}|\mathbb{V}) J^\mu (\kappa^\mu - z^\mu) \quad (1.31)$$

Uma forma de facilitar o entendimento desta conta é lembrarmos da regra para probabilidades condicionais, que nos diz que:

$$P(\mathbb{H}|\mathbb{V}) = \frac{P(\mathbb{H}, \mathbb{V})}{P(\mathbb{H})}$$

Onde o fator de normalização $P(\mathbb{H})$ pode ser obtido por meio do teorema de Bayes. Nos deixando com:

$$P(\mathbb{H}|\mathbb{V}) = \frac{P(\mathbb{H}, \mathbb{V})}{\int d\mathbb{V} P(\mathbb{H}, \mathbb{V})}$$

Portanto, temos que a função ótima de modulação será dada pela integral:

$$\bar{W}_{opt} = \frac{\int d|b| P(b, h) W_{opt}}{\int d|b| P(b, h)} \quad (1.32)$$

Lembrando da definição de valor absoluto:

$$|b| = \begin{cases} b & \text{se } b > 0 \\ -b & \text{se } b < 0 \end{cases}$$

Chegamos em:

$$\int d|b|[\dots] = \int_0^\infty db[\dots] - \int_{-\infty}^0 db[\dots] \quad (1.33)$$

Utilizando novamente o fato de que h e b são variáveis gaussianas correlacionadas e realizando as médias necessárias, chegamos na função de modulação ótima, dada por:

$$W(\rho_\mu, J_\mu, z_\mu) = \frac{1}{\sqrt{2\pi}} J_\mu \lambda_\mu \exp\left(-\frac{z_\mu^2}{2\lambda_\mu^2}\right) \frac{1}{H(-z_\mu/\lambda_\mu)} \quad (1.34)$$

Onde

$$\lambda = \tan(\pi\varepsilon_g) = \frac{\sqrt{1-\rho^2}}{\rho} \quad \text{e} \quad H(x) = \frac{1}{2} \operatorname{erfc}\left(\frac{x}{\sqrt{2}}\right) \quad (1.35)$$

Substituindo a equação 1.34 na equação diferencial 1.19 e com uma mudança de variável obtemos:

$$\frac{d\rho}{d\Lambda} = \frac{1-\rho^2}{2\pi\rho} \int Dx \frac{e^{-x^2/\lambda^2}}{H(x/\lambda)} \quad (1.36)$$

Onde Dx é a medida gaussiana. Integrando numericamente esta equação e substituindo na relação $\varepsilon_g = \frac{1}{\pi} \arccos(\rho)$ podemos encontrar o erro de generalização. Chegamos então a $\varepsilon_g(\Lambda) = 0.88/\Lambda$, resultado que é superior ao algoritmo de Hebb ($\varepsilon_g \propto \Lambda^{-1/2}$) e ao perceptron ($\varepsilon_g \propto \Lambda^{-1/3}$). É interessante notar que a função de modulação depende explicitamente de ρ , de forma que o aprendizado possui diferentes amplitudes de acordo com a fase na qual a aprendizagem se encontra. Quanto mais próximo do professor (maior ρ), menos o estudante aprende com exemplos classificados corretamente (o comportamento de W para alguns valores de ρ está exemplificado na figura 1.7). De tal forma, é tentador considerar uma distribuição não uniforme de exemplos, de forma a maximizar o ganho de informação a cada passo, e evitar passos que tragam poucas mudanças em J . Portanto, iremos agora partir para a última situação, aprendizado ótimo com seleção de exemplos.

Podemos notar que a dependência da função de modulação com ρ se assemelha muito as diferenças cognitivas entre conservadores e liberais apresentadas na seção 2.2.

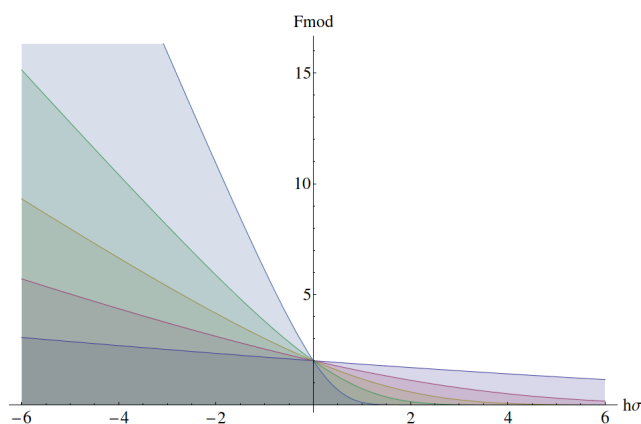


Figura 1.7: Gráfico da função de modulação em função de $z = h\sigma$ para diferentes valores do parâmetro de aprendizagem ρ .

Seleção de Exemplos

Partindo dos resultados anteriores, podemos então nos perguntar: Qual distribuição de exemplos maximiza o desempenho da função de modulação ótima? É interessante perceber que independente da distribuição dos exemplos a função de modulação deve ser a mesma. Este método também foi desenvolvido por Caticha e Kinouchi [33].

O primeiro passo para obtermos a melhor estratégia de apresentação de exemplos é notar que a equação 1.36 pode ser escrita como:

$$\frac{d\rho}{d\Lambda} = \frac{1 - \rho^2}{4\pi\rho} \int dh P(h) \exp(-h^2/\lambda^2) \left[\frac{1}{H(h/\lambda)} + \frac{1}{H(-h/\lambda)} \right] \quad (1.37)$$

Pois a média restante no lado direito da equação 1.36 depende apenas das variáveis visíveis. Devemos portanto definir a distribuição $P(h)$ que maximiza esta média. É simples notar que o termo $\exp(-h^2/\lambda^2)$ possui um máximo em $h = 0$ e também notamos que $H(0) = 1$, portanto a distribuição de exemplos ótima é $P(h) = \delta(h)$. Ou seja, os exemplos localizados na região conhecida como borda da dúvida são os que no total acrescentam mais informação à nossa rede neural. Neste caso, temos que:

$$h = \frac{\vec{J} \cdot \vec{\xi}}{N} = 0 \quad \text{e} \quad W_{opt} = \frac{1}{\sqrt{2\pi}} J_{\mu} \lambda_{\mu}$$

A princípio, podemos nos perguntar como estes exemplos devem ser escolhidos. A resposta metafórica para este problema é a de que o aluno deve

questionar o professor sobre assuntos perpendiculares a seu vetor de acoplamento sináptico, por isto esta técnica também leva o nome de questionários (*queries*). Substituindo $P(h)$ em 1.37 ficamos simplesmente com:

$$\frac{d\rho}{d\Lambda} = \frac{1 - \rho^2}{2\pi\rho}$$

Equação que pode ser integrada analiticamente, fornecendo uma solução do tipo:

$$\rho(\Lambda) = \sqrt{1 - e^{-2\Lambda/\pi}} \quad (1.38)$$

O que leva a um comportamento assintótico para o erro de generalização do tipo

$$\varepsilon_g(\Lambda) \sim \frac{1}{\pi} \exp\left(-\frac{\Lambda}{\pi}\right) \quad (1.39)$$

Resultado extremamente superior a todos os decaimentos algébricos apresentados anteriormente. O principal objetivo deste trabalho será analisar a influência da estratégia de seleção de exemplos em um modelo de dinâmica de opiniões. Serão utilizadas a função ótima com exemplos uniformes e a função ótima com exemplos perpendiculares.

É interessante definir uma energia/custo de aprendizado, tal que possua uma relação com a função de modulação análoga à relação entre força/energia, do tipo:

$$F^\mu = -\frac{\partial E^\mu}{\partial z^\mu} \quad (1.40)$$

que pela definição, deve ser igual a:

$$E^\mu = -\int F^\mu dz^\mu = \lambda^2 \log(P(\sigma|h)) \quad (1.41)$$

onde

$$P(\sigma|h) = H\left(-\frac{z}{\lambda}\right).$$

Esta energia será extremamente importante para nosso modelo, pois será utilizada como energia de interação social.

Capítulo 2

Modelo e métodos

Agora que já entendemos as motivações empíricas deste trabalho e temos todo o ferramental teórico que será necessário, podemos enfim começar o desenvolvimento do modelo a ser utilizado. Este trabalho pode ser considerado como uma extensão do modelo desenvolvido por Nestor Caticha, Renato Vicente e colaboradores [1], sendo composto de uma sociedade de agentes que interagem de maneira conformista, aprendendo sobre sua vizinhança social de maneira ótima. O modelo desenvolvido na próxima seção já foi estudado em um caso mais simples [1, 2], utilizando o perceptron como base para a aprendizagem social, possibilitando assim um estudo analítico, por meio de aproximações de campo médio, que quando comparados com dados empíricos permitiram uma associação entre ideologias políticas e estilos cognitivos [34], onde agentes com perfil cognitivo que valoriza a corroboração foram estatisticamente semelhantes a conservadores, e agentes que valorizam a novidade foram semelhantes aos liberais. Ao analisar dados de pesquisas de opinião foi também possível constatar que os padrões na distribuição dos dados não se alteram quando sua dimensionalidade é reduzida a 5, sugerindo assim que a utilização da teoria dos fundamentos morais é justificável e não deve ser fonte de erro sistemático. Nas páginas seguintes serão discutidas a estrutura formal do modelo e os métodos que serão utilizados para a obtenção das grandezas de interesse.

2.1 Modelo

A primeira imposição fundamental de nosso modelo está relacionada com a Teoria dos Fundamentos Morais. Devido a característica dimensional da moralidade, devemos levar em conta que em nosso estudo deve haver um número finito N de dimensões morais. Entretanto, por questões de genera-

lidade, não precisamos necessariamente considerar um valor fixo para N , ou seja, $N \in \mathbb{Z}$ tal que $0 < N < \infty$. Em seguida, levamos em consideração a proximidade entre o aprendizado sequencial de redes neurais artificiais com o aprendizado por reforço apresentado na seção 2.1. Nosso sistema será portanto constituído de agentes que aprendem/navegam pelo universo moral utilizando a mesma estratégia de aprendizado ótimo, apresentada anteriormente. Cada agente será representado por um vetor moral $\vec{\omega}$ (análogo ao vetor de acoplamento sináptico \vec{J}), onde o vetor moral de um agente i será dado por $\vec{\omega}_i = (\omega_{i1}, \dots, \omega_{iN})$, normalizado de forma que $|\vec{\omega}_i| = 1$. Cada componente ω do vetor pode ser considerada como a representação de uma dimensão moral (porém sem carga semântica associada) e de sua importância relativa, de forma que cada atributo moral pode possuir no máximo um valor entre $\pm \frac{1}{\sqrt{N}}$. Desta forma, cada agente possuirá uma matriz moral própria, e que será adaptada de acordo com a interação com parceiros sociais.

Cada agente i possui também uma estratégia cognitiva ρ_i que define o grau de importância dado ao erro e à corroboração, onde agentes com $\rho \approx 1$ dão grande importância a exemplos mal-classificados (novidade), e agentes com $\rho \approx 0$ dão igual relevância a todos os assuntos discutidos (conforme apresentado na figura 1.7). Este parâmetro ρ é análogo ao *overlap* das redes neurais artificiais e será considerado como uma representação do estilo cognitivo do agente. Consideraremos que os agentes já congelaram seus estilos cognitivos, e que estes serão fixos ao longo do tempo, pois acredita-se que esta característica é desenvolvida na infância e não se altera ao longo do período de vida do indivíduo. Portanto, teremos uma sociedade uniforme em seus estilos cognitivos. Em estudos anteriores [34], ficou evidenciado que o comportamento de sociedades não homogêneas em ρ é qualitativamente o mesmo. O que é de se esperar, dado que a dependência com o parâmetro de aprendizagem é assimétrico com relação a troca de índices, fazendo com que suas estatísticas mantenham o mesmo comportamento, gerando apenas uma multimodalidade na distribuição total de opiniões h .

A vizinhança social do agente i será definida por um grafo auxiliar $\mathcal{G}(\mathcal{V}, \mathcal{A})$, onde \mathcal{V} é o conjunto dos vértices do grafo e \mathcal{A} é o conjunto de arestas. Este grafo define a estrutura da sociedade, e pode possuir diversos formatos. Para este trabalho será utilizada uma rede do tipo Barabási-Albert [35], porém podem ser utilizados qualquer tipo de rede (quadrada, Watts-Strogartz, pequeno mundo, etc). O modelo de Barabási-Albert (BA) consiste em uma rede do tipo *livre de escala* na qual a probabilidade de se acrescentar uma conexão à um nó é igual a

$$p_i = \frac{k_i}{\sum_j k_j},$$

onde k_i é o número de arestas partindo do nó i . Portanto, sendo proporcional ao número de conexões já realizadas, propriedade conhecida também como conectividade preferencial. O número de vizinhos em uma rede BA é portanto inhomogêneo, podendo-se controlar apenas o número médio de vizinhos K . Além disso, o coeficiente de clustering, segue uma lei de potências do tipo

$$C \sim n^{-0.75} ,$$

escalando assim com o tamanho da rede n , diferentemente de outras redes, como a de pequeno-mundo onde o coeficiente de clustering depende apenas do grau médio da rede. Este tipo de propriedades já foram observadas em diversas redes sociais reais, como por exemplos na internet, Facebook[©], crescimento urbano etc, sendo portanto uma escolha interessante para a aplicação de nosso modelo. Na figura 2.1 é possível observar um exemplo de rede gerada utilizando o algoritmo BA, na qual fica evidente a conectividade preferencial, onde certos nós da rede possuem um grau muito mais elevado que outros.

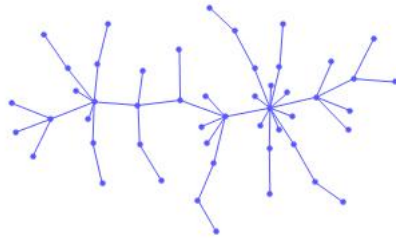


Figura 2.1: Exemplo de rede gerada utilizando o algoritmo de Barabási-Albert.

A partir dos parâmetros estabelecidos, a dinâmica ocorre da seguinte forma: sorteia-se uma aresta, e desta escolhe-se um dos agentes (vértices) como *aluno*. Cada agente poderá conversar com seus parceiros sociais sobre assuntos $\vec{X} = (X_1, \dots, X_N)$ contidos dentro do mesmo espaço moral de N dimensões. Por motivo de simplicidade, e também devido as escalas de tempo da dinâmica, consideraremos inicialmente que os agentes discutem entre si apenas o assunto médio \vec{Z} , chamado de *Zeitgeist* (termo alemão que significa “o espírito da época”). O vetor \vec{Z} é definido por:

$$\vec{Z} \propto \frac{1}{P} \sum_{\vec{X}_\mu \in \mathcal{X}} \vec{X}_\mu$$

Devido a simetria esférica do problema, podemos definir o vetor de Zeitgeist como direção preferencial do espaço, ou seja, $\vec{Z} = (1, \dots, 1)/\sqrt{N}$. Para compreender o efeito das estratégias de convencimento, consideraremos também a existência de um vetor *oráculo* \vec{Z}_O , tal que $\vec{Z}_O = -\vec{Z}$. Este vetor possuirá um efeito semelhante ao de um agente fixo na sociedade, servindo como uma segunda direção de quebra de simetria, competindo assim com o vetor de Zeitgeist.

Com uma probabilidade α os agentes consultarão o oráculo sobre assuntos da forma $(\vec{X}^\mu, \sigma_{Z_O})$, extraídos de uma distribuição de probabilidades $P(\vec{X})$ que será chamada de *estratégia de convencimento*. Com probabilidade $1 - \alpha$ as conversas se darão dentro da sociedade, e serão do tipo $(\vec{Z}, \sigma_{\vec{\omega}_j})$. O parâmetro α será chamado de *exposição*, e está relacionado com a quantidade de tempo que os agente estão expostos a uma influência externa.

Uma outra grandeza importante que devemos definir é a chamada *convicção*, sendo

$$h_i^\mu = \vec{\omega}_i \cdot \vec{X}^\mu,$$

para os exemplos discutidos com o oráculo, e

$$h_i = \vec{\omega}_i \cdot \vec{Z},$$

para as conversas intra-sociais. Análoga ao campo local no caso das redes neurais artificiais. A convicção representa o quanto o agente está convicto de sua classificação sobre o assunto \vec{X} .

Outra grandeza importante é a opinião z , dada por:

$$z = h_i \sigma(h_j) \quad \text{ou} \quad z^\mu = h_i^\mu \sigma(h_{Z_O}^\mu).$$

Onde:

$$\sigma(x) = \text{sgn}(x) = \begin{cases} 1 & \text{se } x > 0 \\ -1 & \text{se } x < 0 \end{cases},$$

Quando a classificação do agente i é a mesma que a do oráculo e/ou de seu parceiro social a opinião será positiva, e quando ambos discordarem a opinião será negativa. Portanto, $(z > 0)$ representa a concordância entre os agentes, $(z < 0)$ representa a discordância, e quando $(z \approx 0)$ a classificação é ambígua. A energia de interação entre os agentes e seus parceiros sociais (outros agente e/ou oráculo) será idêntica a energia ótima de aprendizado, dada por:

$$E = -\lambda^2 \log \left(H\left(-\frac{z}{\lambda}\right) \right).$$

Onde novamente $H(x) = \frac{1}{2}\text{erfc}(x/\sqrt{2})$. No entanto, para a interação em uma vizinhança, a forma total da energia deve levar em conta a opinião de todos os parceiros sociais, sendo portanto:

$$E_{i,\mathcal{V}} = -\lambda^2 \sum_{(i,j) \in \mathcal{V}} \log \left(H\left(-\frac{z_{ij}}{\lambda}\right) \right).$$

Novamente $\lambda = \sqrt{1 - \rho^2}/\rho$.

A Hamiltoniana total do sistema será dada por:

$$\mathcal{H} = (1 - \alpha) \sum_{\mathcal{V} \in \mathcal{G}} E_{i,\mathcal{V}} + \alpha \sum_i E(\vec{\omega}_i, \vec{Z}_O, \langle \vec{X} \rangle)$$

É fácil notar que o valor mínimo da energia de interação ε ocorre quando há concordância entre os interlocutores, ou seja, o termo z é igual a 1. No entanto, consideraremos apenas situações nas quais existem flutuações em torno do estado fundamental. Portanto, será utilizado o mesmo ferramental teórico da mecânica estatística clássica.

2.2 Métodos

Devido ao princípio de máxima entropia (Apêndice A), podemos afirmar que a distribuição de probabilidades dos estados de nosso sistema será a distribuição de Boltzmann, dada por:

$$P(\mathcal{H}) = \frac{1}{Z} e^{-\beta \mathcal{H}},$$

onde Z é uma constante de normalização (função de partição). O parâmetro β é um multiplicador de Lagrange derivado do vínculo de que \mathcal{H} mantenha um valor médio constante, representando o grau de flutuação aceito em torno do estado de mínima energia (estado fundamental), equiparável ao inverso da temperatura termodinâmica. Este parâmetro pode ser interpretado como uma *pressão social*, pois quanto mais elevado o valor de β , menos um agente poderá divergir de sua vizinhança social.

Para entender o comportamento de nosso sistema e extrair as grandezas de interesse será necessário amostrar a distribuição de probabilidades associada. Para realizar esta tarefa utilizaremos a técnica de amostragem de Monte Carlo. Mais especificamente, utilizaremos um algoritmo do tipo Metropolis [36].

O algoritmo de Metropolis consiste em aceitar mudanças no estado do sistema de acordo com a regra:

$$P_A(\vec{\omega}'|\vec{\omega}) = \min(1, e^{-\beta\Delta\varepsilon}) = \begin{cases} 1 & \text{se } \Delta\varepsilon < 0 \\ e^{-\beta\Delta\varepsilon} & \text{se } \Delta\varepsilon > 0 \end{cases}$$

Desta forma, assegura-se que haja uma boa amostragem em torno das regiões de alta probabilidade, porém com a possibilidade de flutuações na energia, eliminando o problema de armadilhamento em máximos locais.

Para realizarmos as mudanças no estado do sistema, altera-se o vetor moral $\vec{\omega}_i$, do agente sorteado como aluno no passo μ , sorteando-se um novo vetor $\vec{\omega}'_i$ dentro de um cone N dimensional de raio u , ou seja:

$$\vec{\omega}'_i = \frac{\vec{\omega} + \vec{u}}{|\vec{\omega} + \vec{u}|}$$

Onde \vec{u} é um vetor com componentes sorteadas uniformemente no intervalo $[-u, u]$, em todas as N dimensões. O tamanho de u é tal que a taxa de aceitação se mantenha sempre próxima de 0.5, considerado como um valor ideal para uma boa amostragem do sistema. A partir desta dinâmica é possível extrair as grandezas de interesse, ou seja, os momentos da distribuição de Boltzmann associada. Veremos que os parâmetros β e ρ são extremamente importantes para caracterizarmos transições de fase de nosso sistema.

O caso α igual a 0 (nenhuma exposição) já foi amplamente estudado [34], onde foi observada uma transição de fase no espaço ρ x β , que será apresentada mais a frente em nosso trabalho. Seus resultados foram também comparados com assinaturas estatísticas retiradas dos questionários desenvolvidos por Jonathan Haidt, mostrando ótima concordância entre teoria e dados empíricos. Desta comparação, foi possível confirmar a influência do estilo cognitivo na afiliação política dos agentes, onde agentes com um parâmetro de aprendizado (ρ) mais elevado puderam ser associados a assinaturas estatísticas de liberais, enquanto que os conservadores puderam ser associados a um perfil cognitivo mais corroborativo. Na próxima seção serão apresentados os resultados obtidos para o caso geral ($\alpha > 0$), onde estes comportamentos estudados ainda devem se fazer presentes, assim como algumas características oriundas do termo adicional na Hamiltoniana.

O principal objetivo deste trabalho será analisar o efeito do vetor oráculo (campo externo) no sistema, bem como a influência das estratégias de convencimento na eficácia das mudanças coletivas de opinião. Serão utilizados para este trabalho duas estratégias: uma distribuição de exemplos uniforme na superfície da hiper-esfera N -dimensional e uma distribuição de exemplos perpendiculares aos agentes. Vimos na seção 2.3 que a aprendizagem sem ruído na borda da dúvida leva a um decaimento exponencial com o número de exemplos no

cenário professor aluno, enquanto que exemplos distribuídos na esfera levam a uma queda do erro proporcional ao inverso do número de exemplos.

A pergunta que podemos fazer, inspirados por este resultado, é se haverá mudança mais rápida do consenso de uma população se houver uma preocupação em expô-las às suas dúvidas? Possivelmente não ocorrerá diferença tão significativa quando a observada na aprendizagem sequencial, dado que a dinâmica de aprendizado social ocorre na presença de ruído ($\beta > 0$).

Para sortear um exemplo perpendicular são extraídos vetores da distribuição uniforme, depois tomada a projeção destes vetores na direção perpendicular ao agente sorteado. Todos os vetores considerados neste modelo são normalizados, portanto todas as medidas de similaridade consistem apenas no produto escalar. Serão obtidas as magnetizações com relação ao oráculo, os tempos de auto-correlação e as respectivas distribuições de opinião. Como apresentado anteriormente, a estratégia na qual os exemplos estão localizados na borda da dúvida possui um decaimento exponencial do erro de generalização, enquanto que a distribuição uniforme apresenta um erro que decai de forma polinomial. Espera-se portanto que a estratégia na qual há seleção de exemplos seja superior à estratégia uniforme. Queremos estudar a dependência com α , na qual ocorrerá uma competição entre os termos de convencimento (Oráculo) e o termo de interação social (Zeitgeist).

Para melhor entender o comportamento macroscópico de nosso sistema é preciso definir os chamados parâmetros de ordem, grandezas extraídas da distribuição que nos fornecem uma compreensão do grau de organização de nosso sistema. Normalmente são definidos de forma que quando o parâmetro de ordem está próximo de zero o sistema se encontra em uma fase desordenada, e quando o parâmetro de ordem é não nulo o sistema está em uma fase ordenada. Em nosso modelo, assim como em um modelo de magnetismo, os parâmetros de ordem principais são a magnetização:

$$m = \langle h \rangle = \frac{1}{n} \sum_{i=1}^n \vec{\omega}_i \cdot \vec{Z} , \quad (2.1)$$

e a magnetização com relação ao oráculo

$$m_O = \langle h_O \rangle = \frac{1}{n} \sum_{i=1}^n \vec{\omega}_i \cdot \vec{Z}_O . \quad (2.2)$$

As magnetizações (ou opiniões médias) podem ser vistas como uma medida de coesão social, e fornecem informações sobre a existência de uma direção de quebra de simetria moral bem como sobre a distribuição dos agentes em seu entorno.

Os tempos de relaxação τ do sistema, que podem ser extraídos da auto-correlação do sistema, dada por:

$$c(t) = \langle m(t)m(t+t') \rangle - \langle m(t) \rangle \langle m(t+t') \rangle , \quad (2.3)$$

sendo que a auto-correlação pode ser parametrizada por:

$$c(t) \propto \exp(-t/\tau) \quad (2.4)$$

Portanto, é possível extrair os tempos de relaxação do sistema por meio de uma regressão linear do logaritmo de $c(t)$, calculado utilizando-se a equação 2.3. Como será mostrado na próxima seção, a auto-correlação possui um decaimento aproximadamente exponencial, justificando assim esta parametrização.

Capítulo 3

Resultados

O principal objetivo deste trabalho é a caracterização e o estudo da influência das estratégias de convencimento em mudanças coletivas de opinião. Para isto, serão analisadas as dependências com os parâmetros ρ , β e α , bem como a influência da estrutura da rede por meio do número médio de vizinhos K . Em trabalho anterior de Nestor Caticha e Jonatas César [34] foram observadas transições de fase, que como veremos se mantêm para o caso $\alpha \neq 0$, com grande influência no estado final da sociedade.

Primeiramente, devemos garantir que a distribuição convirja para um estado estacionário. Para isto simplesmente observamos o comportamento de um dos parâmetros de ordem por um longo período de tempo. Espera-se que este deva permanecer constante (ou pelo menos flutuando em torno do valor médio). A convergência do algoritmo está representada na figura 3.1. É possível notar que o valor de m converge rapidamente, se mantendo aproximadamente constante por um número grande de passos, indicando assim a convergência da amostragem. Cada passo da dinâmica representa o número de movimentos necessários para que todos os agentes sejam atualizados uma vez em média.

Os primeiros passos das simulações não são utilizados para o cálculo das grandezas termodinâmicas, e são popularmente chamados de *passos de termalização* (ou burn-in em inglês). Esta etapa é importante para garantir que não haja memória do estado inicial, bem como que a amostragem possua o grau de flutuação correspondente à temperatura simulada.

É interessante também avaliar o decaimento das auto-correlações, para que não cometamos erros grosseiros ao estimar os tempos característicos utilizando o modelo exponencial (eq. 2.4). Para tal, devemos analisar o comportamento da auto-correlação em uma escala monolog ($\log c(t) \times t$) apresentado na figura 3.2.

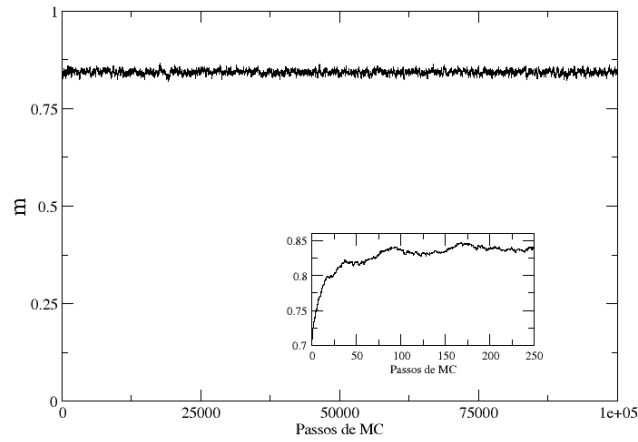


Figura 3.1: Convergência do valor da magnetização em função do número de passos de Monte Carlo realizados. O sub-gráfico representa uma ampliação dos 250 primeiros passos.

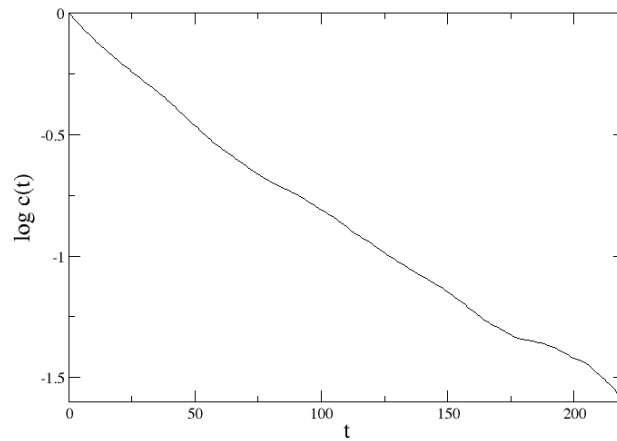


Figura 3.2: Gráfico de $\log c(t) \times t$ indicando o decaimento exponencial para as autocorrelações. Para a obtenção desta curva foram mantidos os parâmetros constantes em $\beta = 1$, $\rho = 0.3$ e $\alpha = 0$.

Como podemos notar, a curva possui um comportamento bem próximo do linear. Portanto, realizar uma regressão linear com relação a reta obtida e extrair o coeficiente linear é um bom método de determinação para o valor do tempo de readaptação τ , lembrando que este tempo característico depende de ρ , β e K .

3.1 Dependência com ρ e β

Para estudar a dependência da dinâmica com o parâmetro de aprendizagem (ρ) e a pressão social (β), foram realizadas simulações com 400 agentes, distribuídos em uma rede Barabasi-Albert com um número médio de vizinhos igual a 8. Em todas as realizações da dinâmica foram realizados 2000 passos iniciais com $\alpha = 0$, de forma que a sociedade entrasse em equilíbrio em torno do vetor de Zeitgeist, para só então ser realizada a dinâmica completa e a obtenção das medidas. Isto foi feito para que pudessem ser obtidos os tempos de readaptação do sistema no caso em que há mudança coletiva de opinião. Lembrando que a termalização garante que não haja informação sobre o estado inicial nas medidas de magnetização.

O primeiro resultado a ser analisado são os diagramas de fase no espaço $\rho \times \beta$ para o caso limite $\alpha = 0$ (Figura 3.3), no qual se recuperam os resultados de J. César, e o caso $\alpha = 1$ (Figura 3.4a e 3.4b), para garantir que não há mudança qualitativa de comportamento, e que o sistema ainda apresenta o mesmo tipo de transição de fase. Neste diagrama de fase podemos observar uma transição contínua, na qual a sociedade passa de um estado desordenado ($m = 0$) para um estado ordenado ($m \neq 0$). As faixas com magnetização aproximadamente constante apresentadas neste diagrama podem ser associadas também com diferentes afiliações políticas da sociedade, como mostrado em [2] e [34]. Essas afiliações políticas são tais que sociedades com menor magnetização estão usualmente associadas aos liberais, ou seja, liberais tendem a possuir uma maior dispersão de opiniões.

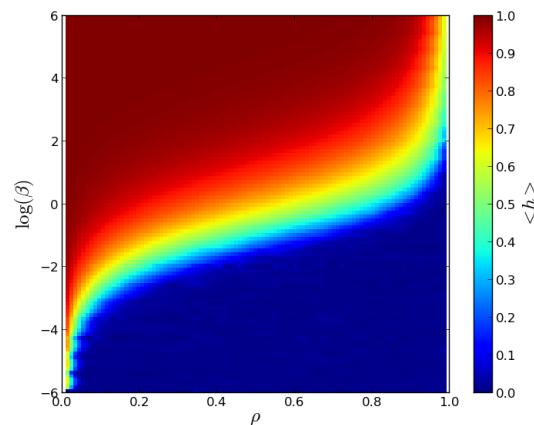


Figura 3.3: Diagrama de fase para a magnetização no espaço $\rho \times \beta$ no caso em que $\alpha = 0$.

Na situação em que $\alpha \neq 0$ a transição de fase se mantém para as duas estratégias de convencimento. No entanto, a linha de transição ocorre em locais distintos. Os exemplos sorteados na região da borda da dúvida ($h_i = 0$) apresentaram uma magnetização maior do que a de exemplos sorteados uniformemente para um mesmo valor de β e ρ , portanto podemos afirmar que a apresentação de exemplos na borda implica em uma menor mudança de dispersão de opiniões (resultado que será reforçado mais a frente). Também é possível afirmar que a fase ordenada ocorre em uma região de menor pressão social, portanto a estratégia de convencimento na qual $h_i = 0$ possui uma maior proximidade com o caso $\alpha = 0$. Lembrando que, quando $\alpha = 1$ não há interação intra-social, e os agentes não possuem vínculo direto entre si, portanto o estado final depende apenas do vetor oráculo. É de esperar que os casos no qual $0 < \alpha < 1$ a magnetização seja uma espécie de combinação do resultado apresentado na figura 3.3 com o resultado apresentado na figura 3.4.

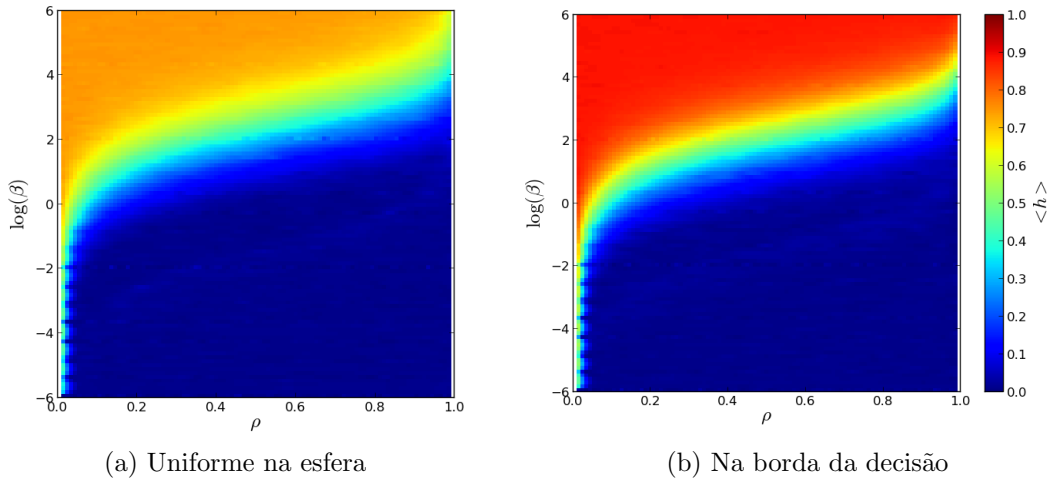


Figura 3.4: Diagrama de fase para o valor absoluto da magnetização no espaço $\rho \times \beta$ no caso em que $\alpha = 1$ para as duas estratégias estudadas. (a) Exemplos sorteados da distribuição uniforme na N-esfera. (b) Exemplos sorteados tais que $h_i = 0$.

Para compreender a influência dos parâmetros na coesão social, foram avaliados também os histogramas de sobreposição intra-social ρ_{ij} , onde

$$\rho_{ij} = \vec{\omega}_i \cdot \vec{\omega}_j \quad (3.1)$$

Devido a normalização dos vetores ω , o valor da sobreposição fica res-

trito ao intervalo $[-1, +1]$. O formato do histograma obtido nos fornece informações sobre a dispersão de vetores morais da sociedade, sendo que uma sociedade homogeneamente distribuída possuirá um histograma com uma média centrada em 0 e simétrico nas duas direções, enquanto que uma sociedade com opiniões semelhantes possuirá um histograma assimétrico, com uma maior número de ocorrências perto de $\rho_{ij} = 1$. Os histogramas obtidos estão apresentados nas figuras 3.5a e 3.5b.

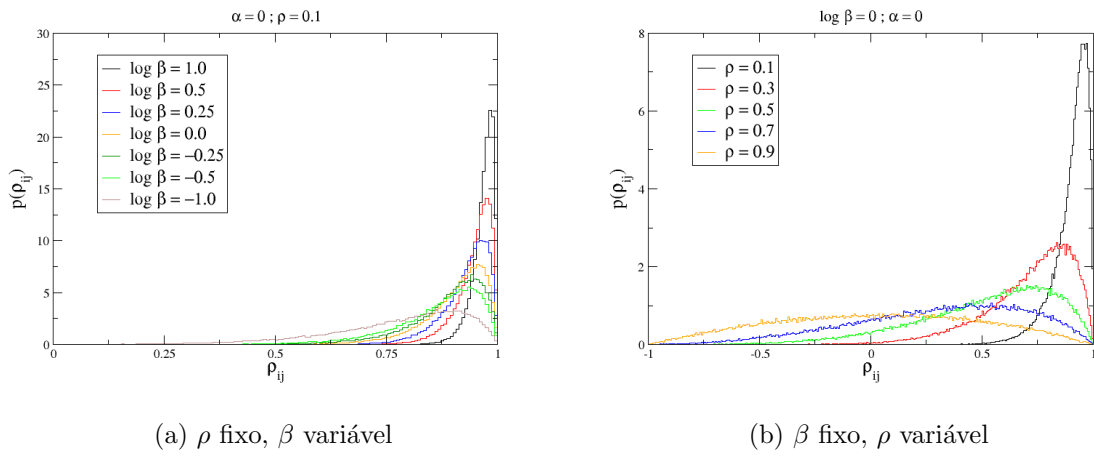


Figura 3.5: Histograma para as sobreposições intra-sociais ρ_{ij} para diversos valores do parâmetro de aprendizagem ρ e da pressão social β .

Podemos notar neste que agentes com um menor valor para o parâmetro de aprendizagem (associados aos conservadores) possuem uma maior coesão, com um histograma estreito e distribuído próximo de $\rho_{ij} = 1$. Já os agentes que possuem um perfil cognitivo que valoriza a novidade, apresentam um histograma mais largo, com uma maior dispersão de vetores morais. Também é possível notar o efeito da pressão social no comportamento coletivo: quanto maior o valor de β , mais os agentes se aproximam. Portanto, já neste resultado, podemos observar que um aumento da pressão social eleva o grau de conservadorismo da sociedade.

Analisamos também a dependência destes histogramas com a variável de exposição ao oráculo (α), apresentado na figura 3.6. Neste caso, é possível notar que o aumento da influência do termo do oráculo na energia implica em uma maior dispersão de matrizes morais. No entanto, na região próxima a transição 3.6b há uma tendência ao reagrupamento da sociedade, correspondendo ao máximo da magnetização m_O , apresentada na figura 3.7. O decaimento na coesão social com o aumento de α está relacionado com o for-

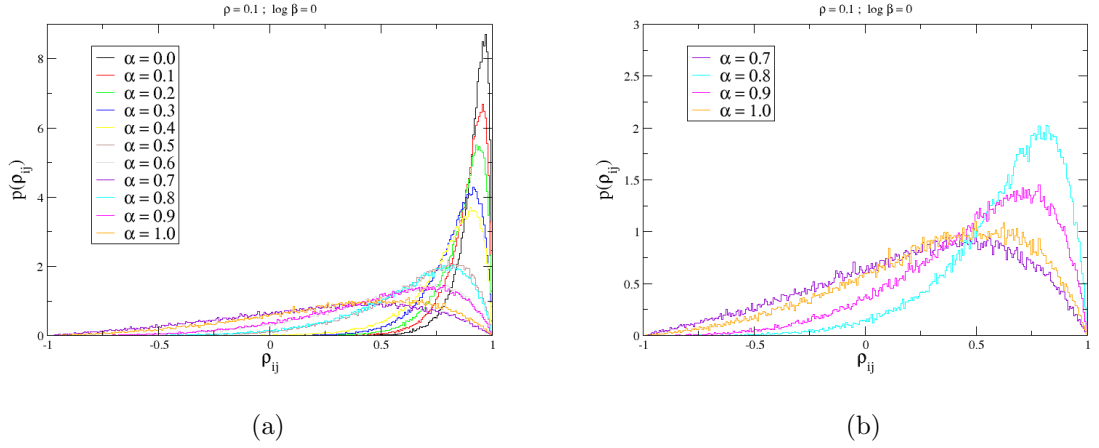


Figura 3.6: (a) $p(\rho_{ij})$ em diversos valores de α . (b) $p(\rho_{ij})$ na região próxima a transição de fase em α .

mato do termo da energia de interação com o oráculo, no qual os agentes não interagem entre si, fazendo com que a sociedade seja mais decorrelacionada e mais suscetível à influência da variância da estratégia de convencimento, cujo caráter probabilístico introduz mais uma escala de flutuação (pois são apresentados exemplos distintos a cada passo).

Em seguida, foi analisada a semelhança da sociedade com o oráculo, através das curvas de $m_O = \langle \vec{\omega} \cdot \vec{Z}_O \rangle$, em função do parâmetro de exposição ao oráculo (α), cobrindo todo o intervalo $0 < \alpha < 1$. Foi possível então observar uma transição de fase com relação a este parâmetro, na qual a sociedade altera abruptamente sua orientação moral; passando a ter o oráculo como direção preferencial. Os resultados obtidos estão representados na figura 3.7. Também é possível notar que a transição de fase observada nas figuras 3.3 e 3.4 se mantém, fazendo com que a magnetização máxima diminua com o aumento de ρ , passando continuamente da fase ordenada para a fase desordenada. Outra característica interessante é a diminuição progressiva da magnetização do sistema para valores de $\alpha \approx 1$, que pode ser entendida como uma mudança gradual do diagrama de fase da figura 3.3 para o diagrama da figura 1.2. Esta diminuição da magnetização está também representado nos histogramas da figura 3.6, onde fica evidenciado que a magnetização é altamente relacionada à coesão social. Devemos portanto atribuir esta diminuição da magnetização a influência do termo de interação do oráculo, que depende explicitamente da estratégia de convencimento e de sua variância.

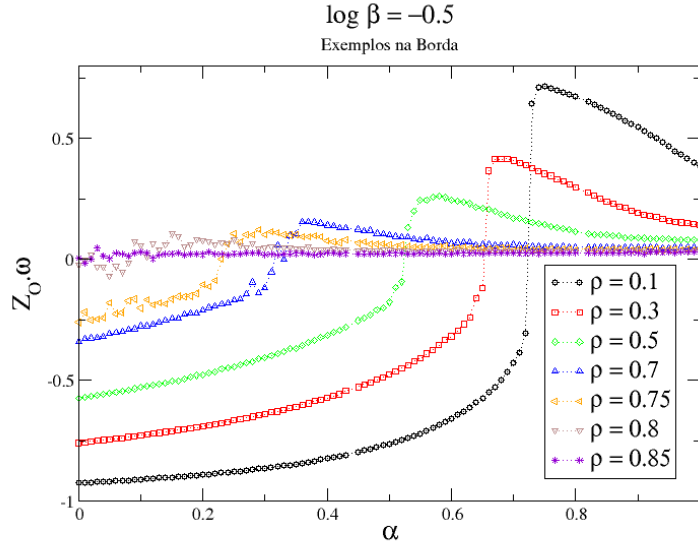


Figura 3.7: Gráfico m_O x α para diversos valores de ρ . As barras de erro são menores que os símbolos e foram omitidas por questão de clareza. Para este β o valor de ρ_c no qual a sociedade se encontra em um estado desordenado é de aproximadamente $\rho = 0.85$ (figura 3.3). Para este gráfico foram utilizados exemplos na borda da dúvida, entretanto o comportamento é o mesmo para exemplos distribuídos na esfera.

É possível também notar que o α_c no qual ocorre a transição de fase é dependente do parâmetro ρ , sendo que agentes associados com indivíduos liberais necessitam de uma menor exposição ao oráculo para que esta transição ocorra. É possível determinar o valor de α_c utilizando-se do efeito conhecido como desaceleração crítica (critical slowing down), fazendo com que os tempos característicos de readaptação diverjam na transição, como mostrado na figura 3.8. Portanto, para determinar α_c , precisamos apenas de $\operatorname{argmax}(\tau(\alpha))$.

Podemos então estudar a dependência de α_c com os parâmetros do modelo (ρ, β) . Os resultados obtidos podem ser utilizados para a obtenção de um outro diagrama de fase. Devido ao grande custo computacional, é possível apenas estudar cortes deste diagrama. Para que seja possível analisar a dependência é necessário realizar um corte em dois sentidos do diagrama: transversal (ρ fixo) e outro longitudinal (β fixo), possibilitando assim uma compreensão das tendências existentes neste espaço. Observando os resultados obtidos fica evidente que agentes com menor valor do parâmetro de

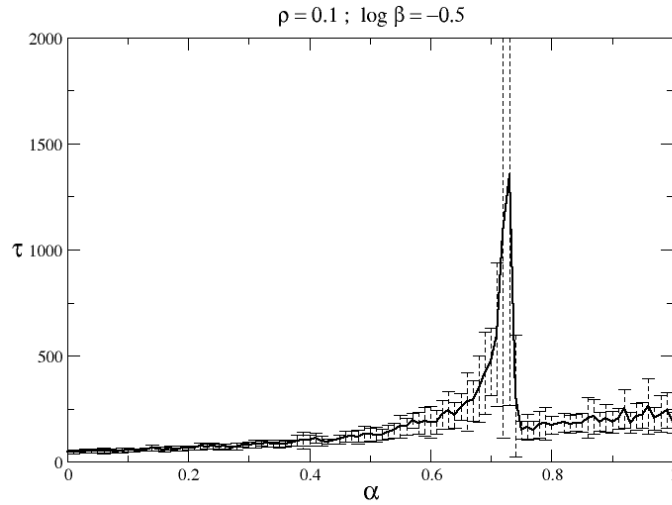


Figura 3.8: Curva τ x α , onde é possível observar o efeito da desaceleração crítica próximo a transição de fase.

aprendizagem (associados aos conservadores) possuem um α_c maior do que agentes que procuram novidades (associados aos liberais). Assim como ocorre em sociedades reais, um aumento de pressão social eleva o grau de conservadorismo da sociedade, legitimando o comportamento observado em nosso modelo. Esta tendência ao conservadorismo pode ser notada por meio da aproximação das curvas em direção a maiores valores da exposição crítica.

Nestes cortes, fica também evidenciado que agentes liberais possuem uma maior sensibilidade às variações de pressão social, tendo as maiores mudanças de α_c dentro do espectro em ρ . É também notável que α_c diminui rapidamente com relação à diferentes estratégias cognitivas, possuindo valores consideravelmente menores independente da pressão social. O mesmo ocorre para o máximo da magnetização em relação ao oráculo, sofrendo rápida diminuição em função de ρ , e uma maior tendência ao conservadorismo com um aumento da pressão social. Podemos então definir um diagrama de fase no espaço ρ x β , porém com relação ao valor máximo de m_O . Ambos os cortes nos espaços α_c e $\max m_O$ estão apresentados na figura 3.9.

Na figura 3.10 estão apresentados os diagramas de fase, obtidos por meio de uma interpolação de primeiros vizinhos dos cortes da figura 3.9. Esta interpolação é feita da seguinte maneira: cada ponto é dividido em um número menor de pontos, e então o valor de cada ponto intermediário é determinado por meio de uma média com os primeiros vizinhos. Esta interpolação não fornece o formato exato da curva, porém nos provém de informações com relação ao comportamento da grandeza estudada. Podemos notar que a tendência apresentada nos cortes é mantida nos dois diagramas.

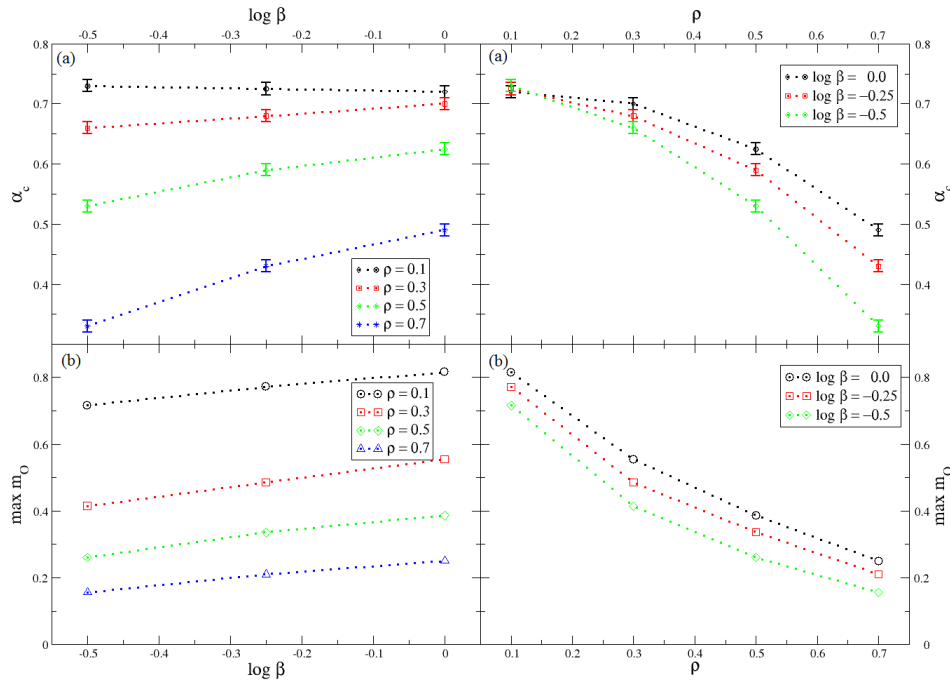


Figura 3.9: (a) Cortes do diagrama de fase para α_c (b) Cortes do diagrama de fase para $\max m_O$. O tamanho das barras de erro é de 0.01 (α) e está relacionado com a precisão utilizada na variação do parâmetro α . O erro em $\max m_O$ é menor que os símbolos.

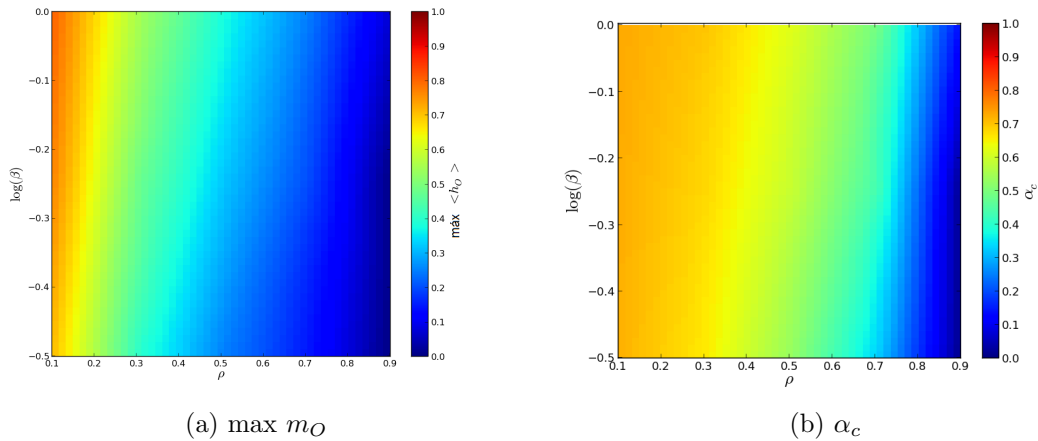


Figura 3.10: (a)Máximo da magnetização em torno do oráculo, em função dos parâmetros ρ e β . (b) Exposição crítica em função dos parâmetros ρ e β .

Resta agora entender a dependência do sistema com o número de vizinhos da rede, para posteriormente comparar as duas estratégias.

3.2 Dependência com o Número de Vizinhos

Para poder compreender o efeito da vizinhança na dinâmica de mudanças de opinião, foram simuladas sociedades com 400 agentes, em uma rede do tipo Barabasi-Albert, tendo uma pressão social constante $\beta = 1$, variando-se então o número de vizinhos. Foram obtidas as mesmas grandezas que na seção anterior: a magnetização com relação ao oráculo, os histogramas da sobreposição intra-social, e os cortes do diagrama de fase.

Devemos esperar que o aumento do número de vizinhos possua um efeito semelhante ao aumento da pressão social, pois implica em uma maior rigidez para mudanças na rede, causando assim uma menor dispersão de matrizes morais dentro da sociedade. Devemos lembrar que o termo para a interação do oráculo é independente da estrutura da rede, portanto em $\alpha = 1$ é esperado que todas as redes apresentem exatamente o mesmo resultado, pois este deve depender apenas dos parâmetros β e ρ .

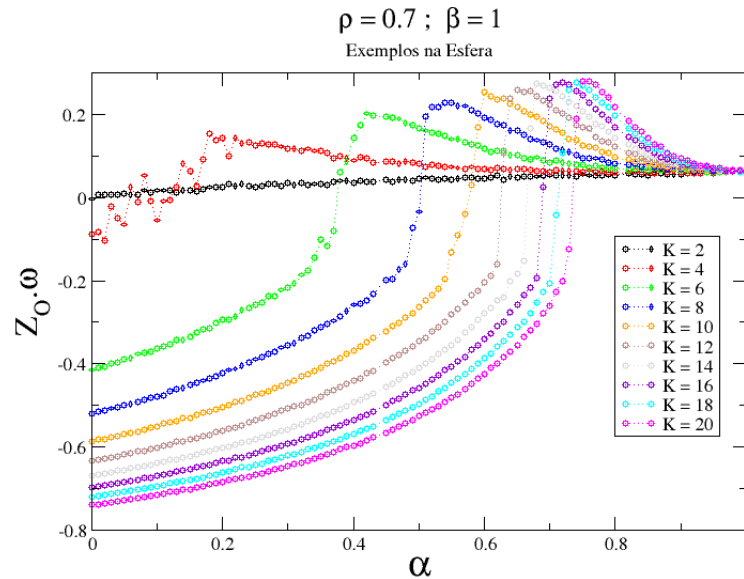


Figura 3.11: Gráfico de m_O x α para diferentes valores para o número médio de vizinhos K .

É possível notar que assim como previsto, o número de vizinhos na rede

tem um efeito semelhante ao da pressão social, onde um valor maior de K implica em uma atitude mais conservadora, ocorrendo também uma transição de fase nesta dimensão. É visível que o estado para $\alpha = 1$ é independente do número de vizinhos, o que era esperado (dado a forma da função Hamiltoniana). Pudemos também comparar o valor de α_c para diferentes valores de K utilizando o mesmo critério da seção anterior para a obtenção de α_c , resultado exibido na figura 8.

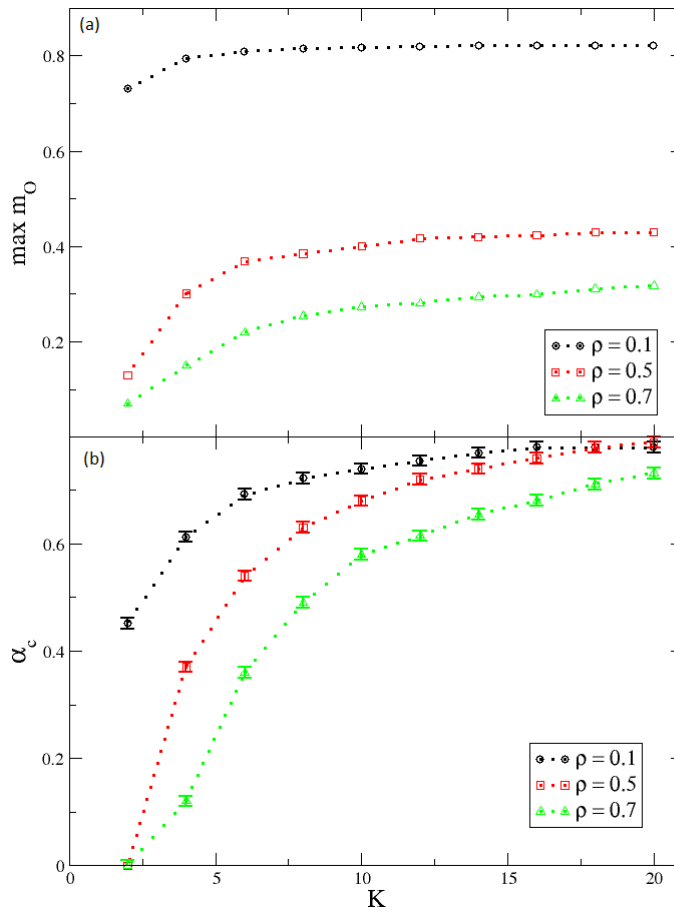


Figura 3.12: Gráfico para três diferentes estilos cognitivos de (a) $\max m_O$ x K (b) α_c x K . É possível notar uma transição de fase na figura (b) para valores mais elevados de ρ , onde α_c passa de 0 para $\neq 0$ com o aumento do número médio de vizinhos. Para este gráfico foram utilizados exemplos na borda, no entanto o comportamento para exemplos uniformemente distribuídos é semelhante.

Esta transição de fase encontrada com relação ao parâmetro K é seme-

lhante à encontrada em [2], na qual foi notado que a escala correta para esta transição é βK , ou seja, o aumento do número de vizinhos é análogo e possui o mesmo efeito que um aumento de pressão social. Notamos que o efeito da rede é progressivamente menor, e que agentes com um estilo cognitivo considerado como liberal possuem uma maior influência do grupo, com uma maior variação nas grandezas medidas ($\max m_O$ e α_c).

3.3 Comparação entre as estratégias

Depois de entender a influência dos parâmetros do modelo nos resultados obtidos, começaremos agora a análise da eficiência das estratégias estudadas. Para tal objetivo, serão feitas comparações diretas das figuras apresentadas nas seções anteriores, de forma a analisar qualitativamente a variação nas grandezas de interesse. Uma análise quantitativa precisa de algumas diferenças não pode ser feita devido à demanda computacional de tal tarefa. Como a estratégia de seleção de exemplos é melhor em uma situação de aprendizado sequencial, é esperado que esta seja a maneira mais eficiente de se convencer a sociedade. Fazendo uma analogia com uma situação real: em uma discussão, é mais eficiente insistir em argumentos que estejam na borda da dúvida da pessoa do que alternar entre assuntos com grande grau de discordância/concordância. Assuntos com certa ambiguidade moral fazem com que o agente em média altere mais sua matriz moral na direção desejada. Os dados utilizados para esta comparação são os mesmos obtidos nas seções anteriores. Foram portanto utilizadas redes do tipo Barabasi-Albert com $n = 400$ agentes e $K = 8$ vizinhos em média (exceto no caso em que K é variável).

Primeiramente foram avaliadas as diferenças entre as curvas de magnetização com relação ao oráculo (m_O). Espera-se que os exemplos na borda apresentem uma exposição crítica (α_c) menor do que uma distribuição uniforme de exemplos, bem como uma maior magnetização máxima. O valor da exposição crítica está relacionado com a eficiência da apresentação de exemplos, enquanto que a variação na magnetização está relacionada com a variância da estratégia de convencimento, que acaba por inserir um grau extra de flutuação no sistema. No entanto, é necessário frisar que mesmo com este grau maior de flutuação, garantimos que a temperatura se mantenha constante por meio do critério de Metrópolis e pela convergência das grandezas amostradas. Os resultados estão apresentados na figura 3.13.

Analisando os resultados obtidos é possível notar que, para o caso em que são selecionados exemplos perpendiculares aos agentes, a magnetização é sempre superior se comparadas aos exemplos distribuídos uniformemente.

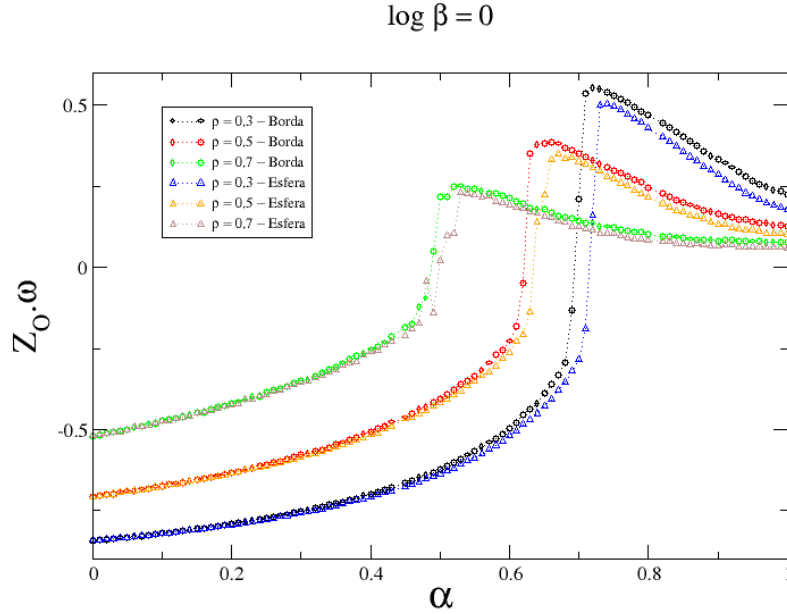


Figura 3.13: Magnetização m_O em função da influência do oráculo α para duas estratégias de convencimento distintas. Os círculos representam exemplos na borda da dúvida, os triângulos os exemplos distribuídos uniformemente. Para este gráfico a pressão social foi mantida constante em $\log \beta = 0$.

Há também uma pequena diferença na ocorrência de α_c , onde para todos os valores do parâmetro de aprendizagem a transição entre as direções de quebra de simetria ocorre antes para os exemplos selecionados, se comparados aos exemplos na hiperesfera. A influência da variância da estratégia de convencimento pode ser observada em $\alpha = 1$, onde a magnetização depende apenas dos parâmetros ρ e β e dos exemplos apresentados. Este mesmo resultado é visto nos diagramas de fase da figura 3.4, onde a magnetização para os exemplos na esfera é sempre inferior (para um mesmo valor dos parâmetros do modelo).

Em seguida foram observadas as curvas de m_O para diferentes pressões sociais. Novamente, espera-se uma maior eficiência da estratégia de seleção de exemplos, que novamente pode ser explicada pelo diagrama de figura 3.4. Mantendo-se o estilo cognitivo fixo, e variando-se a pressão social, nota-se que o mesmo tipo de comportamento ocorre. Novamente, a estratégia na qual são selecionados exemplos ambíguos ($h = 0$), a magnetização é mais elevada e o valor de α_c menor, independente da pressão social.

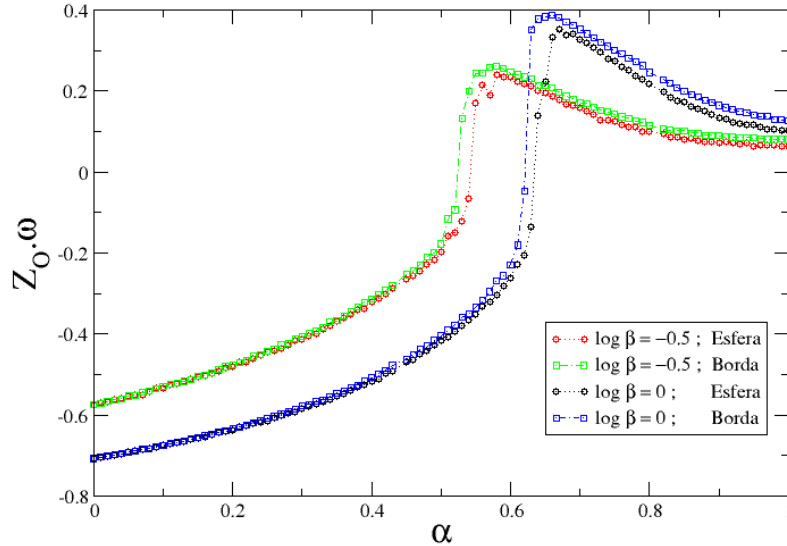


Figura 3.14: Magnetização m_O em função da influência do oráculo α para duas estratégias de convencimento distintas. Os quadrados representam exemplos na borda da dúvida, os círculos os exemplos distribuídos uniformemente. Para este gráfico o estilo cognitivo foi mantido constante em $\rho = 0.5$.

Partimos agora para uma comparação entre os cortes apresentados nas seções 4.1 e 4.2, onde podemos comparar a dependência dos valores de α_c e $\max m_O$ como os parâmetros do modelo. Primeiramente são apresentadas as dependências com a pressão social (β) e o parâmetros de aprendizagem (ρ) (figura 3.15). Em seguida são apresentados os cortes para a dependência com o número médio de vizinhos na rede (K) (figura 3.16).

Fica claro por estes cortes que a estratégia de exemplos na borda da dúvida é sempre ligeiramente superior, apresentando tanto uma menor exposição crítica ao oráculo α_c quanto uma maior magnetização máxima (maior coesão social) após a transição de fase. Podemos ver que a dependência de α_c e $\max m_O$ com os parâmetros ρ e β não é afetada de forma significativa pela mudança de estratégias de convencimento. Podemos supor que o formato do diagrama de fase nos espaços (ρ, β, α_c) e $(\rho, \beta, \max m_O)$ não se alteraria de maneira significativa. Isso ocorre pois a Hamiltoniana depende efetivamente apenas do valor médio e da variância das distribuições de exemplo, agindo portanto como um fator de escala dentro do termo da energia. Esta

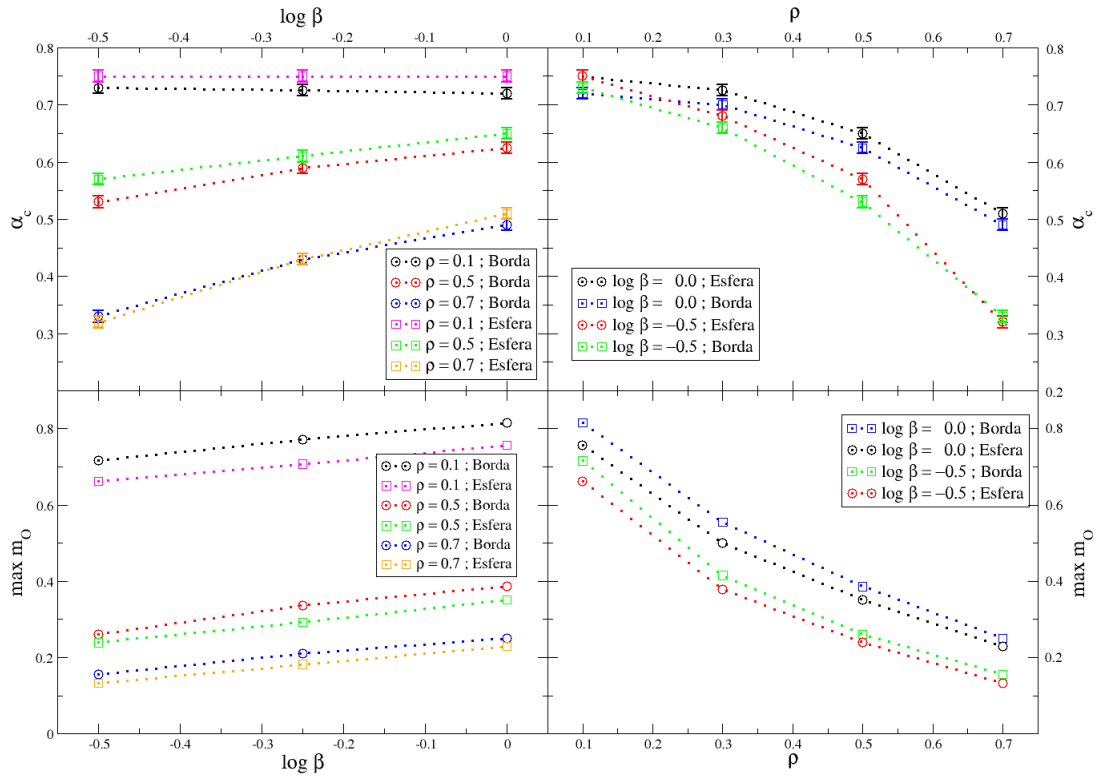


Figura 3.15: Gráficos de (a) α_c x $\log \beta$ (b) α_c x ρ . Os círculos representam a estratégia na qual os exemplos são distribuídos uniformemente, os quadrados representam a estratégia de exemplos perpendiculares aos agentes.

tendência se mantém para a variação do número médio de vizinhos, onde a estratégia de exemplos na borda apresentou o mesmo tipo de superioridade com relação aos exemplos uniformemente distribuídos.

Por último, analisamos a mudança de atitude política das sociedades após a transição de fase. Para isto, simplesmente avaliamos uma medida da importância dada ao conjunto de dimensões. Lembrando dos resultados obtidos por Haidt (seção 2.1), conservadores atribuem o mesmo valor a todas as dimensões morais, enquanto que os liberais valorizam apenas um pequeno número de dimensões. Para realizar tal medida, simplesmente multiplica-se o valor absoluto da magnetização pelo número total de dimensões (5 para estas simulações). Observamos então o valor da atitude política antes de ser exposta ao oráculo ($\alpha = 0$) e em $\operatorname{argmax}_{\alpha} m_O(\alpha_c)$. Pudemos notar que após a transição ocorre uma mudança em direção a atitudes consideradas como mais liberais, associadas com uma menor magnetização. Novamente foi observado que a estratégia na borda da dúvida é ligeiramente superior a de

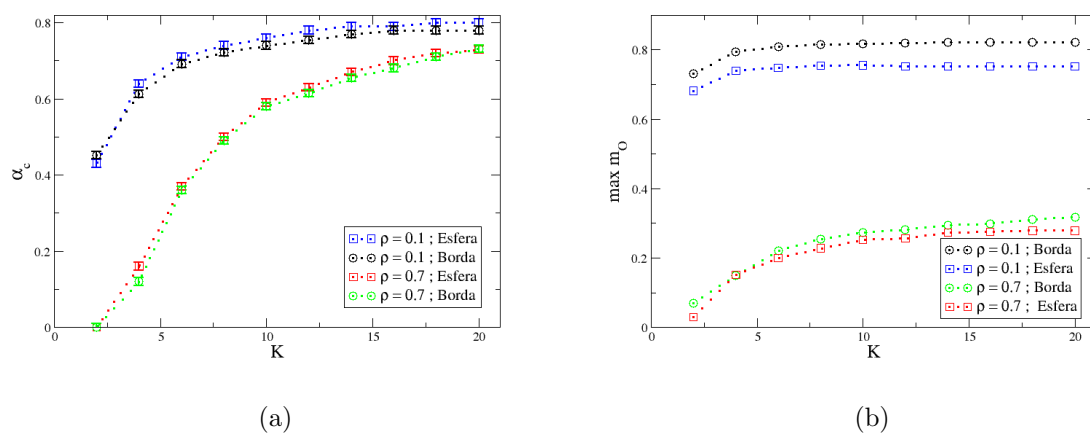


Figura 3.16: Gráficos de (a) $\alpha_c \times K$ (b) $\max m_O \times K$. Para os dois gráficos a pressão social foi mantida constante em $\log \beta = 0$.

exemplos uniformes, necessitando (ou implicando) em uma menor mudança de atitude política após a transição (figura 3.17).

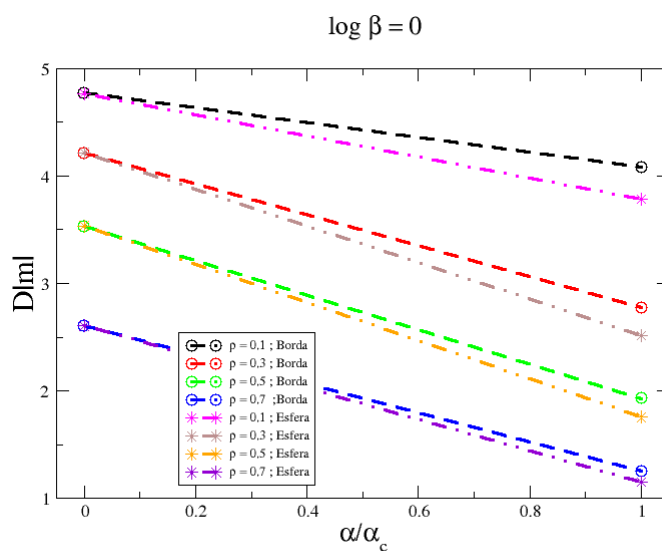


Figura 3.17: Medida do grau de conservadorismo da sociedade (número de dimensões \times magnetização média) para as duas estratégias estudadas.

3.4 Zeitgeist livre

Após analisar a dependência do modelo com os parâmetros no caso em que o Zeitgeist possuía uma direção fixa no espaço, é interessante partirmos para uma situação mais realista: o Zeitgeist livre, ou seja, com a possibilidade de deriva da direção de quebra de simetria. Para as simulações apresentadas nesta seção foram utilizadas novamente redes de Barabasi-Albert com $n = 400$ agentes e $K = 8$ vizinhos em média (com a óbvia exceção do caso em que K é variável) e $N = 5$ dimensões morais.

Para que o vetor de Zeitgeist pudesse se deslocar foi necessário introduzir uma nova definição para \vec{Z} , como sendo a média das matrizes morais em um dado instante. Ou seja;

$$\vec{Z}(t) = \frac{1}{n} \sum_{i=1}^n \vec{\omega}_i(t) = \langle \vec{\omega}_i(t) \rangle . \quad (3.2)$$

Portanto, este vetor deverá ser inicialmente muito próximo de $\vec{Z}(0)$, devendo se alterar conforme forem aceitas flutuações com relação ao estado fundamental. Mantendo-se o vetor oráculo ainda fixo e na mesma direção, i.e $\vec{Z}_O = -\vec{Z}(0)$, espera-se que o vetor de Zeitgeist instantâneo convirja para a mesma direção do oráculo. Esta mudança na direção de \vec{Z} deve ocorrer para qualquer valor de α , contanto que seja fornecido tempo suficiente para a readaptação do sistema. A grandeza interessante a ser analisada neste caso é uma medida de similaridade de $\vec{Z}(t)$ com $\vec{Z}(0)$. Uma medida de similaridade simples e já muito utilizada neste trabalho é a sobreposição entre os dois vetores, definida simplesmente por:

$$R_Z = \vec{Z}(t) \cdot \vec{Z}(0) , \quad (3.3)$$

onde $R = 1$ indica que os vetores possuem a mesma direção e sentido, $R = 0$ indica que estes são perpendiculares, e $R = -1$ indica que estes possuem sentidos opostos. Portanto, podemos acompanhar a evolução do Zeitgeist com relação ao tempo, e inferir sobre a eficiência das estratégias na mudança da direção preferencial, bem como a influência dos parâmetros do modelo na dinâmica de mudança de opinião.

Primeiramente, devemos compreender a influência dos parâmetros ρ , β , K e α na deriva do Zeitgeist. Estes resultados estão apresentados na figura 3.18.

É possível notar que há uma grande influência do parâmetro α , onde em $\alpha = 0.1$ a mudança de opinião é extremamente lenta, necessitando de mais de 100.000 passos de Monte Carlo para a convergência da matriz moral média. No entanto, assim que o valor de α passa para 0.3, o sistema converge

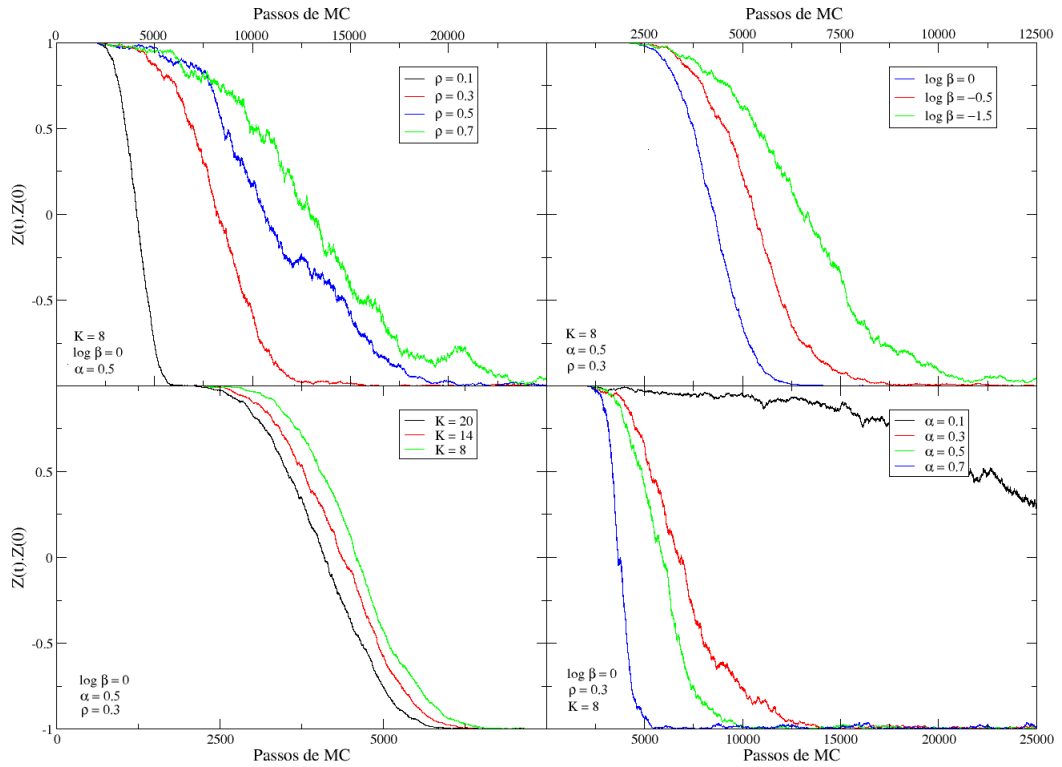


Figura 3.18: Dependência de R_Z com os parâmetros do modelo. Para estes gráficos foi utilizada apenas a estratégia de exemplos na borda, no entanto o comportamento é idêntico para exemplos distribuídos uniformemente.

rapidamente. Não há grande dependência do tempo de convergência com o número médio de vizinhos, sofrendo apenas um leve deslocamento.

Para os parâmetros ρ e β , foram encontrados resultados razoavelmente contra-intuitivos. Segundo o que foi obtido anteriormente, era de se esperar que os agentes associados aos liberais possuíssem um menor tempo de readaptação, porém foi observado que os conservadores convergem consideravelmente mais rápido. Tal fato pode ser compreendido a partir dos resultados para a pressão social, que indicam que a convergência aumenta em conjunto com β . Isto ocorre pois o estado de mínima energia é aquele no qual o vetor de Zeitgeist está alinhado com o oráculo, porém, como os agentes conservadores são mais fortemente correlacionados, uma pequena mudança em um único agente traz consigo toda a sociedade. Ou seja, a sociedade muda de direção de maneira coesa. Tal fato pode ser observado na figura 3.19, que apresenta a magnetização em função do tempo, juntamente com o valor de R_Z .

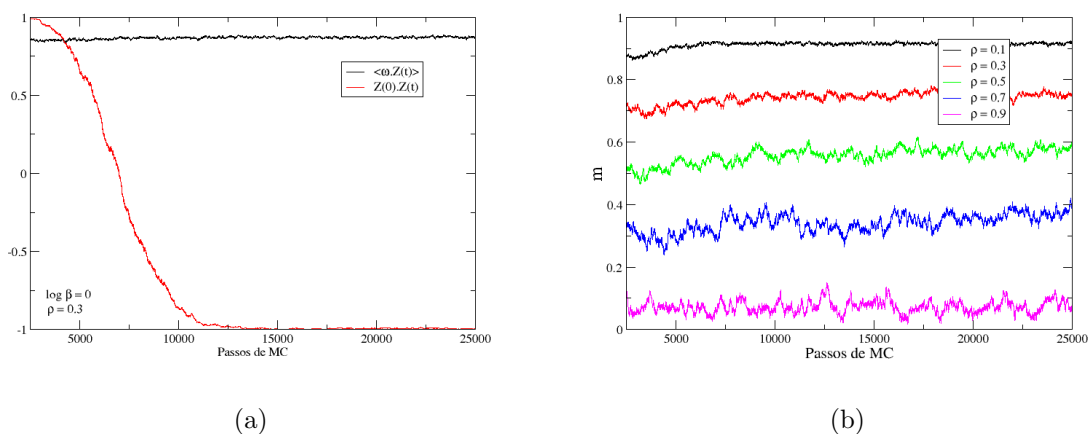


Figura 3.19: (a) Magnetização (curva preta) e sobreposição da matriz moral (curva vermelha) em função do tempo de simulação. (b) Magnetização instantânea para diferentes valores de ρ em função do tempo de simulação.

Na figura 3.19 a curva em preto representa a magnetização do sistema com relação ao Zeigeist instantâneo, e a curva vermelha representa o valor de R_Z . Neste gráfico fica claro como a sociedade altera sua matriz de maneira coesa, a magnetização se mantém praticamente constante (exceto pelas flutuações térmicas) durante o processo de deriva do vetor de Zeigeist. Isto pode ser explicado da seguinte maneira: dado que um agente i foi sorteado para se comunicar com o oráculo, sua energia de interação deverá diminuir na direção de \vec{Z}_O , deslocando levemente o vetor $\vec{Z}(t)$ em direção ao oráculo. Porém, como a magnetização é constante, a sociedade será toda deslocada na mesma direção nos passos subsequentes. Levando em conta que os agentes conservadores possuem uma maior magnetização e uma maior correlação (ρ_{ij}), então este processo de readaptação do vetor de Zeigeist ocorrerá de maneira mais rápida. Podemos notar também que as curvas de R_Z apresentadas na figura 3.18 para agentes mais liberais possuem uma escala maior de flutuação. Isto ocorre devido à maior tolerância de agentes liberais com relação a opiniões divergentes, o que explicaria também a convergência mais lenta. Portanto, podemos afirmar que o parâmetro ρ está relacionado com o grau de acoplamento dos agentes com o vetor de Zeigeist. Na figura 3.19 estão apresentadas as magnetizações com relação ao Zeigeist instantâneo para diversos valores do parâmetro de aprendizagem, corroborando com o que foi dito anteriormente.

Tendo compreendido a influência dos parâmetros na dinâmica de rea-

daptação, bem como o comportamento do vetor de *Zeitgeist* instantâneo, podemos então partir para a comparação entre as estratégias de convencimento. Para a obtenção destes resultados foi utilizada uma rede do tipo Barabasi-Albert com $n = 400$ agentes e $K = 8$ vizinhos em média, a pressão social foi mantida constante em $\beta = 1$ e o valor de exposição ao oráculo em $\alpha = 0.5$ de maneira a garantir uma rápida convergência da matriz moral média. Os resultados estão apresentados na figura 3.20.

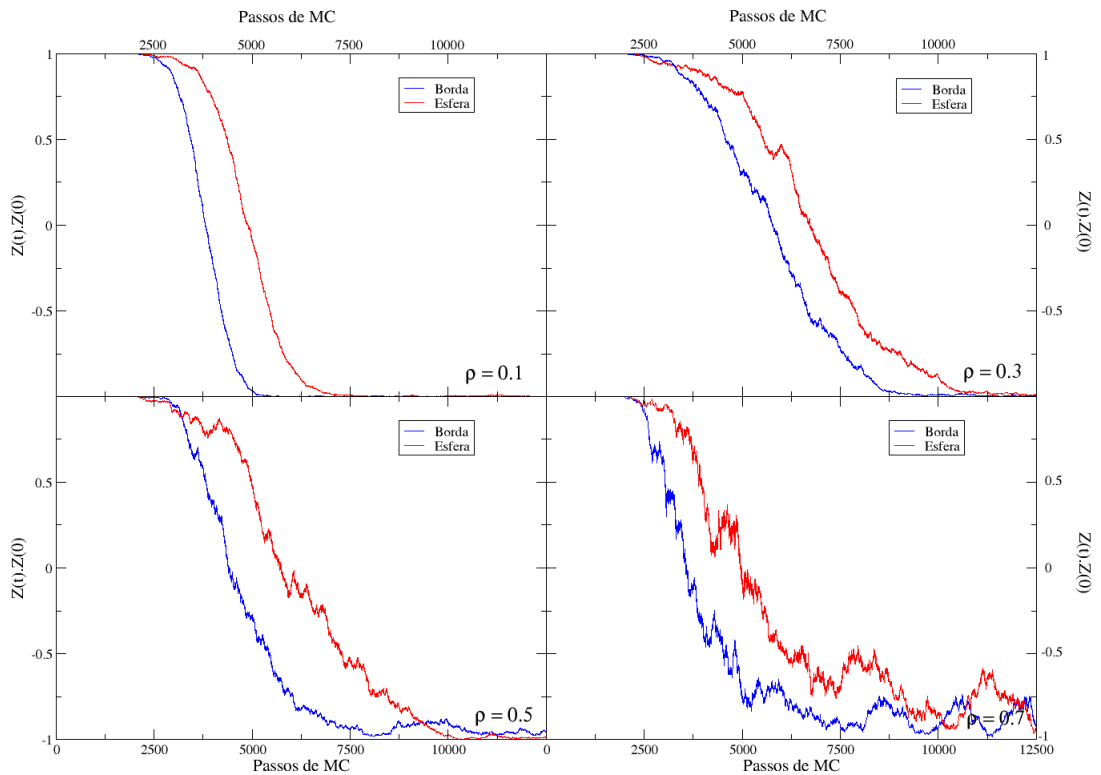


Figura 3.20: Comparação entre a convergência das estratégias para 4 valores distintos do parâmetro de aprendizagem ($\rho = 0,1 ; 0,3 ; 0,5$ e $0,7$).

Podemos notar que a estratégia na qual os exemplos estão sempre localizados na borda da dúvida ($h_i = 0$) se apresenta superior em todos os casos. A diferença do tempo de convergência é em torno de 2500 passos de Monte Carlo para todos os valores de ρ . Novamente, a diferença entre as duas estratégias não é tão considerável, porém podemos afirmar ser o suficiente para justificar a utilização de exemplos na borda.

Conclusão

O convencimento de populações é um problema de grande impacto social. Estudar técnicas e estratégias que facilitem mudanças de opinião podem oferecer possibilidades de solução de conflitos, bem como técnicas de manipulação de massas. Nosso estudo visa avaliar o funcionamento e a influência das estratégias de convencimento/diálogo em casos de mudanças coletivas. Partindo de um sistema minimalista, pudemos começar a análise de como o perfil cognitivo, a pressão social e a estrutura da sociedade influenciam na dinâmica de mudanças coletivas.

Inicialmente, foi possível notar que a inserção de uma segunda direção de quebra de simetria não alterou significativamente o comportamento observado em trabalhos anteriores [1, 2, 34]. Pudemos observar que as sociedades possuem uma faixa de existência dentro do espaço de fase, na qual regiões de ultra-liberalidade e de baixa pressão social não correspondem a sociedades organizadas, e portanto são insustentáveis. Esta mesma categorização das sociedades modelo estudadas se mantém mesmo quando há uma influência externa. No entanto, a existência de um agente de inserção de informação (o oráculo), implica em uma maior flutuação das matrizes morais. Como é possível notar na figura 3.21. Estes resultados foram reforçados no histograma apresentado na seção 4.1, no qual é possível notar uma maior dispersão das sobreposições intra-sociais (figura 3.22).

Com ressalvas, podemos fazer analogias para este resultado, utilizando como exemplo a influência da mídia, que ao fornecer seu ponto de vista/opinião e enquadrar (*framing* em inglês) o assunto de diferentes maneiras, tende a inserir certo grau de variação para opiniões reais [41–44]. No caso modelo estudado, há uma transição de fase associada a exposição à este agente externo. Esta mudança de direção de quebra de simetria é dependente do estilo cognitivo da sociedade. Agentes associados com liberais ($\rho \approx 1$) possuem uma maior influência do oráculo no caso em que o *Zeitgeist* é fixo. Isto ocorre pois sua maior tolerância a flutuações faz com que a sociedade passe a perceber o termo de interação com o oráculo para menores exposições. Podemos entender este comportamento ao notar que a diminuição progressiva da mag-

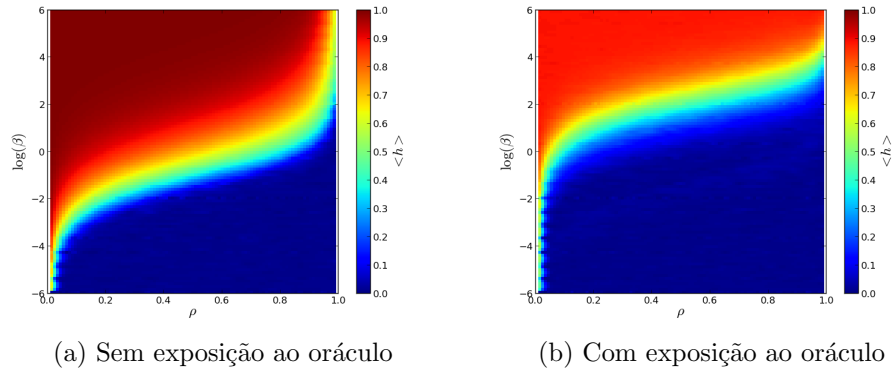


Figura 3.21: Diagramas de fase obtidos na seção 4.1, ilustrando as diferenças nos casos com/sem oráculo.

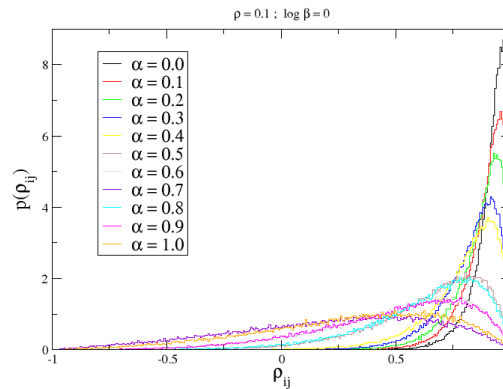


Figura 3.22: Histograma apresentado na seção 4.1, no qual é possível observar o aumento na dispersão de opiniões com o acréscimo da influência do oráculo

netização, indicando que a sociedade deve se abrir (aumentar sua dispersão moral) a fim de se reorganizar em torno da nova opinião preferencial. Como os agentes liberais já possuem uma maior dispersão, esta transição se torna mais fácil.

Ao analisar os efeitos da pressão social na dinâmica, notamos que há um aumento do conservadorismo da sociedade. Este resultado já havia sido observado no estudo realizado por J. César [34]. No entanto, foi possível notar que este deslocamento em direção a atitudes mais conservadoras eleva a exposição necessária para que a sociedade passe a utilizar o oráculo como direção preferencial, como é possível notar nas curvas apresentadas na figura

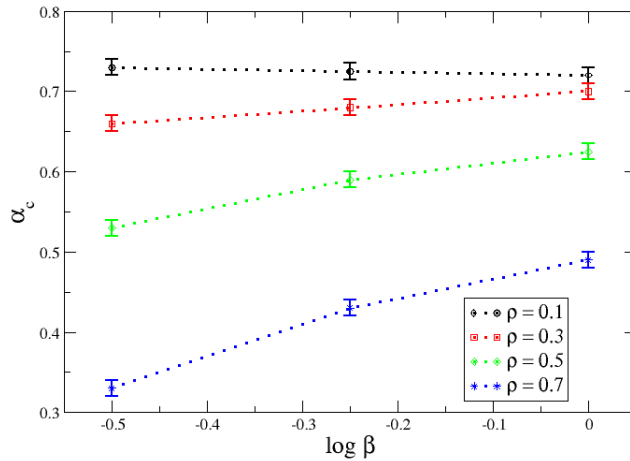


Figura 3.23: Gráfico da exposição crítica contra a pressão social, indicando que um maior valor de β dificulta a transição em direção ao oráculo, elevando o valor de α_c .

3.23. Neste mesmo gráfico (figura 3.23) é possível notar também que agentes com perfil explorador sofrem maior influência do aumento de pressão social. Tal comportamento pode ser entendido ao se considerar que o forte acolhimento de agentes corroborativos e sua baixa dispersão geram um efeito semelhante ao aumento de β . O mesmo tipo de comportamento ocorreu com o aumento do número médio de vizinhos. Então, como era de se esperar, as características fundamentais do modelo de mantém para este caso geral. Pudemos também notar uma transição de fase com relação ao parâmetro K , onde para sociedades com poucos vizinhos a mudança de opinião em direção ao oráculo não ocorre, ou fica extremamente atenuada (figura 3.24). O que indica que também é necessário um certo grau de coesão social/troca de informação para que esta mudança ocorra. Isto ocorre devido ao formato da energia, que é simétrico com relação a troca de sinal. Portanto, para um certo valor de α , o mínimo da energia de interação se desloca para a direção do oráculo, e a troca de informação entre os agentes passa a reforçar esta mudança. Há portanto um valor ótimo (ou crítico) de α , onde ocorre um máximo da magnetização na direção do oráculo. Este comportamento pode ser observado no histograma da figura 3.22, no qual o comportamento para $\alpha \approx \alpha_c$ (figura 3.6b). apresenta um leve aumento no consenso da sociedade, corroborando com a existência de um valor ótimo para a exposição ao agente externo.

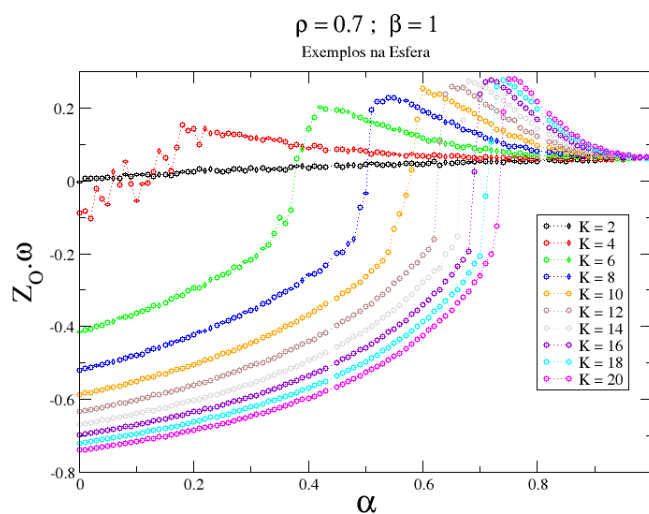


Figura 3.24: Figura obtida na seção 4.2, indicando a diminuição da magnetização para sociedades com um menor número médio de vizinhos.

Na segunda etapa do trabalho, analisamos as diferenças qualitativas entre as estratégias de convencimento empregadas. Analisando os resultados para as magnetizações, foi possível compreender que há uma influência da variância das estratégias na coesão da sociedade. Por possuir uma maior carga informacional a cada passo, a estratégia que apresenta exemplos na região denominada *borda da dúvida* apresentou uma maior eficiência na realização de mudanças coletivas. Foi observada uma menor necessidade de exposição ao oráculo, bem como as magnetizações se mantiveram em valores superiores após a transição. A interpretação dada para estes resultados pode estar relacionada com a mudança de atitude observada em pessoas após a utilização de táticas de persuasão [40], indicando que há um deslocamento em direção à atitudes mais liberais após a transição. Espera-se portanto que a estratégia de seleção de exemplos implique e/ou necessite de uma menor mudança de atitude política (figura 3.25), o que em casos reais seria uma vantagem.

Por último, foi realizado um estudo considerando o *Zeitgeist* livre. Este caso nos forneceu informações relevantes sobre a eficiência das estratégias, bem como alguns resultados contra-intuitivos. Fugindo do que havia sido observado anteriormente, os agentes considerados como conservadores possuíram uma maior facilidade de alteração. A conclusão tomada com relação a estes resultados foi de que agentes com uma estratégia cognitiva corroborativa possuem um maior acoplamento com relação ao vetor de *Zeitgeist*, indicando

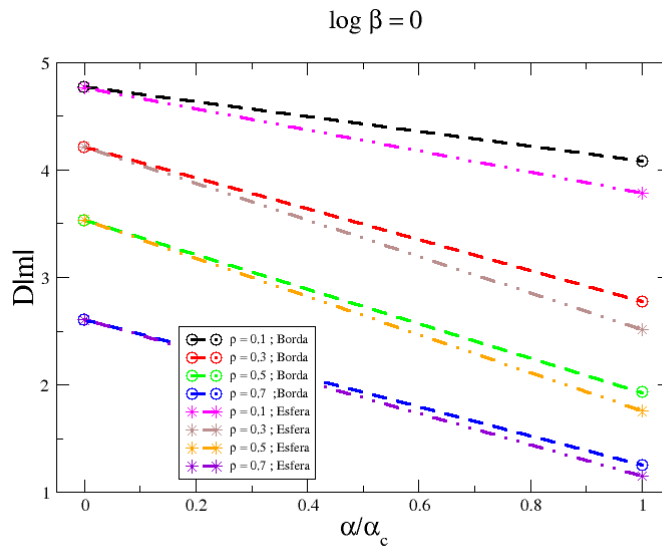


Figura 3.25: Medida de conservadorismo da sociedade, indicando uma mudança em direção a atitudes mais liberais (menor utilização das dimensões morais) após a transição.

que a sociedade se move de maneira coesa em direção ao vetor Oráculo, como pode ser observado no comportamento da magnetização, constante mesmo com a alteração da orientação moral média dos agentes (figura 3.26). Em um caso real, o vetor de Zeitgeist livre corresponderia a uma mudança gradual de opinião geral da sociedade, juntamente com o agente externo. Portanto, parece plausível que em algumas situações a coesão social causada pelo conservadorismo sirva de agente mutante. Estes resultados foram observados utilizando-se a magnetização da sociedade e da dependência da orientação relativa da opinião média com os parâmetros do modelo, que nos forneceu um indicativo da importância da coesão social para este tipo de mudança coletiva. É possível notar que o comportamento de agentes que valorizam a novidade acaba por permitir uma maior flutuação do vetor de Zeitgeist instantâneo, portanto quando alguma mudança em direção ao Oráculo ocorre em um dos agentes ela é atenuada pelas flutuações naturais da sociedade. Podemos notar que a pressão social exerce o mesmo efeito que a diminuição do valor de ρ (perfil corroborativo), assim como ocorria com o Zeitgeist fixo, corroborando assim a importância da coesão social como fator de mudanças em uma situação de Zeitgeist livre. Podemos afirmar que o Zeitgeist livre, apesar de não alterar o comportamento qualitativo da sociedade bem como sua natureza organizacional, altera de maneira significativa a dinâmica de

mudança de opiniões na sociedade. De tal forma, podemos compreender que tal comportamento deva ocorrer de maneira distinta em situações empíricas. É necessário portanto, refletir sobre as possíveis consequências e situações na qual cada caso ocorreria, para que possamos inferir sobre a validade deste modelo quando confrontado com resultados empíricos.

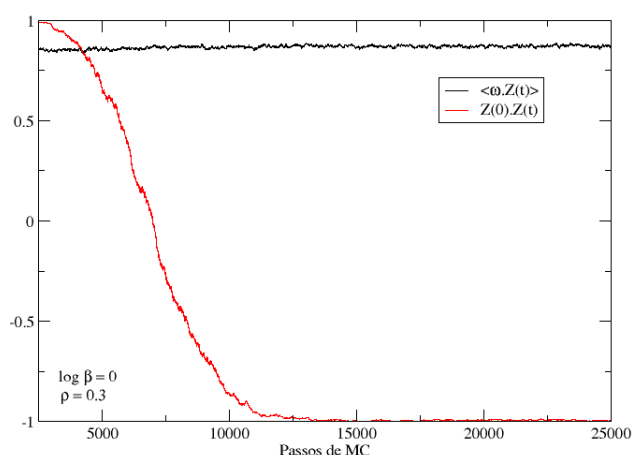


Figura 3.26: Figura apresentada na seção 4.4, ilustrando o comportamento constante da magnetização (em preto), mesmo em conjunto com a deriva da opinião média (em vermelho).

Após compreender a influência dos parâmetros de aprendizagem e pressão social na deriva do Zeitgeist, estudamos a influência das estratégias. Conforme o esperado, a estratégia de convencimento com seleção de exemplos se apresentou mais eficiente para todas situações (figura 3.27). Esta eficiência é traduzida em um menor tempo de readaptação, observado na figura 3.20. É possível notar que o comportamento para diferentes ρ 's não se altera com as estratégias (conforme apresentado na figura 3.20), possuindo um estado final muito semelhante e com um grau de flutuação próximo para ambas as estratégias de convencimento, no entanto o estado transiente é de menor duração no caso dos exemplos na borda da dúvida. Podemos então concluir que mesmo em uma situação com diferenças na dinâmica de mudanças de opinião, como ocorre no caso em que o vetor de Zeitgeist é livre, a estratégia de convencimento com exemplos na borda da dúvida possui maior eficiência.

Após a análise cautelosa do comportamento do modelo, pudemos observar que algumas características fundamentais permanecem em situações mais

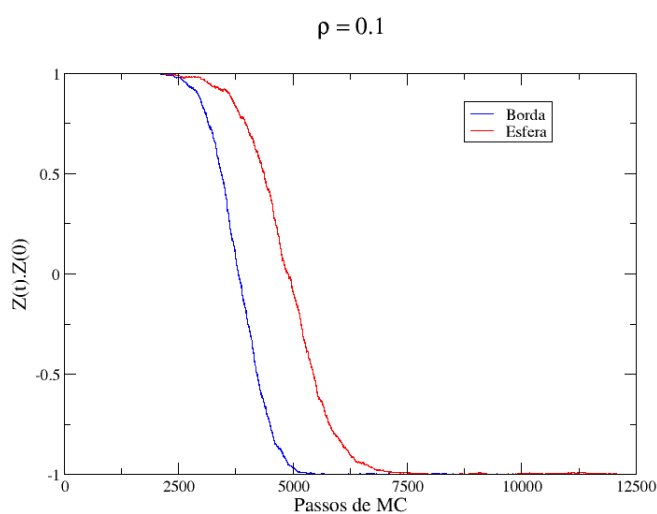


Figura 3.27: Comparação entre as estratégias de convencimento no caso em que o vetor de Zeitgeist corresponde a opinião média da sociedade.

gerais. Os diagramas de fase obtidos quando não há a existência de um vetor oráculo permanecem mesmo quando ocorre a influência de um agente externo, alterando apenas o valor da magnetização, porém mantendo seu comportamento geral. Podemos também notar que a existência das diferentes estratégias altera o grau de flutuação das sociedades, bem como sua velocidade de adaptação em situações de Zeitgeist livre, porém não exerce influência qualitativa no comportamento do modelo. É possível afirmar que apesar de distintos, ambos os casos estudados (Zeitgeist livre/fixo) apresentaram uma maior eficiência da estratégia com exemplos na borda da dúvida, portanto, em estudos futuros, é esperado que este tipo de comportamento também se mantenha. Esta estratégia de convencimento já havia sido cogitada por teóricos, sendo chamado de *processo de correção flexível*, como pode ser observado na seguinte frase “Da mesma maneira que a confiança em pensamentos aumenta seu grau de confiabilidade, o aumento da dúvida faz com que os indivíduos descartem suas ideias. Algumas vezes, os indivíduos podem criar tantas dúvidas de seus pensamentos que podem passar a acreditar que o oposto é verdade” [45], corroborando com a hipótese de que exemplos na borda da dúvida possuem maior potencial para alterar as matrizes morais de indivíduos. Outro estudo, liderado por Phillip Fernbach sugere que extremismos surgem de uma falsa sensação de entendimento do assunto [46], sugerindo portanto que uma forma eficiente de convencimento se dá no enfraquecimento de tais convicções, ou seja, é necessário minar as certezas em uma

certa direção de pensamento, enquadrando assim a estratégia de exemplos na borda da dúvida. Fernbach afirma também que a exposição a ideias que ferem as certezas dos indivíduos levam a atitudes mais moderadas, reforçando assim os resultados encontrados por nossas simulações.

Pudemos então notar que modelos consideravelmente simples são capazes de reproduzir características de sistemas complexos como a dinâmica de opiniões e a composição moral. Utilizando ferramentas de ramos distintos do conhecimento, pudemos desenvolver e estudar modelos matemáticos que possibilitam o vislumbre de tendências em situações de grau de complexidade mais elevado, fornecendo assim potenciais caminhos e estudos posteriores a serem realizados. Mesmo com uma análise com caráter mais qualitativo, é possível afirmar que este trabalho vai no mesmo sentido que muitos estudos de psicologia social, bem como abre a possibilidade de novos estudos empíricos, possivelmente facilitando o entendimento de certos processos por trás da tomada de decisões morais e de mudança de opinião. Futuramente esperamos estender ainda mais o modelo, abrindo a possibilidade de conectá-lo com outros estudos sociológicos e que, com certo grau de otimismo, esperamos que nos forneçam evidências de como se dá a estruturação da natureza humana de maneira mais quantitativa.

Apêndice A

Informação e Maximização da Entropia

Foi necessário um avanço no entendimento de algumas teorias para que pudesse ser realizada a conexão entre as ciências ditas exatas e humanas. Um destes avanços foi a compreensão e a interpretação das probabilidades como a ferramenta correta para se lidar com crenças sobre asserções em situações de informação incompleta. Iremos tratar rapidamente sobre o tema neste apêndice, para então aplicar a técnica de maximização da entropia a fim de mostrar que a distribuição de Boltzmann é a escolha mais sensata para descrever os estados de nosso sistema.

A teoria da probabilidade é um assunto de interesse razoavelmente antigo, porém só no último século passou a receber a devida atenção. Seu surgimento, no século XVII, possuía como principal intuito facilitar vitórias e aprimorar estratégias para jogos de azar. Tal aplicação fez com que a teoria fosse mal vista, e ganhasse ares de conhecimento de rua. Grandes matemáticos da época como Pierre de Fermat e Blaise Pascal já se aventuravam na área, porém sem exposição pública. Fermat era jurista, e portanto tinha receio de ser considerado um vigarista, ou de que estaria envolvido com jogos de azar. Alguns séculos depois, Laplace voltou a se voltar para as probabilidades como forma de melhor inferir sobre valores de parâmetros obtidos a partir de um conjunto de dados. O problema tratado por Laplace era complicado, determinar a massa de Saturno com base nas órbitas dos planetas de sistema solar, ou seja, uma medida indireta. Em sua busca por respostas, acabou por desenvolver o que seria a primeira formalização do chamado teorema de Bayes, uma espécie de regra para a atualização de probabilidades. No entanto, apenas no século XX com os trabalhos de Kolmogorov, Fisher e outros, que a estatística passou a ser formalmente considerada como uma área do conhecimento. Porém, apenas com o seminal trabalho de Claude

Shannon [37], sobre a capacidade de transmissão de canais via rádio, é que surgiram as primeiras interpretações sobre o significado das probabilidades. Em seu artigo de 1948 intitulado "A teoria matemática da comunicação", Shannon procurou responder qual a capacidade máxima de informação que poderia ser transmitida em um canal ruidoso. A resposta para esta questão é famosa equação para a entropia:

$$S = - \sum_i p_i \log p_i$$

Onde pela primeira vez, probabilidades recebiam uma interpretação relacionada a informação. Anos mais tarde, com o advento dos axiomas de Cox e Jaynes [38], foi possível afirmar que probabilidades são uma ferramenta para se descrever o grau de confiança em uma asserção em situações de informação incompleta, utilizando apenas um conjunto mínimo de desejos matemáticos (ou *desideratas*). Atualmente, Ariel Caticha demonstrou o poder da teoria de probabilidades, ao derivar a equação de Schrödinger como um caso específico da chamada inferência entrópica [39]. Iremos agora rapidamente derivar a equação da entropia e sua interpretação.

Suponha um conjunto de alternativas mutuamente exclusivas i . O estado do sistema é desconhecido e a informação \mathcal{I} é incompleta. Utilizando o conceito de probabilidade condicional e a interpretação de Cox/Jaynes podemos atribuir uma probabilidade tal que $p(i|\mathcal{I}) = p_i$. Para podermos atribuir qual a alternativa i é a correta com total certeza ($p_i = 1$) seria necessário acrescentar mais informação ao sistema. A pergunta a ser feita é: quanta informação?

Para responder esta pergunta, utilizaremos o desenvolvimento apresentado por Ariel Caticha em seu livro "Entropic inference and the foundations of physics".

Considere agora um conjunto discreto de n variáveis mutuamente exclusivas descritas por alternativas exaustivas i , cada qual com sua respectiva probabilidade p_i . Precisamos então de uma função S que nos fornece a quantidade faltante de informação. Esta função deve obedecer à três axiomas propostos por Shannon:

- **Axioma 1:** S deve ser uma função contínua das probabilidades p_i .
- **Axioma 2:** Se p_i for igual para todos os valores de i , e portanto $p_i = 1/n$, então S deve ser uma função crescente com relação a n . Ou seja, $S = S(1/n, \dots, 1/n) = F(n)$ onde $F(n)$ é crescente.
- **Axioma 3:** Para todos os possíveis agrupamentos $g = 1, \dots, N$ dos estados $i = 1, \dots, n$, tal que $g \in G$, temos que a função S deve satisfazer

a relação

$$S[p] = S_G[p] + \sum_g P_g S_g[p_{\cdot|g}]$$

Este terceiro axioma é conhecido como *propriedade de agrupamento* e é facilmente demonstrado.

Podemos agora partir para solução destes axiomas. Vamos primeiramente supor que todos os estados i são igualmente prováveis, de forma a utilizar o axioma 2. Portanto $p_i = 1/n$. Suponhamos então que todos os N grupos g possuem uma mesma quantia m de estados, sendo $m = n/N$. Portanto, $P_g = 1/N$ é a probabilidade de se escolher um grupo e $p_{i|g} = p_i/P_g = 1/m$ é a probabilidade de se escolher um estado i dado que se escolheu um grupo g . Pelo axioma 2, temos que:

$$S[p_i] = S(1/n, \dots, 1/n) = F(n) \quad (\text{A.1})$$

$$S[P_G] = S(1/N, \dots, 1/N) = F(N) \quad (\text{A.2})$$

$$S[p_{i|g}] = S(1/m, \dots, 1/m) = F(m) \quad (\text{A.3})$$

Utilizando a propriedade de agrupamento do axioma 3, temos que:

$$F(mN) = F(m) + F(N) \quad (\text{A.4})$$

Como nenhuma restrição foi imposta aos valores de m e N no desenvolvimento dos axiomas, esta relação deve valer para qualquer valor inteiro de m e N . É fácil notar que

$$F = k \log n \quad (\text{A.5})$$

é solução para esta equação. No entanto, para garantir a unicidade desta solução precisamos também impor que $F(m)$ seja monotonicamente crescente com relação a m [38]. Consideremos então dois s e t , ambos maiores que 1. É possível aproximar a razão de seus logaritmos de maneira arbitrária, ou seja, é possível encontrar dois números inteiros α e β tais que:

$$\frac{\alpha}{\beta} \leq \frac{\log s}{\log t} < \frac{\alpha + 1}{\beta} \quad (\text{A.6})$$

ou

$$t^\alpha \leq s^\beta < t^{\alpha+1} \quad (\text{A.7})$$

Como F é monotonicamente crescente, a relação se mantém em F , ficando

$$F(t^\alpha) \leq F(s^\beta) < F(t^{\alpha+1}) \quad (\text{A.8})$$

Como $F(a.b) = F(a) + F(b)$, então $F(a^c) = cF(a)$. Portanto

$$\alpha F(t) \leq \beta F(s) < (\alpha + 1)F(t) \quad (\text{A.9})$$

ou

$$\frac{\alpha}{\beta} \leq \frac{F(s)}{F(t)} < \frac{\alpha + 1}{\beta} \quad (\text{A.10})$$

Então $F(s)$ e $F(t)$ devem se aproximar pela mesma razão que $\log s$ e $\log t$. Comparando então estas duas razões, devemos ter:

$$\left| \frac{F(s)}{F(t)} - \frac{\log s}{\log t} \right| = \left| \frac{F(s)}{\log s} - \frac{F(t)}{\log t} \right| \leq \frac{F(t)}{\beta \log s} \quad (\text{A.11})$$

O lado direito da inequação pode ser feito arbitrariamente pequeno utilizando-se um β arbitrariamente grande (já que seu tamanho não foi limitado no começo da dedução). Portanto, $F(s)/\log s$ é deve ser uma constante. Ou seja $F(s) = \log s$, o que garante a unicidade da solução apresentada.

Vamos agora supor que os grupos g não são uniformemente divididos, portanto possuem um tamanho m_g . Então $P_g = m_g/n$ e $p_{i|g} = 1/m_g$. Substituindo no axioma 3 temos

$$S[p] = S_G[P] + \sum_g P_g S[p_{\cdot|g}] \quad (\text{A.12})$$

ou

$$F(n) = S_G[P] + \sum_g P_g F(m_g) \quad (\text{A.13})$$

então

$$S_G[P] = F(n) - \sum_g P_g F(m_g) = \sum_g P_g [F(n) - F(m_g)] \quad (\text{A.14})$$

Substituindo a solução $F(x) = k \log x$, temos:

$$S_G[P] = k \sum_g P_g [\log n - \log m_g] = -k \sum_g P_g \log \frac{m_g}{n} \quad (\text{A.15})$$

Que é simplesmente

$$S_G[P] = -k \sum_{g=1}^N P_g \log P_g \quad (\text{A.16})$$

Portanto, substituindo $S_G[P]$ na equação A.12, obtemos a quantidade de informação faltando sobre a distribuição das probabilidades p_i, \dots, p_n , a entropia de Shannon:

$$S[p] = -k \sum_{i=1}^n p_i \log p_i \quad (\text{A.17})$$

Onde k é uma constante arbitrária que define a escala/unidade da entropia. Como $0 \leq p_i \leq 1$, então $S[p]$ será sempre positiva.

Sabemos portanto que a entropia representa de certa maneira a ignorância que possuímos sobre um certo conjunto de distribuições de probabilidades. Partindo desta informação, devemos nos perguntar qual a melhor escolha para uma distribuição desconhecida de um sistema/modelo. A resposta é vem pelo princípio da mínima informação discriminatória: devemos optar por aquela que nos leva ao menor número de suposições, i.e a que possui máxima entropia. Desta maneira evitaremos informações falsas e/ou duvidosas, inserindo em nossa distribuição apenas o conjunto de informações conhecidas. Este tipo de abordagem é conhecida como Máxima Entropia (MaxEnt). Entretanto, devemos levar em consideração que a maximização da entropia é sempre relativa, ou seja, devemos maximizar a entropia com relação à uma distribuição de referência. Devemos portanto maximizar a entropia relativa:

$$S[p_i||q_i] = - \sum_i p_i \log \frac{p_i}{q_i} \quad (\text{A.18})$$

Esta entropia relativa é conhecida também como *divergência de Kullback-Leibler* (porém com o sinal invertido), e pode ser interpretada como uma distância relativa no espaço paramétrico das distribuições.

No entanto ainda resta uma pergunta, como inserir as informações relevantes? O máximo que podemos fazer nesta situação é considerar os vínculos impostos ao problema (que usualmente é a única informação que possuímos). Levando em consideração que a solução para a máxima entropia no caso sem vínculos é a distribuição uniforme. Por ser o caso mais simples, a distribuição uniforme costuma ser a referência para a projeção no espaço paramétrico, portanto assumindo o lugar de q_i . Para a maioria dos problemas envolvidos, este vínculo surge como a imposição de que uma certa função arbitrária mantenha um valor médio constante. No caso dos problemas físicos: energia, potencial químico, etc. No caso mais geral queremos que $\langle f^k \rangle = \sum_i p_i f_i^k = F^k$,

onde temos k funções f com seus respectivos valores médios. Podemos então utilizar o método dos multiplicadores de Lagrange para maximizar a entropia sujeita aos vínculos. Dando origem a:

$$\frac{\delta}{\delta p_i} \left[S[p] - \alpha \left(\sum_i p_i - 1 \right) - \lambda^k \left(\sum_i p_i f_i^k - F^k \right) \right] = 0 \quad (\text{A.19})$$

$$0 = \log p_i + 1 + \alpha + \lambda^k f_i^k \quad (\text{A.20})$$

Renomeando $1 + \alpha = \lambda_0$, temos:

$$p_i = \exp[-\lambda_0 - \lambda^k f_i^k] \quad (\text{A.21})$$

Notando que a constante λ_0 surge do vínculo de normalização das distribuições, temos que $e^{\lambda_0} = Z$, onde Z é a função de partição. Ficamos assim com:

$$p_i = \frac{1}{Z} e^{-\lambda^k f_i^k} \quad (\text{A.22})$$

A distribuição de probabilidades p_i é conhecida como *distribuição canônica*. No caso em que as funções f^k são grandezas físicas, a distribuição é conhecida por *distribuição de Boltzmann*.

Para nosso modelo $f_i^k = \mathcal{H}(h_{ij}, h_{iO})$, $F^k = E$ e portanto $\lambda^k = \beta$, com $k = 1$.

Especificações Computacionais

Todas as simulações foram realizadas em um computador com processador Intel i7[®] 3.2GHz 8 core, 32Gb de Memória DDR3.

Foram realizadas simulações utilizando todos os núcleos de processamento para a variação dos parâmetro. O tempo computacional médio para a obtenção de uma figura como a da 3.7 foi de uma semana, e fora utilizados em torno de 4Gb de memória por simulação. A quantidade de dados gerada se aproxima a 18Gb. Foi necessário ao todo em torno de 1 mês de simulação para a obtenção de todos os resultados apresentados.

O código foi desenvolvido na linguagem C, utilizando como base o código de J. César para seu desenvolvimento, e foi compilado utilizando o compilador g++.4.4. O código será posteriormente fornecido em repositórios online para download e poderá ser utilizado livremente.

Referências Bibliográficas

- [1] N. Caticha, R. Vicente; *Agent-based Social Psychology: from Neurocognitive Processes to Social Data*. Advances in Complex Systems 14(5):711–731, 2011.
- [2] R. Vicente, A. Susemihl, J. P. Jericó, N. Caticha; *Moral foundations in an interacting neural networks society*. Physica A 400: 124–138, 2014.
- [3] R. Axelrod; *The Dissemination of Culture: A Model With Global Convergence and Local Polarization*. J. of Conf. Res. 41(2): 203-226, 1997.
- [4] A. H. Rodriguez e Y. Moreno; *Effects of mass media action on the Axelrod model with social influence*. Phys. Rev. E Stat. Nonlin. Soft. Matter. Phys 82(1 Pt 2), Epub 2010 Jul 21.
- [5] K. Klemm, V. M. Eguíluz, R. Toral e M. San Miguel; *Role of dimensionality in Axelrod's model for the dissemination of culture*. Phys. A 327:1-5, 2011.
- [6] T. E. Forland; *Mentality as a social emergent: can the zeitgeist have explanatory power*. History and Theory, 47: 44–56, 2008.
- [7] J. Haidt e J. Graham; *Planet of the Durkheimians, where community, authority and sacredness are foundations of morality*. Social and psychological bases of ideology, Cap. 15, p. 371–401, Oxford University Press, 2009.
- [8] J. Haidt e C. Joseph; *Intuitive ethics: how innately prepared intuitions generate culturally variable virtues*. Daedalus 133 (4): 55–66, 2004.
- [9] J. Haidt e T. Wheatley; *Hypnotic disgust makes moral judgments more severe*. Psychol. Sci. 16(10): 780–784, 2005.
- [10] J. Haidt; *The emotional dog and its rational tail: A social intuitionist approach to moral judgment*. Psychological Review 108 (4): 814–834, 2001.

- [11] J. Haidt e C. Joseph; *The moral mind: How five sets of innate intuitions guide the development of many culture-specific virtues, and perhaps even modules*. The Innate Mind, Vol. 3, Cap. 19, P. Carruthers, S. Laurence and S. Stich, 2006.
- [12] J. Haidt e J. Graham ; *When morality opposes justice: Conservatives have moral intuitions that liberals may not recognize*, Social Justice Research 20: 98–116, Springer, 2007.
- [13] J. Haidt; *The new synthesis in Moral Psychology*. Science, 316: 998–1001, 2007.
- [14] J. Cesar e N. Caticha; *For whom will the bayesians vote?*, pre-print, 2013.
- [15] N. Eisenberger, M. Lieberman, and K. Williams; *Does rejection hurt? an fmri study of social exclusion*. Science, 302: 290–292, 2003.
- [16] P. R. Nail e I. McGregor; *Conservative shift among liberals and conservatives following 9/11/01*. Social Justice Research 22(2-3): 231–240, 2009.
- [17] M. Sherif; *And experimental approach to the study of attitudes*. Sociometry 1 (1): 90–98, 1937.
- [18] *Asch conformity experiments*. Wikipedia, the free encyclopedia.
- [19] C. B. Keasey; *Experimentally induced changes in moral opinions and reasoning*. Journal of Personality and Social Psychology 26(1) :30–38 , 1973.
- [20] J. Turner; *Social categorization and the self-concept: A social cognitive theory of group behavior*. Advances in group processes: Theory and research 2: 77–112, Greenwich, CT: JAI press, 1985.
- [21] J. Turner, M. Hogg, P. Oakes, S. Reicher e M. Wetherell; *Rediscovering the social group: A self-categorization theory*. Oxford: Blackwell, 1987.
- [22] J. Turner; *Social influence*. Milton Keynes: open university press, 1991.
- [23] M. Deutsch e G. Harold; *A study of normative and informational social influences upon individual judgement*. J. of Abnor. and Soc. Psych. 51: 629–636, 1955.
- [24] J. Alford, C. Funk e J. Hibbing; *Are political orientations genetically transmitted?* R. Am. Polit. Sci. Rev. 99: 153–167, 2005.

- [25] J. Block e J. Block; *Nursery school personality and political orientation two decades later*. J. Res. Pers. 40: 734–749, 2006.
- [26] E. K. Miller e J. D. Cohen; *An integrative theory of prefrontal cortex function*. Annu. Rev. Neurosci. 24: 167–202, 2001.
- [27] M. M. Botvinick, T. S. Braver, D. M. Barch, C. S. Carter e J. D. Cohen; *Conflict monitoring and cognitive control*. Psychol. Rev. 108: 624–652, 2001.
- [28] J. T. Jost, Jack Glaser, A. W. Kruglanski e F. J. Sulloway; *Political conservatism as motivated social cognition*. Psychological Bulletin, 129(3):339–375, 2003.
- [29] W. McCulloch e W. Pitts; *A logical calculus of the ideas immanent in nervous activity*. Bulletin of Mathematical Biophysics, 7:115–133, 1943.
- [30] A. Engel e C. Van den Broeck; *Statistical mechanics of learning*. Cap 1: 2, Cambridge Press, 2001.
- [31] D. O. Hebb; *The organization of behavior*. Wiley, Nova York, 1949.
- [32] F. Rosenblatt; *The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain*. Cornell Aeronautical Laboratory, Psych. Rev. 65(6):386–408, 1958.
- [33] O. Kinouchi e N. Caticha; *Optimal generalization in perceptrons*. J. Phys. A Math. Gen. 25: 6243–6250, 1992.
- [34] J. E. César; *Mecânica estatística de sistemas de agentes bayesianos: aplicação à teoria dos fundamentos morais*. Tese de doutoramento, orientador: Nestor Caticha, Insituto de Física, USP, 2013.
- [35] R. Albert e A. L. Barabasi; *Statistical mechanics of complex networks*. Revi. of Mod. Phys. 74: 72–75, 2002.
- [36] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, E. Teller; *Equations of State Calculations by Fast Computing Machines*. Journal of Chemical Physics 21(6): 1087–1092, 1953.
- [37] C. Shannon; *The Mathematical Theory of Communication*. Bell System Technical Journal 27: 379–423, 623–656, 1948.
- [38] E. T. Jaynes; *Probabilty Theory: the logic of science*. Cambridge University Press, 2003.

- [39] A. Caticha; *Entropic inference and the foundations of physics*.
- [40] M. V. Day, S. T. Fiske, E. L. Downing e T. E. Trail; *Shifting liberal and conservative attitudes using moral foundations theory*. : Working Paper (2014)
- [41] D. Chong, J. N. Druckman; *Framing Public Opinion in Competitive Democracies*. American Political Science Review 101(4): 637–655, 2007.
- [42] E. F. Fowler, S. E. Gollust, A. F. Dempsey, P. M. Lantz and P. Ubel; *Issue emergency, evolution of controversy and evolution for competitive framing: the case of HPV vaccine*. The International Journal of Press/Politics 17: 169, 2012.
- [43] A. S. Gerber, J. G. Gimpel, D. P. Green, D. R. Shaw; *How large and long-lasting are the persuasive effects of televised campaign ads? Results from a randomized field experience*. American Political Science Review 105(1) : 135–150 ,2011
- [44] D. Chong e J. Druckman; *A theory of framing and opinion formation in competitive elite environments*. Journal of Communication 57: 99–118, 2007.
- [45] R. F. Baumeister e E. J. Finkel; *Advanced Social Psychology: The State of the Science*. Oxford Press.
- [46] P. Fernbach, T. Rogers, C. Fox, S. Sloman; *Political Extremism Is Supported by an Illusion of Understanding*. Psychological Science 24(6): 939–946 , 2013.