## Discovering Archetypes to Interpret Evolution of Individual Behavior

Kanika Narang\*, Austin Chung\*, Hari Sundaram\*, Snigdha Chaturvedi<sup>§</sup>

\* University of Illinois, Urbana-Champaign <sup>§</sup> University of California, Santa Cruz

#### Abstract

In this paper, we aim to discover archetypical patterns of individual evolution in large social networks. In our work, an archetype comprises of *progressive stages* of distinct behavior. We introduce a novel Gaussian Hidden Markov Model (G-HMM) Cluster to identify archetypes of evolutionary patterns. G-HMMs allow for: near limitless behavioral variation; imposing constraints on how individuals can evolve; different evolutionary rates; and are parsimonious.

Our experiments with Academic and StackExchange dataset discover insightful archetypes. We identify four archetypes for researchers: Steady, Diverse, Evolving and Diffuse. We observe clear differences in evolution of male and female researchers within the same archetype. Specifically, women and men differ within an archetype (e.g. Diverse) in how they start, how they transition and the time spent in mid-career. We also found that the differences in grant income are better explained by the differences in archetype than by differences in gender. For StackOverflow, discovered archetypes could be labeled as: Experts, Seekers, Enthusiasts and Facilitators. We have strong quantitative results with competing baselines for activity prediction and perplexity. For future session prediction, the proposed G-HMM cluster model improves by an average of 32% for different Stack Exchanges and 24% for Academic dataset. Our model also exhibits lower perplexity than the baselines.

## 1 Introduction

In this paper, we develop models to understand how individuals evolve in large social networks. The problem is important: as individuals interact within the context established by social norms, they gain in experience, and behavioral changes reflect the newfound experience. However, despite a significant focus on community discovery and their evolution in social networks, our understanding of individual evolution is limited (McAuley et al. (2013) and Yang et al. (2011) are some notable exceptions). Understanding evolutionary patterns is useful in a variety of applications: language evolution (Danescu-Niculescu-Mizil et al., 2013); expertise evolution (McAuley et al., 2013); journey optimization in digital advertising platforms.

Despite variations in how individuals can evolve, we observe regularities. For instance, for an academic, transition

Copyright © 2018, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

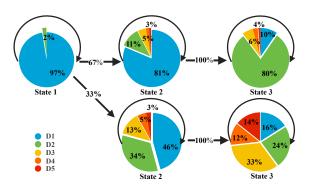


Figure 1: A stylized academic evolutionary trajectory. Each pie chart is a behavior stage in the trajectory. It shows the fraction of papers published in each research area  $D_m$  in that stage. We use a normalized representation focused on change of areas: the label  $D_1$  represents the first research area of every academic,  $D_2$  the second research area etc. Normalized representations allow us to discover commonalities in behavioral changes of academics across seemingly unconnected domains. In this example, top group of researchers evolve to shift their research focus to a new domain while the bottom group becomes increasingly interdisciplinary.

through different stages—PhD Student (focusing on a single research area), being an assistant professor (working on highly related areas) to eventually post-tenure (multiple areas, interests in multidisciplinary collaborations etc.)—mark changes in behavior. These elementary behavioral evolutionary patterns are visible in almost all academic fields, suggesting that surface variations (i.e. area of research for an academic) hide deeper regularities in patterns of behavioral change. We refer to these *latent* regularities in individual behavior as *archetypes* and we plan to explain all individuals' surface variations (the observed research area on which the academic focuses) with a small set of archetypes. Figure 1 shows a stylized example.

Thus a model for learning archetypes needs to: express large observable behavioral variation while exhibiting latent stochastic regularities governing the change of behavior. Furthermore, the model should allow individuals to evolve at different rates. Finally, the results ought to be interpretable in a post-hoc manner.

Our work makes the following contributions:

#### A framework for modeling evolutionary trajectories:

We propose a sophisticated framework to identify

dominant, interpretable, archetypes amongst individuals for modeling the evolution of their research interests. In contrast, prior work on user modeling has either focused on qualitative analysis (e.g., (Ward, 2001)), or use engineered features (Angeletou et al., 2011) or ignored temporal changes (Mamykina et al., 2011). In our work, we assume that an archetype is a probabilistic model that encodes individual progression through stages of distinct behavior. Specifically, we learn a Gaussian Hidden Markov Model (G-HMM) to capture this progression where latent states capture behavioral stages in the evolution. To encode the idea of experience, while we allow individuals to continuously evolve onto the higher stages, we constrain our model to prevent individuals from returning to a stage from which they have evolved. We model all individuals with a small set of archetypes. We jointly learn the mapping of users into their archetype and the archetype's model parameters through an Expectation-Maximization framework.

**Identification of dominant archetypes:** We apply our model to understand the evolution of research interests of Computer Scientists and the users of *Stack Exchange* by studying their activities on the platform. Specifically, for academic dataset, we identify four archetypes: (i) *Steady* researchers who primarily work in their first research area through out their career; (ii) *Evolving* researchers, who continuously shift their dominant area of research; (iii) researchers with *Diverse* research interests; and (iv) researchers who have *Diffused* interests with infrequent contributions in multiple areas. Each archetype is significantly different (p < .001) from the others.

Qualitative analysis—variation by gender: We examined qualitatively, a subset of our data—all full professors (as of Spring 2018) in the top 50 CS departments in United States for gender differences in their academic trajectory. We observe significant differences in the models that explain the evolution of male and female researchers within the same archetype. For example, the models that explain women and men differ (significance level: p < .01) in the *diverse* archetype; we observe differences in where they start, rates of transition and research interests during mid-career.

Qualitative analysis—effect on grant income: We examine grant income (as of Spring 2018) from the National Science Foundation in US for the same subset of CS academics, to understand the relationship between variations in grant income over the course of academic trajectory and how difference in archetype or gender could serve as explanations. We find significant differences in grant income across genders within a behavioral stage of an archetype. For example, for the steady archetype, there are differences (p < .05) in grant income across genders at the starting stage. Regardless of gender, we also find significant differences across behavioral stages within an archetype. For example, we find significant differences (p < .001) in grant income between stages 2 and 3 of all archetypes.

Additionally, we have strong quantitative results with

competing baselines for activity prediction and perplexity on both Academic and Stack Exchange communities. The proposed G-HMM cluster model improves by 24% for Academic and on an average of 32% for Stack Exchange communities for future session prediction. Our model exhibits lower perplexity than the baselines. Our model improves by 149% and 25% for predicting trajectory of unseen users for Academic and Stack Exchange communities respectively.

**Significance:** We propose a sophisticated probabilistic framework to identify dominant, interpretable, evolutionary archetypes. We show that the discovered archetypes are significantly different and are straightforward to use to test hypotheses (e.g. evolutionary variation with gender; effects of gender on income).

We organize the rest of this paper as follows. We discuss related work in Section 2. In Section 3, we formally describe our model followed by dataset description in Section 4. We then discuss discovered trajectories in Section 6. Section 5 describes our experiments and its limitations in Section 7 and we finally conclude in Section 8.

## 2 Related Work

Activity Modeling: Our work is most similar to activity sequence modeling that predicts next action or event in a sequence. We are different from those works as our focus is on modeling user *behavior* that is how a user spends her time among possible actions in each session. Yang et al. (2014) and Knab et al. (2003) proposed generative models that assigns each action to a progression stage and classify event sequences simultaneously. They used their model to predict cancer symptoms or products user would review in the future. However the model did little to provide meaningful and interpretable stages and clusters. The major contribution of our work is in giving an *interpretable* model that helps us to characterize the temporal changes in user *behavior*.

Hidden Markov Model (HMM) have been used to model and cluster time sequences (Bicego et al., 2003; Coviello et al., 2014; Smyth, 1997) in the past. However, most of these models learn an HMM for each user sequence and then employ clustering algorithms to cluster the learned HMMs. These approaches are not scalable and the clusters thus identified are not interpretable.

User Profiling: There also has been work in the past on identifying and characterizing user roles in online social networks (OSNs). Maia et al. (2008) identified five distinct user behaviors of YouTube users based on their individual and social attributes. While Mamykina et al. (2011) identified user roles based on just answer frequency in StackExchange. Similar user behavioral studies are done by (Adamic et al., 2008) and (Furtado et al., 2013) on Yahoo Answers and Stack Overflow respectively. These studies however ignore *temporal changes* in the behavior and use engineered features for behavior modeling.

Some behavioral studies do model the evolution too. Benevenuto et al. (2009) learnt a Markov model to examine transition behavior of users between different activities in Orkut in a static snapshot. Angeletou et al. (2011) constructed hand crafted rules to identify user roles and study change of user roles' composition in the community over time. Our model,

instead, works directly on raw activity data and cluster users with similar pattern of behavioral *evolution*.

Academic Data Mining: There has been extensive interest in mining Academic Data (bibliographic data, researchers' usage of social media etc.). Studies have been done to understand the evolution of research interests on a community level. Biryukov et al. (2010) worked on understanding scientific communities in DBLP dataset while Liu et al. (2014) focused on evolution of research themes in CHI papers over time. Chakraborty et al. (2018) studied trajectories of successful papers in computer science and physics by analyzing paper citation counts. On an individual level, Safavi et al. (2018) studied career transitions between academia and industry for Computer Science researchers. Studies have also explored usage of Twitter by the Academic community (Hadgu et al., 2014; Linek et al., 2017; Mehrazar et al., 2018). In contrast, our focus is on finding commonalities in evolution of research interests of scientists across subdomains.

Recent studies also look at gender differences in funding patterns, productivity and collaboration trends in academia (Way et al., 2016, 2017). Some earlier studies also reported gender differences in academia. Kahn (1993) identified gendered barriers in obtaining tenure for academics in economics, while Ward (2001) found gendered differences in pay related to publication record. On the other hand, we explore *gender differences* in a complementary dimension of change in research interests and its effect on grant income.

## 3 Modeling Evolution Trajectories

In this section, we first describe the model requirements followed by the problem definition. Then, we present our approach that satisfies those requirements.

## 3.1 Model Requirements

An individual's behavior evolves with experience. We operationalize the notion of experience as a *progression through behavioral stages*. A *behavioral stage* is a period of time where individuals exhibits stochastic regularities in their behavior. An individual's *trajectory* is a specific progression through behavioral stages.

We expect individual *evolutionary trajectories* to be unique. However, as the stylized example from the introduction on academic life suggests, despite differences in subfields (say HCI vs Data Mining), individuals show latent regularities over time—that is, we observed latent regularities in *how they change*, and *not in what they do*. Latent regularities suggests that we can represent *all* individual trajectories with a finite set of latent, dominant evolutionary trajectories denoted as *archetypes*. Therefore, models that discover archetypes should satisfy these requirements:

**Large observed variation:** The model should allow for large (possibly infinite) observed individual behavioral trajectories; each with different behavioral stages.

**Stochastic regularities:** The model should reflect two empirical observations: individuals appear to progress through a series of *distinct* behavioral stages reflecting

gain in experience; individuals exhibit regularities in how they evolve through these stages.

**Different rates of evolution:** Individuals evolve at different rates and can skip stages.

**Parsimonious:** We should be able to approximate individual trajectories by small number of *dominant* trajectories.

We want the model to be post-hoc interpretable (Lipton, 2016). That is, we would like model outcomes, including the definition of stages and how individuals progress amongst them to allow for meaningful interpretation (e.g. 'the person started to work on multidisciplinary research after tenure').

#### 3.2 Problem Definition

Now, we discuss the problem formally. We represent an individual i, as a time ordered sequence,  $\mathbf{X_i}$ , of their activities, over a time granularity appropriate to that domain. For example, for the Stack Exchange data, the granularity is that of a single visit to the network, while in case of the Academic corpus, the granularity is that of a year since most conferences occur annually.

Without loss of generality, we refer to the fundamental temporal unit of analysis as a session. Thus,  $\mathbf{X_i}$  is a sequence of sessions,  $\mathbf{X}_{ij}$ , where  $j \in \{1,2,\ldots t_i\}$  and  $t_i$  is the number of sessions for an individual i. In general, lengths of sequences will vary across individuals depending on their activity level. A session,  $\mathbf{X}_{ij}$ , is a vector  $\langle o_1, o_2, \ldots, o_M \rangle$ , where M is the number of different actions possible in that network. Each element  $o_m$  of the vector  $\mathbf{X}_{ij}$ , denotes the fraction of time the individual performs the m-th action during a single session. For example, for the Stack Exchange data, if the set of possible actions include 'posting a question', 'answering a question', and 'commenting on answers or on questions', then a session is a distribution over these three actions during a single visit to the social network.

The problem addressed in this paper is to associate an archetype with each individual. We assume that there exist C different archetypes, and given a sequence of sessions for an individual  $\mathbf{X}_i = \{\mathbf{X}_{i1}, \mathbf{X}_{i2} \dots \mathbf{X}_{it_i}\}$ , the goal is to assign the sequence to one of the C archetypes—each associated with a set of K latent behavioral stages. During this assignment, we also identify how the individual evolves through its archetype's distinct stages by outputting the sequence  $Y_i = \{Y_{i1}, Y_{i2} \dots Y_{it_i}\}$ , where  $Y_{ij}$  represents the behavioral stage  $k \in [1, K]$  assigned to j-th session in individual i's evolutionary trajectory. We constrain the number of stages  $K \ll t_i$ , and allow skipping of stages, while disallowing return to earlier stages.

#### 3.3 A Framework for Identifying Archetypes

We use a Gaussian-Hidden Markov Model (G-HMM) based approach to model individual behavior. To capture broad variations amongst individuals, we learn a set of C G-HMMs where each G-HMM represents an archetype. We jointly learn the partitioning of the individuals into different archetypes and the model parameters for each archetype.

Each Gaussian HMM, associated with an archetype c, has K discrete latent states and  $\pi^c$  is a K dimensional prior vector for its latent states. The model makes a first order Marko-

vian assumption between state transitions using the transition probability matrix  $\tau^{\mathbf{c}}$ ; where  $\tau_{kl}^c$  represents the probability of transitioning from state k to l in the c-th archetype. Lastly, the model assumes that given a latent state, k, from an archetype c, the M dimensional session vector,  $X_{ij}$ , is Normally distributed with mean  $\mu_k^c$  and covariance  $\Sigma_{\mathbf{k}}^{\mathbf{c}}$ .

In the above generative process, the G-HMM associated with different archetypes do not share latent states. In other words, each G-HMM has its own set of discrete latent states. However, we fix the number of states (K) to be the same for each archetype.

Encoding Experience & Variable Evolutionary Rates: To encode the idea of experience, as well as to allow variable evolutionary rates, similar to Yang et al., 2014, we allow only forward state transitions (including self loop) within a G-HMM that represents an archetype. This choice appears sensible to us since semantically, each latent state of the G-HMM represents a behavioral stage of evolution, and its corresponding mean vector encapsulates behavior distribution in that stage. Then, forward transition denotes progression through behavioral stages. We operationalize this idea by using an upper triangular state transition matrix.

**Training:** We train our Gaussian HMM archetype model using a (hard) Expectation Maximization (Dempster et al., 1977) based iterative procedure described in Algorithm 1. During training, the goal is to learn the G-HMM parameters,  $\lambda^{\mathbf{c}}$ , for each archetype c, where  $\lambda^{\mathbf{c}} = \langle \mu^{c}, \Sigma^{\mathbf{c}}, \pi^{c}, \tau^{c} \rangle$ and archetype assignments for each user,  $c_i$ . We begin with initializing the Gaussian HMMs with initial parameters,  $\lambda_0^1, \lambda_0^2, \dots, \lambda_0^C$ . Thereafter, in the iterative training process, in the Expectation step, we use current estimates of  $\lambda^{c}$  's to assign an archetype to each user sequence in the data. In the Maximization step, we use current archetype assignments to learn the corresponding G-HMM's parameters  $\lambda^{c}$ . We use a modified version of the Baum-Welch algorithm (Rabiner, 1990) allowing for forward-only transitions. Thus, this method jointly partitions the input sequences into different archetypes as well as learns the parameters of the associated G-HMMs.

**Implementation Details:** Our iterative training procedure requires initialization for G-HMM parameters  $\lambda_0^c$ . We perform k-means clustering on all sessions of all user sequences in our corpus, treating the sessions as independent of each other (thus losing the sequential information). The cluster centers, thus obtained are used as the initial means,  $\mu_0^c$ , for the latent states. We fix each  $\Sigma_k^c$  as an identical diagonal covariance matrix  $\sigma I$  with  $\sigma=0.01$  based on preliminary experiments. We initialize transition matrices,  $\tau_0^c$ , and states' prior probabilities,  $\pi_0^c$ , for each archetype randomly.

Our implementation is based on Kevin Murphy's HMM Matlab toolbox <sup>2</sup>. Also, we implement a parallelized version of our EM algorithm to reduce computation time. We test our model on Intel Xeon Processor with 128 Gb RAM and clock speed of 2.5 GHz. Our model takes around 10 minutes on our biggest dataset Stack Overflow (570K users) using

## Algorithm 1: Gaussian HMM archetype

Input:  $\mathbf{X_i}$  and  $\lambda_\mathbf{0}^\mathbf{c}$   $\forall i \in \{1, 2, \dots N\} \ \forall c \in \{1, 2, \dots C\};$ Output:  $\mathbf{Y_i}$  and  $\lambda^\mathbf{c}$   $\forall i \in \{1, 2, \dots N\} \ \forall c \in \{1, 2, \dots C\};$ Initialize the  $c^{th}$  archetype with initial parameters,  $\lambda_\mathbf{0}^\mathbf{c}$   $\forall c$ ; while not converged do

M-Step: Re-assign archetypes to sequences  $\mathbf{X_i}$  as:  $c_i = argmax_c P(\mathbf{X_i}|\lambda^\mathbf{c}) \ \forall i \in \{1, 2, \dots N\};$ E-Step: Re-estimate the G-HMM parameters,  $\lambda^\mathbf{c} \forall c \in \{1, 2, \dots C\},$  using modified Baum-Welch

#### end

#### Convergence Criteria;

algorithm.;

- · Log Likelihood difference falls below threshold; or
- Number of iterations is greater than threshold; or
- Number of sequences re-assigned in an iteration is less than 1% of the data

parallelization with 30 threads.

#### 4 Dataset

For our analysis, we use the Microsoft Academic Dataset and Stack Exchange Dataset. Table 1 shows the data statistics, and we provide a brief description of the two corpora in the rest of this section.

| Dataset        | N      | $\bar{t}$ | $t_{\rm max}$ | M |
|----------------|--------|-----------|---------------|---|
| Academic       | 4578   | 24.15     | 47            | 6 |
| StackOverflow  | 561937 | 47.13     | 750           | 5 |
| English        | 3828   | 44.01     | 729           | 5 |
| Money          | 873    | 44.41     | 706           | 5 |
| Movies         | 678    | 48.40     | 598           | 5 |
| CrossValidated | 3728   | 38.94     | 738           | 5 |
| Travel         | 1000   | 56.14     | 736           | 5 |
| Law            | 195    | 47.79     | 584           | 5 |

Table 1: Dataset statistics for the Academic and Stack Exchange datasets. N: number of users; M: possible actions in each session;  $t_{\max}$ : maximum session length;  $\bar{t}$ : mean session length. For authors,  $\bar{t}$  is their average career length (in years).

## 4.1 Microsoft Academic Dataset

We use the Microsoft Academic dataset to study evolutionary patterns of researchers with a focus on Computer Scientists. To this end, we extract publication history of authors in Computer Science (CS) using the Microsoft Academic Knowledge Service API <sup>3</sup>. Microsoft Academic Service additionally annotates each publication with the year of publication, publication venue and the CS subfield (out of 35 identified fields) to which it belongs.

For this study, we decide to focus only on *influential* scientists with sufficient publication history. We identify *influential* authors based on *prominence* of the conference venues in which they publish. To quantify *prominence* of a conference, we construct a conference-conference citation graph where each conference in our dataset forms a node and

<sup>&</sup>lt;sup>1</sup> Experiments with tied-states of archetypes led to worse results

<sup>&</sup>lt;sup>2</sup>bit.ly/hmmtoolbox

<sup>3</sup>http://bit.ly/microsoft-data

the weighted edges represent inter-conference citation frequency. Specifically, the weight of a directed edge from conference  $C_1$  to conference  $C_2$  is proportional to the fraction of papers published in  $C_2$  cited by papers published in  $C_1$ . We then use the Pagerank algorithm (The PageRank Citation Ranking: Bringing Order to the Web. 1999) on this directed graph and define conference prominence as the Pagerank of the corresponding conference-node. Thereafter, we define an author's influence as the weighted sum of prominences of the conferences (s)he has published in. Here, conference-prominences are weighted by the fraction of the author's papers published in that venue.

We rank authors in decreasing order of their *influence* and extract top 750 most-influential authors from each of the 35 CS areas in the dataset. Note that authors can be *influential* in more than one subfield. We then extract unique authors from this set who have at least 15 years of publication history. We only consider publication history from 1970 to 2016 to avoid missing data. The resulting dataset consists of records of 4578 authors.<sup>4</sup>

We now describe how we represent an author's academic life-cycle as a sequence,  $X_i$ , comprising of session-vectors,  $X_{ij}$ . We chose each session to be an year long as most CS conferences occur annually. For this dataset, a session-vector represents the fraction of papers an author publishes in various area-of-interests (AoIs) in that year.

For defining an AoI of an author, we consider all papers published by the author in her academic life. We identify her primary AoI,  $D_1$ , as the *first* subfield (out of 35 subfields) in which she publishes cumulatively at least 3 papers in the first 3 years. Usually, an author's  $D_1$  is about their PhD dissertation work and we expect students to settle down after a few years. Thus, after identification of  $D_1$ , hopefully with a steady paper count, we define her secondary AoI,  $D_2$ , as the subfield in which she publishes at least 3 papers in *one* year. Similarly, we also define tertiary  $(D_3)$ , quaternary  $(D_4)$ , and quinary  $(D_5)$  AoI. We do not define AoIs beyond  $D_5$  because 80% of authors do not explore more than 5 subfields in our dataset. Also, in a given year, if an author publishes fewer than 3 papers in an unexplored subfield, these papers count towards a sixth dimension AoI called Explore. This denotes that the author has started exploring new subfields but they are not yet significant enough to be one of the  $D_m$ 's  $(m \in [1, 5])$ , and indicate a possible shift in research interests. To summarize, each session is a 6 dimensional vector (M = 6), and its elements are the fraction of the author's publications in one of the 5  $D_m$ 's or the  $6^{th}$  Explore dimension. This normalized representation for sessions allows our model to discover behavioral patterns of author's changing research interests in a domain independent manner.

#### 4.2 Stack Exchange Dataset

Our second dataset consists of activity logs of users of Stack Exchange <sup>5</sup>(as of Feb 2017), a popular online question-answering platform. In this paper, we work on 7 diverse communities of the platform: Stack Overflow, English,

Money, Movies, CrossValidated, Travel and Law. These communities have varied sizes and cater to different audiences. For each user, the data contains details about their activities on the community. Stack Exchange allows 5 different activities (M = 5): post a Question; Answer a question; Comment on a question or an answer; Edit operations like assign tags, edit body or title of a post; and Moderator operations like voting. Like before, we represent a user by a sequence,  $X_i$ , of session vectors,  $X_{ij}$ . We split the activitysequence of a user into sessions using a time threshold similar to session definitions in web search (Narang et al., 2017). Specifically, we create a new session if the difference between two consecutive activities is more than 6 hours. A gap longer than this marks a new visit to the community. Hence, a session is a subsequence of the user's activity-sequence and is formally represented as a distribution over the M possible activities; where its  $m^{th}$  element represents the fraction of total activity spent in the  $m^{th}$  activity in that session.

Lastly, to focus on users who have spent enough time in the network to exhibit behavioral changes, we filter users with less than 10 sessions, and also remove outliers with more than 750 sessions.

## 5 Quantitative Experiments

In this section, we evaluate our model on two different tasks: Future Prediction and Perplexity. We describe the baselines in Section 5.1 and report results in Section 5.2.

#### 5.1 Baselines

**Distance GHMM**: Our first baseline uses the GHMM clustering model as defined in Ghassempour et al., 2014. In this baseline, we learn a GHMM for each user and then cluster the models using distance metric  $\delta$ , the symmetric KL divergence  $(d_{kl})$  between two G-HMMs (Juang et al., 1985).

$$d_{kl}(\lambda^p, \lambda^q) = \frac{1}{N_p} \sum_{i \in N_p} log \frac{P(X_i | \lambda^p)}{P(X_i | \lambda^q)}, \tag{1}$$

We use k-medoids clustering; since this method doesn't give a representative model for each cluster, we learn a GHMM per cluster. For fair comparison, we set k, the number of clusters to be the same as our model.

**Vector AutoRegressive Model (VAR):** VAR models are used to model multivariate time series data (Ltkepohl, 2007). It assumes that each variable in the vector is a linear function of it's own past values as well as other variables. For each user sequence  $\mathbf{X_i}$ , jth session is modeled as,

$$\mathbf{X}_{ij} = A_1 \mathbf{X}_{ij-1} + \ldots + A_p \mathbf{X}_{ij-p} + u_j \tag{2}$$

where  $A_i$  is  $M \times M$  matrix,  $u_j \sim \mathcal{N}(0, \Sigma_u)$  and we set p = 1 as in first-order Markov models.

**Gaussian clusters (GCluster):** In this baseline, we assume that individuals *do not evolve* in their lifespan. This is a simplified version of our model. It assumes that there are different archetypes but that each archetype has only one state. Hence, all sessions of a sequence are generated from a single multivariate Gaussian.

We can not compare with other sequence prediction baselines (Knab et al., 2003; Yang et al., 2014) as they assume a discrete set of activities while in our case, each session is a probability distribution over possible activities.

<sup>&</sup>lt;sup>4</sup>This data will be made available upon publication

<sup>5</sup>https://data.stackexchange.com/

#### 5.2 Tasks

Future Prediction: In this task, we predict future behavior of an individual given her history. We assign the first 90% sessions of each sequence for training, and predict the behavior in future sessions (the remaining 10% of the sequence). We first use all the training sessions to learn parameters of our model. Then, for each sequence, we run Viterbi algorithm to decode state assignment of its test sessions,  $t_i'$ . The test sessions of the i-th user will have same archetype assignment  $c_i$  determined in the training session for that user.

We compute Jensen-Shannon $(d_{js})$  divergence between the mean  $\mu^{c_{ij}}$  of the assigned state  $Y_{ij}$  and the observed vector  $X_{ij}$ .  $d_{js}$  is a symmetric K-L divergence between two vectors. We report the average  $\bar{\Delta}$  over all test sessions:

$$\bar{\Delta} = \frac{1}{|T|} \sum_{i \in N, j \in t'_i} d_{js}(\mu^{c_{ij}}, X_{ij}), \qquad (3)$$

$$d_{js}(\mu^{c_{ij}}, X_{ij}) = \frac{1}{2} d_{kl}(\mu^{c_{ij}}, p) + \frac{1}{2} d_{kl}(X_{ij}, p),$$
 (4)

where,  $p = \frac{1}{2}(\mu^{c_{ij}} + X_{ij})$  and  $d_{kl}$  measures KL divergence distance. For VAR, we use the model learnt on training sessions of user i to make prediction for her future sessions.

Table 2 shows our results on this task for different datasets. Our model outperforms the baselines for all Stack Exchange datasets with an average improvement of about 32% and 24% on the Academic dataset. Hence, learning archetypes can help us to accurately predict an individual's future behavior in the social network.

| Dataset        | Our Model | VAR  | Distance<br>HMM | Gaussian<br>Cluster |
|----------------|-----------|------|-----------------|---------------------|
| Academic       | 0.22      | 0.31 | 0.42            | 0.29                |
| StackOverflow  | 0.23      | 0.36 | NA              | 0.37                |
| English        | 0.19      | 0.29 | 0.26            | 0.31                |
| Money          | 0.19      | 0.52 | 0.32            | 0.32                |
| Movies         | 0.23      | 0.35 | 0.35            | 0.37                |
| CrossValidated | 0.21      | 0.38 | 0.33            | 0.35                |
| Travel         | 0.19      | 0.30 | 0.25            | 0.29                |
| Law            | 0.19      | 0.26 | 0.33            | 0.27                |

Table 2: Average Jensen-Shannon divergence of future sessions using 90-10% split of each user sequence. Lower values are better. Distance HMM did not converge on StackOverflow dataset.

**Perplexity** Perplexity measures how surprised the model is on observing an unseen user sequence. A lower value of perplexity indicates low surprise and hence a better model.

$$P_x = -\frac{1}{|T|} \sum_{i \in T} \max_{c \in C} (\log P(\mathbf{X_i^T} | \lambda^c))$$
 (5)

where,  $\mathbf{X_i^T}$  represents a test sequence in Test Set T, and  $\lambda_c$  represents the parameters of the GHMM corresponding to the c-th archetype. We assign  $\mathbf{X_i^T}$  to the archetype c with maximum likelihood. Perplexity is then computed as the average likelihood of all test sequences.

Table 3 reports average perplexity after five fold cross validation. Note that for this experiment, model predicts entire trajectory of a new user. We could not use the regression baseline (VAR) as it is not a generative model. Our

model beats best performing baseline by 149% on Academic and by around 25% on average for StackExchange datasets. Hence, our model also effectively predicts behavior of future individuals joining the social network.

| Dataset        | Our Model | Distance<br>HMM | Gaussian<br>Cluster |
|----------------|-----------|-----------------|---------------------|
| Academic       | -18.37    | 37.73           | 100.79              |
| StackOverflow  | 487.68    | NA              | 678.62              |
| English        | 306.38    | 559.65          | 471.14              |
| Money          | 415.85    | 557.69          | 570.51              |
| Movies         | 596.10    | 724.15          | 743.73              |
| CrossValidated | 398.44    | 514.74          | 554.31              |
| Travel         | 494.06    | 645.64          | 666.97              |
| Law            | 368.89    | 508.08          | 482.27              |

Table 3: Average Perplexity on unseen user sequences after 5 fold cross validation. Lower values are better. DistanceHMM did not converge on StackOverflow dataset.

**Discussion:** For future prediction, our model performs better than VAR model. It shows that modeling cluster of sequences gives a better estimate than modeling each user sequence separately. Also, if we assume no behavior evolution and just cluster users according to their behavior i.e. *GCluster* model, we obtain worse results. Our model also outperforms similarity distance based clustering method: DistanceHMM (Ghassempour et al., 2014), which is also the strongest baseline. It first estimates G-HMM model for each user sequence and then cluster these models. Estimating model for each sequence can be noisy, specially if the user sequence has short length. Instead, when we jointly learn G-HMM model parameters and cluster sequences, we learn a better approximation.

Full vs. Left-Right Transition Matrix: We also test our model with unconstrained full transition matrix where users can jump from a state to any other state in the HMM. We compare our results with this model for future prediction task. On academic dataset, we obtain similar values while for StackExchange communities, the full transition matrix gives better results. This can be because a full matrix has more degrees of freedom but then, it is also more expensive to learn. Also, the states learnt are not interpretable. As (Yang et al., 2014) and (Knab et al., 2003) noted, forward state transitions accurately models the natural progression of evolution, we thus, chose to work with a forward transition matrix.

#### **6** Qualitative Analysis

In this section, we perform qualitative analysis of archetypes identified by our model. Due to space constraints, we analyze the archetypes discovered by our model only for the Academic Dataset and Stack Overflow (programming based Q&A); largest and most popular community in StackExchange. We first describe the discovered archetypes of all researchers in Section 6.1. Then, we examine gender variation in academic trajectory in Section 6.2 and effect of archetype and gender on grant income in Section 6.3. Finally, we details archetypes for Stack Overflow in Section 6.4.

|         | Steady                       | Diverse                      | Evolving                     | Diffuse              |
|---------|------------------------------|------------------------------|------------------------------|----------------------|
| State 1 | {3Y, 5m}                     | {3Y, 3m}                     | {2Y, 9m}                     | {2Y, 7m}             |
|         | $\mathbf{\hat{D}_{1}}$ (87%) | $\mathbf{\hat{D}_{1}}$ (88%) | $\mathbf{\hat{D}_{1}}$ (72%) | $\mathbf{D_1}(76\%)$ |
|         | Ex (11%)                     | Ex (11%)                     | Ex (24%)                     | Ex (22%)             |
| State 2 | $\{4Y, 2m\}$                 | {2Y, 6m }                    | {2Y, 9m }                    | {3Y, 7m}             |
|         | Ex (74%)                     | Ex (80%)                     | Ex (83%)                     | Ex (91%)             |
|         | $D_1$ (23%)                  | $D_1$ (16%)                  | $D_1$ (12%)                  |                      |
| State 3 | $\{7Y, 5m\}$                 | $\{5Y, 6m\}$                 | $\{6Y, 2m\}$                 | $\{8Y, 5m\}$         |
|         | $D_1$ (62%)                  | $D_1$ (73%)                  | $D_1$ (33%)                  | $D_1$ (50%)          |
|         | Ex (32%)                     | Ex (17%)                     | Ex (28%)                     | Ex (39%)             |
|         |                              |                              | $D_2$ (24%)                  |                      |
| State 4 | {5Y, 9m}                     | $\{5Y, 6m\}$                 | {5Y}                         | $\{3Y, 9m\}$         |
|         | $D_2$ (49%)                  | Ex (46%)                     | $D_2$ (66%)                  | $D_2$ (43%)          |
|         | $D_1$ (27%)                  | $D_2$ (20%)                  | Ex (18%)                     | Ex (26%)             |
|         | Ex (17%)                     | $D_1$ (17%)                  |                              |                      |
| State 5 | $\{2Y, 6m\}$                 | $\{6Y, 3m\}$                 | $\{6Y, 5m\}$                 | $\{4Y, 1m\}$         |
|         | $D_1$ (49%)                  | $D_4$ (29%)                  | $D_3$ (43%)                  | Ex (74%)             |
|         | Ex (18%)                     | Ex (20%)                     | Ex (19%)                     |                      |
|         | $D_2$ (14%)                  | $D_3$ (14%)                  | $D_2$ (14%)                  |                      |
|         |                              | $D_1$ (14%)                  |                              |                      |

Table 4: Learned mean vector for each state for four archetypes in the Academic Dataset. We list the *Area-of-Interests* (AoI) in sorted order and annotate them with their % contribution in the state. We list main AoI (> 11%) for each state. Each state is also labeled with it's average duration in {Years (Y), months (m)}. The labels given to these clusters reflect our own interpretation of the user behavior and make disambiguating the behavior easier in the text.

## 6.1 Academic Archetypes

Our analysis reveals four archetypes: Steady, Diverse, Evolving and Diffuse. We chose the number of clusters C=4 using the elbow method Tibshirani et al., 2001: data log likelihoods increased rapidly till four clusters with much slower increase beyond that. Further, we chose number of states per cluster, K=5: beyond five states, KL divergenceKullback et al., 1951 between mean vectors of new states with previous states started reducing rapidly indicating redundant states.

We also conducted t-test to validate differences among the identified archetypes. Specifically, paired-sample t-test Goulden, 1949 is conducted between likelihood values of data points assigned to an archetype with their likelihood values obtained from rest of the archetypes. For instance, for each archetype pair (p,q), we conduct paired t-test between  $\log P(X_i|\lambda^p)$  and  $\log P(X_i|\lambda^q) \ \forall i \ni c_i = p$ . Note that test results for archetype pair (p,q) are not symmetric. We observed that all archetype pairs are significantly different (p < .001). Now, we first discuss what is common to these discovered archetypes before examining each one in detail.

Commonalities in Archetypes: Table 4 summarizes the trajectories (state sequences) learned for the four different archetypes in this dataset. Each archetype is labeled according to our own interpretation of the user behavior, looking at the learned mean vector of GHMM states. We observe that all archetypes exhibit similarities, especially in the first two stages. Across all archetypes, the first stage typically spans around 3 years, and more than 72%

|   | Steady | Diverse | Evolving | Diffuse |
|---|--------|---------|----------|---------|
| Male Professors   | 247    | 206     | 241      | 263     |
| (Top-50 US schools) Female Professors (Top-50 US schools) | 30     | 32      | 26       | 39      |
| Total Authors   | 1329   | 1080    | 1107     | 1062    |

Table 5: Statistics for discovered archetypes. Identified Male and Female professors in Top-50 US schools. Total authors includes all researchers and professors in the dataset.

of the published research is in the author's primary  $A \circ I$ :  $D_1$ . As noted before, this is most likely their PhD dissertation area and hence, the research is more focussed. After gaining some research experience, most authors move to the second stage where they start exploring other research areas denoted by marked increase in their Explore  $A \circ I$  (more than 74%). However, in state 3 and beyond, authors from different archetypes follow different trajectories where they differ in how they change their dominant  $A \circ I$  over time while exploring other domains. Below, we describe each of the four trajectories in more details.

Steady: The first major archetype is of steady researchers, who mainly work in *one* AoI (i.e. their  $D_1$ ) throughout their career. Fig 2 shows the states of this archetype. We can see that most people start in their primary AoI,  $D_1$  (state 1), which possibly reflects their PhD education. After graduation, they spend some time exploring other areas while continuing to publish in  $D_1$  (state 2), but move back to publishing in  $D_1$  for a significant portion of their careers, about 7.5 years (state 3). This is often again followed by a phase where they start working in another area,  $D_2$ , while continuing to publish in  $D_1$  (state 4), they eventually revert to publishing in  $D_1$  (state 5) towards the latter part of their careers. In the last state, they also publish widely in other areas (indicated by almost half of the pie divided between other  $D_m$ 's), but their main interest remains  $D_1$ . Michael Jordan, professor at University of California, Berkeley exhibits this research trajectory. He is a Machine Learning expert; his primary AoI  $D_1$ , and has secondary interests in Data Mining, Optimization and Bioinformatics. Theory professor at University of Illinois, Urbana-Champaign, Jeff Erickson is also assigned to this cluster; he also publishes in his primary AoI  $D_1$  (Theory) with auxiliary interests in the field of mathematical optimization.

Diverse: The second archetype consists of researchers with diverse research interests as they make significant contributions in multiple  $D_m$ 's. Similar to steady researchers, they research in their primary  $A \circ I$   $D_1$  while exploring other domains in the initial 3 states as shown in Table 4. They, then, publish in  $D_2$  and  $D_1$  while spending half time exploring other possible interests (state 4). They evolve to have strong research presence in all 5 AoIs (state 5). This behavior suggests that authors of this archetype tend to work in interdisciplinary areas; or projects with broader scope which gain acceptance by different research communities. One notable example is Prof. Jiawei Han at University of Illinois, Urbana-Champaign, who started his academic career study-

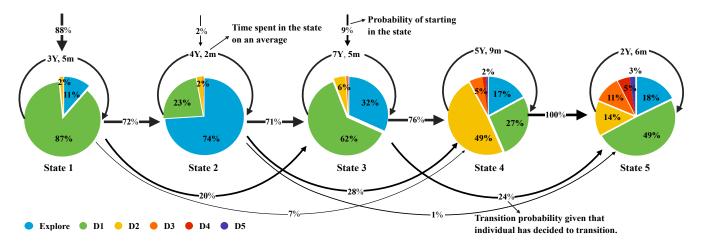


Figure 2: Trajectory (state sequence) for *Steady* archetype in the Academic Dataset. Each pie is a *latent state* or *behavior stage* in the trajectory. It denotes the proportion of papers published in each *Area of Interest*'s in the latent state. Each state is also labeled with the average amount of time spent in the state. For example, in this cluster, 87% of publications in the first 3.5 years are in author's primary  $A \circ I$  while rest 11% are in exploring other areas. The arrows on the top of each pie show prior probability for starting in that state. As we learn a left-to-right G-HMM, author can transition to its immediate next state or any later latent states. Each transition is labeled with the corresponding conditional transition probability i.e. transition probability given that the user has decided to transition. The arrows thickness is proportional to it's weight. Authors in this cluster exhibit *steady* research interest in their primary  $A \circ I$   $D_1$ . Some authors start dominantly contributing in their secondary  $A \circ I$ ,  $D_2$  in State 4. Though, they return to spending around half of their effort in  $D_1$  in State 5.

ing Databases and Data Mining, also making significant contributions in Machine Learning and Bioinformatics lately. Another professor who started in Databases, Jaideep Srivastava of the University of Maryland, evolved on to research distributed implementation of databases, and also data mining and AI related research simultaneously.

Evolving: These researchers have one dominant area of interest (AoI) in each state which *changes* with time. Their dominant area of interest (AoI) *evolves* from  $D_1$  (72%) in state 1 to  $D_2$  (66%) in state 4 to  $D_3$ (43%) in state 5. Even though their AoI shifts across stages, in any given stage, they remain focused on one area and do not publish much in other areas. James Foley, professor in Georgia Tech, started in Computer Graphics and later switched to research on user-computer interfaces and recently, User Modeling. Natural Language Processing (NLP) expert Daniel Jurafsky at Stanford University, also steadily moved from pure NLP based research problems to Speech processing, and later to Machine Learning (ML). Also note, for Jurafsky, this evolution can be attributed to the broader shift of using sophisticated ML models to solve NLP problems.

Diffuse: Authors of this archetype stay focussed in one dominant area in each stage; while in the last stage their research interests are diffused. Authors publish considerably in one dominant area in first 3 stages;  $D_1$  (state 1, 3) to  $D_2$  (state 4). In the last state, which lasts around 4 years, the authors are infrequently publishing (less than 3 papers a year) in new subfields accounting for 74% of their publications. Hence, these authors have diffused research interests after they gain experience. Gerhard Weikum, professor at MPI Germany started in Databases area made a brief transition to Information Retrieval work and later started publishing in Machine Learning and Data Mining fields too. These area

evolutions are more natural transitions as they are highly interrelated which explains contributions in all fields. Anind Dey is a professor at Carnegie Mellon University who initially worked on sensor technology and then switched to Web mining and Human Computing related research problems.

## 6.2 Archetype variations across Gender

In this section, we analyze the evolution of research interests of male and female researchers. To this end, we manually annotate gender of all current and emeritus professors in top 50 Computer Science Universities as reported by U.S. News & World Report<sup>6</sup>. We consider only current and emeritus *Full* Professors as they typically have 15 or more years of publication history. This results in a total of 1084 authors, 127 of whom are women.

While researchers from both genders in the same archetype will traverse the same set of stages, they may differ in *how* they transition  $\tau^c$ , and at *which* stage they start  $\pi^c$ . For this analysis, we estimate separate model parameters for female  $\lambda_f^c$  and male  $\lambda_m^c$  researchers for each archetype c.

| Gender | Steady  | Diverse | Evolving | Diffuse |
|--------|---------|---------|----------|---------|
| Male   | 2.10*** | 2.63**  | 1.15     | 1.10    |
| Female | 1.80*** | 1.64**  | 1.60***  | 1.38*** |

Table 6: Likelihood ratio for academics across archetypes. It measures odds of a researcher being better explained by model for their gender than by model for the other gender.

$$* = p < .05, ** = p < .01, *** = p < .001$$

<sup>&</sup>lt;sup>6</sup>bit.ly/usnews-cs

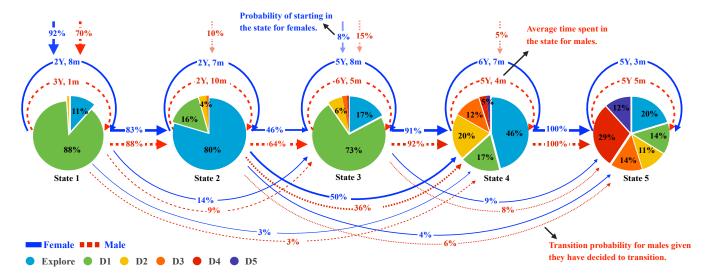


Figure 3: Gender wise representation of trajectory for researchers belonging to the *diverse* archetype in the Academic Dataset. The transitions in blue denote transition probabilities of female professors in the archetype while those in red represents probabilities for their male counterparts. Men start their career from later evolved stages while women make long term state transitions.

To quantify the difference between two models  $(\lambda_f^c, \lambda_m^c)$  for archetype c, we compute their *likelihood ratio*. Likelihood ratio  $R_f^c$  of female researchers in archetype c is:

$$R_f^c = \exp\left(\frac{1}{|N_f^c|} \sum_{i \in N_f} \log \frac{P(X_i | \lambda_f^c)}{P(X_i | \lambda_m^c)}\right)$$
(6)

where  $N_f^c$  represents all female researchers in c-th archetype. The equation simplifies to say that  $\log R_f^c$  is the average difference between log likelihoods of a trajectory of a female researcher generated from model for female researchers with male model of the same archetype. Thus, for instance value of  $R_f^c=2$  denotes that female researchers are twice more likely to be generated by the model of their own gender than of the opposite gender. We compute a similar ratio  $R_m^c$  for men

Table 6 shows the likelihood ratio and *p*-value for the paired-sample t-test Goulden (1949) between the likelihood values. Since most of the values are statistically significant, all researchers are better explained by the model for their gender, than by the model for the opposite gender. Male researchers are distinct for the steady and diverse archetypes, but not for the evolving and diffuse archetypes. For women, on an average, the effect is larger, with the strongest effects seen for the steady, diverse, and evolving archetypes.

For the sake of brevity, we examine gender difference in only the *diverse* archetype in some detail. Figure 3 shows three interesting variations. First, we observe that women are much more likely to start in state 1 (92%), with a dominant area of interest ( $D_1$ ) than in any other state. In contrast, men start in states 1, 2, 3 and 4, with only 70% starting in state 1. Both men and women skip stages, but women are more likely to skip a stage than men. For example, 50% of women skip stage 3, while only 36% men do. Longer skips of two stages are more rare, and both women and men make

these long skips at the same rate. Finally, there are clear differences between mid-career men and women (states 3, 4): women spend more time *exploring* mid-career (state 4) than men, and mid-career men spend more time in their starting area of interest  $(D_1, \text{ state 3})$  than women.

# 6.3 Grant income variability across Archetypes & Gender

Now, we examine the relationship between variation in the academic trajectories and gender to research grants awarded at different stages of academic career. We extract historical information of grants from the National Science Foundation, a large federal funding agency for Science & Engineering in the United States 7. A PI leads the grant, while a Co-PI collaborates in the grant. We analyze the same subset of CS professors in top-50 US universities as in Section 6.2. We collect information for 1062 professors and manually disambiguate names and identify gender by crossvalidating with the researcher's webpage. Then, we compute the average grant money awarded to a researcher, at each stage in their trajectory. Figure 4, which shows letter-value plots of average grant size awarded as PI's, broken down by archetypes (steady, diverse, evolving or diffuse), stage within an archetype and gender, summarizes our findings.

Additionally, we conducted Kruskal-Wallis Htest Kruskal et al. (1952) to establish statistical significance of differences in grant money across latent states within an archetype. This test affirms that at least one latent state is different from another latent state within an archetype. We then conducted Welch's t-test Welch (1947) between consecutive states to find the exact pair of states which are significantly different. We only tested with consecutive latent states as we are only interested in grant income changes

<sup>&</sup>lt;sup>7</sup>bit.ly/nsfgrants

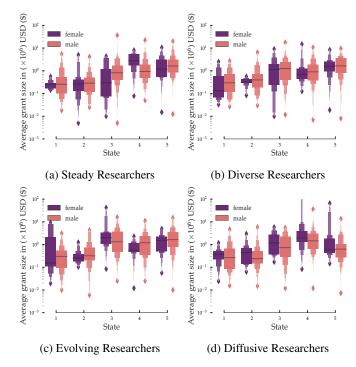


Figure 4: Letter value plots of total grant money awarded by NSF when author is a PI in each stage. In general, Professors get more grant money as they gain experience. Regardless of archetypes, grant income in state 3 is significantly higher from state 2 (p; .01). There are also significant differences across genders within a state of an archetype. For instance, for Evolving archetype, male professors get significantly more income than female professors in state 4 (p; .01).

as author progresses through stages. Table 7 reports the state pairs for each archetype that are statistically different. In the rest of this section, we describe these results in detail.

Regardless of archetypes, we observe that in general authors tend to receive more grant money as they gain experience in Figure 4. On average, across archetypes and gender, PI's receive in state 5, four times the amount of grant money than state 1 (p < .001). Also for researchers across archetypes and across genders, we notice an uptick in grant income in state 3 from state 2 (p < .01 - Table 7). Let us qualitatively examine the steady researchers in detail, by comparing Figure 4 with Figure 2. State 2 in Figure 2 shows the researchers exploring different topics, whereas in state 3, they are spending a significant part of their time on their main domain  $D_1$ . Also notice that 36% of the researchers never visit state 2 - 27% skip state 2, and 9% of the researchers start in state 3. Since state 1 typically represents the time spent by the researchers in their PhD, and with 74% time spent in an explore stage in state 2, it is not surprising that we see limited grant income in their first two states. State 3, perhaps reflects a sustained focus on their domain  $D_1$ , and this pays off in terms of grant income. Similar qualitative arguments follow for the other archetypes.

The grant trajectories over states is different for each archetype (p < .001). Let us examine statistically different

| Archetype (H-test)                      | State Pair (t-test)           |
|---|-------------------------------|
| Steady***                               | State 2 vs 3**                |
| Diverse***                              | State 2 vs 3***               |
|   | State 4 vs 5*                 |
| Evolving***                             | State 2 vs 3*** State 3 vs 4* |
| _,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,, | State 4 vs 5**                |
| Diffuse***                              | State 2 vs 3***               |
| Diffuse                                 | State 4 vs 5*                 |

Table 7: Statistical significance tests for the differences in grant money across latent states within an archetype. Shown are only those tests that are statistically significant. H-testKruskal et al., 1952 confirms that at least one state is different from another state of the archetype; t-testWelch, 1947 was then conducted between each consecutive states within the archetype to determine the differing states.

\* = p < .05, \*\* = p < .01, \*\*\* = p < .001

| Archetype | Latent State (t-test) |
|-----------|-----------------------|
| Steady    | State 1*<br>State 4*  |
| Diverse   | State 2*              |
| Evolving  | State 4** State 5*    |
| Diffuse   | Not significant       |

Table 8: Statistical significance tests Welch, 1947 for the differences in grant money across gender in each state within an archetype. Shown are only those tests that are statistically significant.

$$* = p < .05, ** = p < .01, *** = p < .001$$

state pairs from Table 7 in Figure 4. Steady researchers see a big uptick in their grant income in state 3 and their subsequent grant income is similar in magnitude (p < .01). The grant income for diverse researchers (who have more than one dominant area) increases steadily over states(p < .05). For evolving researchers (who change their dominant area), the grant income rises (state 3, p < .001), falls (state 4, p < .05) and rises (state 5,p < .01), reflecting a degree of unpredictability accompanying changing area of interest. Diffuse researchers, have a pattern similar to steady researchers except in state 5 (p < .05), when the income dips, perhaps due to spending time in too many areas.

To determine differences in grant income across gender, we conducted t-test Welch (1947) for each state within an archetype. Table 8 reports significantly different states within each archetype. We again examine these statistically different states from Table 8 in Figure 4. Evolving women receive significantly *lower* income than evolving men when they *switch to new areas* in state 4 and 5 (p < .05). On the other hand, in our dataset, steady women receive significantly *higher* grant income than steady men when they *switch areas* in state 4 (p < .05). In general, we observe that men show greater grant income variability than do women. The variability is statistically significant (p < .05) during early career in state 1 and 2 for Steady and Diverse re-

searchers respectively. We do not observe significant differences in grant income of male and female Diffusive researchers.

## 6.4 StackOverflow Archetypes

We now describe archetypes learned for the Stack Overflow data. Table 9 depicts the latent states for all archetypes. We label the 4 archetypes discovered as *Expert*, *Seekers*, *Enthusiasts* and *Facilitators*. Posting Comments is the most frequent activity in all archetypes as it is a very low cost activity. Moderator actions and Edits are least favored activities by Stack Overflow users. Most of the users spend initial sessions for *posting questions* (state 1) and significant proportion of their later sessions in *posting answers or comments* (state 6).

Experts users join the community to answer queries or post clarifications or edit answers (state 1-6). They spend at least 68% of their sessions in *posting answers* (state 1 & 3). They rarely ask questions of their own. In communities like Stack Overflow, it is vital to have a dedicated group of experts answering queries for it to be sustainable.

Information Seekers join the community for getting answer to their queries accounting for 69% of activities per session (state 1). They briefly start contributing by posting answers to the community (state 3) but they end up again in *commenting* (state 6) or *asking questions* (state 5).

Enthusiasts start by asking questions and posting comments (state 1). They, then, start answering questions and commenting on other answers (state 2). They briefly stay (4 sessions) in edit state (state 3) but end up migrating to either commenting again (state 4) or asking questions and commenting (state 5). We denote them as *Enthusiasts* as they use the platform to post questions while simultaneously answering queries from their acquired knowledge.

Facilitators join for information seeking (state 1) but start posting answers, clarifying and editing in state 2-3. However, later on they take a more subdued approach and only post comments. The reason for this decreased interest is hard to gauge but identifying these users and retaining their interest could be important to sustain the community.

In summary, we identify four archetypes for researchers: steady, diverse, evolving and diffuse. We observe differences in evolution of male and female researchers within the same archetype. When we examine the diverse archetype in detail, we observe that women and men differ in how they start, how they transition and time spent in mid-career. The differences in grant income are salient across states within an archetype. We also observe differences across genders within a stage of an archetype. We also identify archetypes for StackOverflow: experts, information seekers, enthusiasts and facilitators.

#### 7 Limitations

Our proposed model identified insightful archetypes and its variability with gender and grant income of professors. However, it is important to understand certain caveats to the reported findings. First, in terms of the data, the discovered

|         | Experts        | Seekers            | Enthusiasts    | Facilitators       |
|---------|----------------|--------------------|----------------|--------------------|
| State 1 | {14.4S}        | {3.7S}             | {12.6S}        | {15.4S}            |
|         | <b>A</b> (68%) | $\mathbf{Q}$ (69%) | C (42%)        | $\mathbf{Q}(50\%)$ |
|         | C (19%)        | C (19%)            | <b>Q</b> (39%) | C (32%)            |
| State 2 | {17.8S}        | {8.4S}             | {9.5S}         | {26.4S}            |
|         | C (60%)        | C (51%)            | <b>A</b> (49%) | <b>C</b> (44%)     |
|         | A(25%)         | Q(33%)             | C(32%)         | A(28%)             |
|         |                |                    | E(12%)         | E(22%)             |
| State 3 | {9 <b>S</b> }  | $\{2.2S\}$         | {3.7S}         | $\{8.4S\}$         |
|         | A(87%)         | A (84%)            | E(82%)         | A (87%)            |
| State 4 | {12.3S}        | {5.3S}             | {9 <b>S</b> }  | {24.2S}            |
|         | C(82%)         | C(72%)             | C(75%)         | C(68%)             |
|         |                | Q(15%)             | A(10%)         | A (14%)            |
|         |                |                    |                | E(13%)             |
| State 5 | {5S}           | {3.7S}             | $\{10.5S\}$    | {14.3S}            |
|         | E(45%)         | Q(85%)             | Q(40%)         | <b>E</b> (63%)     |
|         | C (28%)        |                    | C (33%)        | C(22%)             |
|         | A(21%)         |                    | E(16%)         |                    |
| State 6 | {11S}          | {7.7S}             | {5S}           | {21S}              |
|         | <b>C</b> (48%) | C(55%)             | <b>A</b> (61%) | <b>C</b> (57%)     |
|         | A (38%)        | <b>Q</b> (23%)     | C (23%)        | E(24%)             |
|         |                | E(11%)             | E(11%)         | A(13%)             |

Table 9: Learned mean vector for each state for four archetypes in the Stack Overflow Dataset. We list the *activities* in sorted order and annotate them with their % contribution in the state. We list main activities (> 11%) for each state. Each state is also labeled with it's average number of sessions. The labels reflect our own interpretation of the user behavior.

archetypes for academics are for the top researchers in their field (we pick *influential* researchers in each of the 35 research subdomains). Thus, our archetypes do not reflect all computer scientists engaged in research. Also our grant analysis was focused on professors from only top-50 Computer Science schools. In our current study, we collected grant history from data publicly available by NSF. The funding analysis can be extended by collecting data from other possible funding sources like National Institute of Health (NIH), gifts and professor's salary. Hence, we believe that our study is a first step in understanding differences in research conducting behavior of academics and its effect on their income.

Second, as with all inductive models, our qualitative results depend on chosen model. Recently, Deep Neural Networks, especially Recurrent Neural Networks have been proposed to model time series data. There has also been great interest in building interpretable models (Lakkaraju et al., 2016; Ribeiro et al., 2016). However, still most of these models remain as a black box and do not provide meaningful results.

Third, in our current version of the model, we do not consider the effect of collaborations, or the role of conferences where researchers publish, and where they may pick up on normative behavior (e.g. areas in which to work) on the discovered archetype. In future work, we plan to understand the role of community interaction on archetypes and address these limitations. Another interesting research direction is to explore correlation of change in research behavior with career transitions and author's citation count.

## 8 Conclusion

In this paper, we aimed to discover archetypical behavioral patterns for individuals in large social networks. The observation that despite near limitless variation in behavior, the change in behavior exhibits regularities, motivated our research. We introduced a novel Gaussian Hidden Markov Model Cluster (G-HMM) to identify archetypes and evolutionary patterns within each archetype. We chose to work with G-HMM's since they allow for : near limitless variation; constraints how individuals can evolve; different evolutionary rates and are parsimonious.

We identified four archetypes for computer scientists: steady, diverse, evolving and diffuse and showed examples of computer scientists from different sub-fields that share the same archetype. We analyzed full professors from the top 50 CS departments to understand gender differences within archetypes. Women and men differ within an archetype (e.g. diverse) in where they start, rate of transition and research interests during mid-career. We further analyzed grant income of these professors to understand the effect of gender and archetype on income. The differences in income are salient across states within an archetype. There also exist significant differences across genders within a state of an archetype. For StackOverflow, discovered archetypes could be labeled as: *Experts*, *Seekers*, *Enthusiasts* and *Facilitators*. We showed strong quantitative results with competing baselines for future activity prediction and perplexity.

#### References

Adamic, Lada A. et al. 2008. Knowledge Sharing and Yahoo Answers: Everyone Knows Something. In. WWW'08.

Angeletou, Sofia et al. 2011. Modelling and Analysis of User Behaviour in Online Communities. In *The Semantic Web-ISWC'11*.

Benevenuto, Fabricio et al. 2009. Characterizing User Behavior in Online Social Networks. In. IMC'09.

Bicego, Manuele et al. 2003. Similarity-based Clustering of Sequences Using Hidden Markov Models. In *Proceedings of the 3rd International Conference on Machine Learning and Data Mining in Pattern Recognition*. Springer-Verlag, 86–95.

Biryukov, Maria and Dong, Cailing. 2010. Analysis of Computer Science Communities Based on DBLP. In *Proceedings* of the 14th European Conference on Research and Advanced Technology for Digital Libraries. Springer-Verlag.

Chakraborty, Tanmoy and Nandi, Subrata. 2018. Universal Trajectories of Scientific Success. In *Knowl. Inf. Syst.* 54.2, 487–509.

Coviello, Emanuele et al. 2014. Clustering Hidden Markov Models with Variational HEM. In *J. Mach. Learn. Res.* 15, 697–747.

Danescu-Niculescu-Mizil, Cristian et al. 2013. No Country for Old Members: User Lifecycle and Linguistic Change in Online Communities. In *Proceedings of the 22Nd International Conference on World Wide Web*. ACM.

Dempster, A. P. et al. 1977. Maximum likelihood from incomplete data via the EM algorithm. In *Journal of the Royal Statistical Society: Series B* 39, 1–38.

Furtado, Adabriand et al. 2013. Contributor Profiles, Their Dynamics, and Their Importance in Five Q&a Sites. In. CSCW'13.

Ghassempour, Shima et al. 2014. Clustering Multivariate Time Series Using Hidden Markov Models. In *International journal of environmental research and public health*.

Goulden, Cyril H et al. 1949. Methods of statistical analysis. In *Methods of statistical analysis*.

Hadgu, Asmelash Teka and Jaschke, Robert. 2014. Identifying and Analyzing Researchers on Twitter. In *Proceedings of the 2014 ACM Conference on Web Science*. WebSci '14. Bloomington, Indiana, USA: ACM, 23–32.

Juang, B.H. and R. Rabiner, Lawrence. 1985. A Probabilistic Distance Measure for Hidden Markov Models. In 64.

Kahn, Shulamit. 1993. Gender Differences in Academic Career Paths of Economists. In *The American Economic Review* 83.2, 52–56.

Knab, Bernhard et al. 2003. "Model-Based Clustering With Hidden Markov Models and its Application to Financial Time-Series Data". In. *Between Data Science and Applied Data Analysis. Studies in Classification, Data Analysis, and Knowledge Organization.* Ed. by Schader, Martin et al. Springer Berlin Heidelberg.

Kruskal, William H. and Wallis, W. Allen. 1952. Use of Ranks in One-Criterion Variance Analysis. In. Vol. 47. 260. Taylor & Francis, 583–621.

Kullback, S. and Leibler, R. A. 1951. On information and sufficiency. In *Ann. Math. Statistics* 22, 79–86.

Lakkaraju, Himabindu et al. 2016. Interpretable Decision Sets: A Joint Framework for Description and Prediction. In *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining.* KDD '16. ACM, 1675–1684.

Linek, Stephanie et al. 2017. Its All About Information? The Following Behaviour of Professors and PhD Students in Computer Science on Twitter. In *The Journal of Web Science* 3.1, 1–15.

Lipton, Zachary Chase. 2016. The Mythos of Model Interpretability. In *CoRR* abs/1606.03490.

Liu, Yong et al. 2014. CHI 1994-2013: mapping two decades of intellectual progress through co-word analysis. In *Proceedings of the 32nd annual ACM conference on Human factors in computing systems*. ACM, 3553–3562.

Ltkepohl, Helmut. 2007. *New Introduction to Multiple Time Series Analysis*. Springer Publishing Company.

Maia, Marcelo et al. 2008. Identifying User Behavior in Online Social Networks. In. SocialNets'08.

Mamykina, Lena et al. 2011. Design Lessons from the Fastest Q&a Site in the West. In. CHI'11.

McAuley, Julian John and Leskovec, Jure. 2013. From Amateurs to Connoisseurs: Modeling the Evolution of User Expertise Through Online Reviews. In *Proceedings of the 22Nd International Conference on World Wide Web*. ACM.

Mehrazar, Maryam et al. 2018. Can We Count on Social Media Metrics?: First Insights into the Active Scholarly Use of Social Media. In *WebSci*.

Narang, Kanika et al. 2017. Large-Scale Analysis of Email Search and Organizational Strategies. In *Proceedings of the 2017 Conference on Conference Human Information Interaction and Retrieval*. New York, NY, USA: ACM.

Rabiner, Lawrence R. 1990. Readings in Speech Recognition. In. Ed. by Waibel, Alex and Lee, Kai-Fu. Morgan Kaufmann Publishers Inc. Chap. A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition, 267–296.

Ribeiro, Marco Tulio et al. 2016. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD '16. San Francisco, California, USA: ACM, 1135–1144.

Safavi, Tara et al. 2018. Career Transitions and Trajectories: A Case Study in Computing. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. KDD '18. ACM, 675–684.

Smyth, Padhraic. 1997. Clustering Sequences with Hidden Markov Models. In *Advances in Neural Information Processing Systems*. MIT Press, 648–654.

*The PageRank Citation Ranking: Bringing Order to the Web.* 1999. Technical Report 1999-66.

Tibshirani, Robert et al. 2001. Estimating the number of clusters in a data set via the gap statistic. In *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 63.2, 411–423.

Ward, Melanie. 2001. The gender salary gap in British academia. In *Applied Economics* 33.13, 1669–1681.

Way, Samuel F. et al. 2016. Gender, Productivity, and Prestige in Computer Science Faculty Hiring Networks. In *Proceedings of the 25th International Conference on World Wide Web*.

Way, Samuel F. et al. 2017. The misleading narrative of the canonical faculty productivity trajectory. In *Proceedings of the National Academy of Sciences* 114.44, E9216–E9223.

Welch, B. L. 1947. THE GENERALIZATION OF STUDENT'S PROBLEM WHEN SEVERAL DIFFERENT POPULATION VARLANCES ARE INVOLVED. In *Biometrika* 34.1-2, 28–35.

Yang, Jaewon and Leskovec, Jure. 2011. Patterns of Temporal Variation in Online Media. In *Proceedings of the Fourth ACM International Conference on Web Search and Data Mining*. New York, NY, USA: ACM.

Yang, Jaewon et al. 2014. Finding Progression Stages in Time-evolving Event Sequences. In *Proceedings of the 23rd International Conference on World Wide Web*. WWW'14. ACM, 783–794.