

The Utility of Skills in Online Labor Markets

Completed Research Paper

Marios Kokkodis

NYU Stern School of Business
44 W 4th Street, New York, NY, USA
mkokkodi@stern.nyu.edu

Panagiotis G. Ipeirotis

NYU Stern School of Business
44 W 4th Street, New York, NY, USA
panos@stern.nyu.edu

Abstract

In this work we define the utility of having a certain skill in an (OLM), and we propose that this utility is strongly correlated with the level of expertise of a given worker. However, the actual level of expertise for a given skill and a given worker is both latent and dynamic. What is observable is a series of characteristics that are intuitively correlated with the level of expertise of a given skill. We propose to build a Hidden Markov Model (HMM), which estimates the latent and dynamic levels of expertise, based on the observed characteristics. We build and evaluate our approaches on a unique transactional dataset from oDesk.com. Finally, we estimate the utility of a series of skills and discuss how certain skills (e.g. ‘editing’) provide a higher expected payoff once a person masters them over others (e.g. ‘microsoftexcel’).

Keywords: Online labor markets, Utility of skills, HMM, Worker’s expertise, Data analysis, Empirical analysis,

Introduction

Online labor marketplaces (OLMs), such as oDesk.com, Elance.com, and Worker.com, allow employers to connect with workers around the globe, to accomplish tasks that span diverse categories such as: web development, writing and translation, accounting, etc. These marketplaces are growing fast and the worker annual earnings are expected to grow from \$1 billion in 2012 to \$10 billion by 2020 (Agrawal et al. 2013). A typical scenario in these workplaces starts with an employer posting a job (which can be a short task, a long-term project, or both). Multiple workers start bidding for the job and eventually the employer chooses to hire one (or several) of them. Finally, the hired worker(s) completes the task and receives a payment.

As with offline workplaces, work in OLMs is an ‘experience good’, meaning it is practically impossible to know the quality of the task outcome (or even the expertise of a worker in a given skill) in advance (Nelson 1970). A key solution to resolve this uncertainty is the use of online reputation systems, which provide signals about the past performance of workers (Danescu-Niculescu-Mizil et al. 2009; Dellarocas 2003; Liu et al. 2008; Lu et al. 2010). However, reputation systems in these marketplaces fail to capture the actual worker quality; they tend to be highly skewed towards high values, and they eventually become uninformative (Hu et al. 2009).

The workers’ value in OLMs resides in a combination of both observable and latent characteristics. The observed characteristics usually include a list of skills, the educational background, the work history, and

the certifications of the applicant. The latent characteristics include the worker's expertise and true ability in the listed qualifications. Very similar to an offline setting, the demand and supply distributions (and as a result, the expected payoff) of each worker, with a given set of skills and a given level of expertise, are very heterogeneous; for example, a Java 'expert' might have a very different expected payoff than an 'expert' in customer service support. Similarly, a c# 'expert' might have a higher expected payoff than a c# 'beginner', etc.

This observation leads to two fundamental questions: (1) how can we estimate the latent level of expertise of a given worker and a given skill? (2) how can we quantify the value of a skill in an online labor marketplace, and how is this value correlated with the level of expertise of a worker?

In this work, we focus on addressing the abovementioned two questions. We first propose that the utility of each skill is strongly correlated with the level of expertise of a specific worker, an assumption that also holds for offline markets. Based on the intuition that an experienced worker should – on expectation – receive higher compensations than a beginner, we formally define the conditional utility of a skill given a certain level of expertise. However, the actual level of expertise for a given skill and a given worker is latent (not directly observable) and dynamic (evolves over time). To overcome this, we use a series of directly observed characteristics that are intuitively related to the level of expertise of a given skill. Based on these observable characteristics, we propose to build a Hidden Markov Model (HMM), which estimates the latent and dynamic levels of expertise of each worker, in a given skill.

For the deployment and evaluation of our methodology, we use a unique transactional dataset of 1.5 million job applications from the biggest (in terms of worker earnings) online labor market, oDesk.com. We compare our proposed approach with two baselines, and show that our framework performs significantly better, predicting workers' levels of expertise with an exceptionally high accuracy rate. Once we compute each worker's level of expertise, we estimate the conditional utility of each skill in our dataset. Lastly, we discuss how certain skills appear to have a much higher expected compensation, once someone masters them over others.

Our study is the first to quantify the value of a series of skills. We strongly believe that both online labor marketplaces and workers can benefit significantly from this analysis. Finally, we acknowledge that this work is the first step in answering a more important question: given a set of skills, with a certain level of expertise, what should the next steps of a utility-optimizing worker in (but also beyond) an online labor marketplace be?

Background

Our work is novel in the sense that we are the first to define and study the utility of a skill in an OLM setting. Broadly related previous work has focused on different aspects of OLMs and paid crowdsourcing, as well as on skills assessment and 'expert' findings. We briefly discuss this work in the following paragraphs.

OLMs & Paid Crowdsourcing

Past research in Online Labor Markets (OLMs) spans across a variety of problems. John J. Horton (Horton 2010) explores market creators' choices of price structure, price level and investment in platforms. He further discusses possible productivity and welfare implications that these markets can have. Horton et al. (Horton et al. 2011) present a model of workers supplying labor to paid crowdsourcing. They find that workers work less when the pay is lower, but they do not work less when the task is more time consuming.

A different stream of work studies the validity of behavioral experiments in these markets. Rand, D. G. (Rand 2012) discusses how Mechanical Turk can be used as a tool for behavioral experimentation. Similarly, Horton et al. (Horton et al. 2011) show that online experiments can be just as valid - both internally and externally- as laboratory field experiments. In addition, Berinsky et al. (Berinsky et al. 2012) assess the internal and external validity of experiments performed using Mechanical Turk.

On a different direction, a lot of work focuses on incentivizing workers as well as finding ways to manage the quality of their outcomes. In particular, Shaw et al. (Shaw et al. 2011) ran an experiment on Mechanical Turk to measure the effectiveness of social and financial incentive schemes on outcome quality. One of their main findings was that when workers had to think about responses of their peers, combined with financial incentives, they provided higher quality results. Mason et al. (Mason et al. 2010) studied the effect of compensation on performance in the context of two experiments conducted on AMT, and found that increased financial incentives increase the quantity but not the quality of work performed by participants. They also observed an anchoring effect, where workers who were paid more also perceived the value of their work to be greater, and thus were no more motivated than workers who were paid less. Furthermore, Chandler et al. (Chandler et al. 2011) ran a natural field experiment on Amazon Mechanical Turk and found evidence that the user interface and the cognitive biases of the workers play an important role in OLMs. Sheng et al. (Sheng et al. 2008) studied repeated-labeling strategies in OLMs. Two of their main findings were that (1) repeated-labeling can improve label quality but not always, and (2) that when processing unlabeled data is not free, even the simple strategy of labeling everything multiple times can give considerable advantage. Ipeirotis et al. (Ipeirotis et al. 2010) presented algorithms that separate workers' ability errors from errors caused by workers' biases. Finally Ipeirotis et al. (Ipeirotis et al. 2011) discussed the need of standardization of basic building block tasks that could make crowdsourcing more scalable.

In 2003, Snir et al. (Snir et al. 2003) studied costly bidding in online markets and found that higher value projects attract significantly more bids, with lower quality, and that greater number of bids raise the cost to all participants, due to costly bidding and bid evaluation. A. Palais (Pallais 2013) ran an experiment on the oDesk.com platform to study the cold-start problem (i.e. hiring inexperienced workers) in an OLM. Her experiment showed that both hiring workers and providing more detailed evaluations substantially improves workers' subsequent employment outcomes. Finally, Kokkodis et al. (Kokkodis et al. 2013) studied what happens when a worker transitions between different task categories. In their study, they provided a static Bayesian approach for quantifying reputation transferability across different categories in OLMs.

Skills Assessment & Expert Search

In the past, a lot of work has been done that deals with skills assessment and 'expert' search. Hambleton (Hambleton et al. 1991) described the fundamental concepts of the Item Response Theory (IRT), a theory widely used in Computer Adaptive Testing. Desmarais (Desmarais et al. 1995) proposed the creation of a network that captures implication relations among knowledge units (KU), and from this network one can learn someone's knowledge state with a limited number of observations or questions. Jurczyk and Agichtein (Jurczyk et al. 2007) presented an analysis of the link structure of a general-purpose question answering (Q&A) community to discover authoritative users. Similarly, Zhang et al. (Zhang et al. 2007) tested a set of network-based ranking algorithms, including PageRank (Brin et al. 1998) and HITS (Kleinberg 1999), on a Java forum, in order to identify users with high expertise. Bouguessa et al. (Bouguessa et al. 2008) took this analysis one step further: they studied the problem of determining the amount of users that one should choose as authoritative from a ranked list. Balog and Rijke (Balog et al. 2007) proposed a technique for automatically determining the expertise profile of a person by including information from the person's areas of skills and from his/her social profile. Macdonald and Ounis (Macdonald et al. 2008) proposed an approach for predicting and ranking candidate expertise with respect to a query. Finally, Petkova and Croft (Petkova et al.) proposed a general approach for representing the knowledge of a potential expert as a mixture of language models from associated documents.

On a different note, Lazear (Lazear) proposed a 'skill-weight' approach to represent firm-specific human capital. Goes et al. (Goes et al.) tried to identify factors that motivate sellers to seek certifications in OLMs, and found that certification status can negatively impact some sellers' abilities to obtain contracts, even when certification exams are free. Varshney et al. (Varshney et al. 2013) proposed a new collaborative filtering approach for skills assessment prediction and recommendation, that accounts for not only HR-collected data, but also for mined data from online technical communities. Handel (Handel 2003) studied skills mismatch in the labor market. The main concern of his study is the absence of a standardized method of collecting information about the actual skills content of jobs, which is a significant obstacle in answering whether or not job demands are actually exceeding worker's capacities. Next, Cunha et al.

(Cunha et al. 2010) formulated multistage production functions for children's cognitive and non-cognitive skills, and they found that substitutability decreases in later stages of the life cycle of the production of cognitive skills, while it remains constant in the production of non-cognitive skills. Finally, Saito et al. (Saito et al. 2014) proposed a framework for developing micro-tasking skills, which consisted of three modules: (1) a tutorial producer, (2) a task dispatcher, and (3) a feedback visualizer.

Utility of a skill in an Online Labor Market

As we discussed in the introduction, the value of each worker in an OLM is strongly connected to the worker's set of skills and the respective level of expertise in these skills. How can we quantify the value of a skill? In this section, we first define the worker's utility in completing a task. Then, by using this definition, we derive the utility of each individual skill, given a certain level of expertise.

The utility of a task

Every individual worker who joins an OLM focuses on completing tasks that are relevant, but also on tasks that have the highest possible payoff. As a result, the worker's gained utility from completing a task correlates with the received hourly wage (W), as well as the hourly cost of completing the task. Based on this observation, we define the utility of completing a task to be proportional to the difference between the average hourly wage of a task \bar{W} and the average cost of the worker's hourly effort \bar{C} :

$$U_t = \bar{W} \bar{t} - \bar{C} \bar{t} \propto \bar{W} - \bar{C} = U_h$$

In the previous equation, \bar{t} is the average duration (in hours) of a task in the marketplace, and U_h is the resulting hourly task utility.

The conditional utility of a skill

The average expected utility per task is an oversimplification; online labor marketplaces are highly heterogeneous in terms of skills, job categories, workers' abilities, etc. As a result, U_h does not contain much information. To clarify this, consider a scenario where we are interested in the expected utility of completing a task that requires knowledge of *.net*, as well as in the expected utility of completing a task that requires *blog-commenting*. On oDesk.com, the average hourly wage of a *.net* task is \$18.2, while the average hourly wage of a *blog-commenting* task is \$4.6. It is clear that U_h cannot represent both (if any) of these tasks.

To delve deeper into estimating a more accurate expected utility of a completed task, in this work we focus on the value of skills. Based on intuition, we propose that the utility of a skill is not independent of the worker's latent abilities. For example, we expect that an 'expert' in essay editing will charge significantly more than a 'beginner' in editing¹. In addition, the utility of each skill is also correlated with the demand for that skill in the OLM. Having these in mind, we define the conditional utility of a skill.

Definition (Conditional Utility of a skill): Given a skill (s), a worker with a level of expertise in that particular skill ($E_s = e$) and the number of jobs in the marketplace that require this skill (D_s), we define the conditional utility of this skill given the level of expertise (e) as follows:

¹ In fact, the specific example is confirmed empirically (see Figure 8).

$$\begin{aligned}
U_{\{h|E_s=e, D_s\}} &= \frac{1}{|\mathbf{I}|} \sum_{i \in \mathbf{I}} (\left[\bar{W}_{\{E_{i,s}=e\}} - \bar{C}_{\{E_{i,s}=e\}} \right] - \frac{1}{|E_s|-1} \sum_{\{e' \in \{E_{i,s}-e\}\}} \left[\bar{W}_{\{E_{i,s}=e'\}} - \bar{C}_{\{E_{i,s}=e'\}} \right]) \\
&= \frac{1}{|\mathbf{I}|} \sum_{i \in \mathbf{I}} (\bar{W}_{\{E_{i,s}=e\}} - \frac{1}{|E_s|-1} \sum_{\{e' \in \{E_{i,s}-e\}\}} \left[\bar{W}_{\{E_{i,s}=e'\}} \right]) \quad (1)
\end{aligned}$$

In the above equation, $\bar{W}_{\{E_{i,s}=e\}} = \frac{\sum_{o \in O_{E_{i,s}, D_s}} w_o}{|O_{E_{i,s}, D_s}|}$ is the average hourly wage on openings (o) that require the given skill s and completed by worker i with a level of expertise $E_{i,s} = e$. Furthermore, \mathbf{I} is the set of workers that have completed tasks with the given skill s in all available levels of expertise. Finally, $\bar{C}_{\{E_{i,s}=e'\}}$ is the average **hourly** cost of effort of worker i , with a level of expertise s .

Intuitively, the conditional utility of a skill represents the expected increase in hourly wage that is associated with each level of expertise in the given market. For instance, in the simplest case where $|E_s| = 2, e \in \{\text{beginner, expert}\}$, $U_{\{h|E_s=\text{beginner}, D_s\}}$ will represent the utility of an entry-level worker in the marketplace, whereas $U_{\{h|E_s=\text{expert}, D_s\}}$ will represent the utility of an expert worker in the marketplace. The extension to K -levels of expertise is straightforward.

Thus far, we have not discussed the hourly cost of effort that appears (and disappears) in Equation 1. The key assumption we make is that the individual cost of effort $\bar{C}_{\{E_{i,s}=e'\}}$ remains unchanged when a worker changes levels of expertise. To understand the reasoning, consider a case where we have three levels of expertise: ‘beginner’, ‘knowledgeable’, and ‘expert’. James is a worker who has started learning Python, so he is now a beginner. His average hourly cost of effort is $\bar{C}_{\{E_{\text{james, python}}=\text{beginner}\}}$. After a few months, James becomes knowledgeable in Python. His average hourly cost of effort now becomes $\bar{C}_{\{E_{\text{james, python}}=\text{knowledgeable}\}}$. Since on average, James puts in the same effort and has the same productivity in both states, we assume that the two costs are approximately equal. The same argument holds when James becomes an expert. Note that by considering the per person and per hour average cost of effort, we avoid the common problem of having different types of workers based on their productivity/effort (e.g. (Spence 1973)).

Estimating workers’ level of expertise

In this section, we present our proposed approach to identify a worker’s level of expertise in a given skill. We begin by presenting the set of observable characteristics that are associated with the level of expertise in question. We then propose a framework that uses all these observable signals and predicts a given worker’s latent level of expertise in a given skill.

The clues

In the previous section, we assumed that the level of expertise for a given skill and worker is known. In practice, the actual knowledge of a skill is latent, and not directly observable. However, in an online labor market, we observe multiple signals of workers’ expertise. Assume, for example, that we are dealing with an opening that requires Java and SQL, for which we have three applicants: Chris, Marina, and Praveen. Chris is lazy and does not list his skills and expertise on his profile. However in the past, he has completed a series of tasks that required Java and SQL, and he has received great feedback scores for these. On the other hand, Marina is very meticulous and lists all her skills on her profile, which include Java and SQL. She also appears to have repeatedly completed tasks (that required these two skills) under the same employer. Finally, Praveen just joined the market, and has no previously completed tasks. He has listed the two skills on his profile, and has also succeeded in the tests related to those skills. Taking into account this information, which one of the three candidates should we pick for the task? In other words, which candidate has the highest expertise for the task?

The point of this hypothetical example is that there is no objective way to assess the level of expertise of an applicant. However, there are multiple observable clues that provide enough information to shape expectations for someone's expertise.

Feature	Symbol	Domain
Certification	C	[0,1]
Feedback Score	F	[0,1]
Hiring Rate	H	[0,1]
Rehire	R	{0,1}
Wage	W	[1,50]
Mentioned	M	{0,1}

Table 1: Observed clues associated with the worker's level of expertise in a given skill.

The first such clue that we observe is whether or not a worker has passed a certification test associated with the skill at hand. oDesk.com, for example, provides a list of tests that are associated with different skills, varying from Social Media Marketing to Microsoft Excel and SQL². The workers can choose whether or not to take a test for a certain skill. If a worker chooses to take a test, the worker's performance is then listed on the worker's profile³ as a percentile score.

The second clue, which provides information about a worker's expertise, is the worker's past performance on tasks that required the skill at hand. As discussed before, in OLMs, every time a task is completed, workers receive an integer feedback rating between 1 and 5, for six dimensions: communication, cooperation, deadlines, quality, skills, and availability. It is natural to assume that people – who repeatedly receive good feedback scores for tasks that require a certain skill – have a good level of expertise in that skill. Now one might assume that by itself, a feedback score is a sufficient clue for assessing someone's expertise. This is not particularly true in OLMs: feedback scores appear to be highly skewed towards high values (Hu et al. 2009), and hence become uninformative. In particular, on the oDesk.com platform, the median is 5 (even though the highest possible rating is 5). Another possibility that explains this is the 'customer death' phenomenon (Jerath et al. 2011): workers that receive bad feedback do not survive long enough in the marketplace, because it becomes practically impossible for them to get hired again. Therefore, they decide to either create a new account or just abandon the workplace. The resulting marketplace becomes filled with workers who have really high ratings – hence the highly skewed distributions.

The third clue that might include information about a worker's expertise in a skill is the worker's hiring rate. We define this rate as the fraction of applications that lead to a hire. The assumption here is based on the intuition that a worker with high hiring rates, for jobs that require a specific skill, must have a high level of expertise in that skill.

The fourth clue we consider is whether or not the task at hand is a 'rehire', i.e. the employer has hired the same worker before, for a task that required the skill at hand. Hiring someone a second time is the purest evidential signal that the worker did a great job, and as a result we can infer that the worker is an 'expert' in the given skill.

² For a full list of the available tests, see: <https://www.odesk.com/tests>.

³ Workers can decide to hide their test performance on their profile. A thorough analysis of workers' decisions on oDesk.com is presented by Ipeirotis: www.behind-the-enemy-lines.com/2013/10/badges-and-lake-wobegon-effect.html.

The fifth clue we consider is the worker's compensation wage for the previously completed tasks that required a certain skill. The premise here is that the higher the expertise of a worker in a specific skill, the higher the worker's hourly wage should be. This intuition is based on the correlation between compensation and outcome quality in offline workplaces; the higher the compensation wage of an expert, the higher is the expectation of the outcome quality.

Finally, we include a binary variable that captures whether or not workers list the skill at hand on their profiles, at the time of application. The assumption here is that workers who have some expertise in a skill are also comfortable enough to list that skill on their profiles. We present all these signals in Table 1.

Input: Transition Matrix $\mathbf{A}(i,j) = \Pr(e_t = j e_{t-1} = i)$, local evidence $\Pr(\mathbf{X}_t e_t = j)$, initial state $\Pr(e_1 = j)$; 1: $\Pr(e_1 = j \mathbf{X}_1) = \frac{\Pr(\mathbf{X}_1 e_1 = j) \Pr(e_1 = j)}{\Pr(\mathbf{X}_1)}$ 2: for $t = 2 : T$ do 3: $\Pr(e_t = j \mathbf{X}_{1:t-1}) = \sum_i \Pr(e_t = j e_{t-1} = i) \Pr(e_{t-1} = i \mathbf{X}_{1:t-1})$ 4: $\Pr(e_t = j \mathbf{X}_{1:t}) = \frac{\Pr(\mathbf{X}_t e_t = j) \Pr(e_t = j \mathbf{X}_{1:t-1})}{\Pr(\mathbf{X}_t \mathbf{X}_{1:t-1})}$ 5: end for 6: return $\Pr(e_1 = j \mathbf{X}_1), \dots, \Pr(e_T = j \mathbf{X}_{1:T})$
Algorithm 1: Forwards algorithm for our scenario. It predicts the latent state of a given worker for a given skill, given a sequence of observations.

The Framework

Given all these clues, our goal is to estimate the latent level of expertise for each worker in a given skill. To achieve this, we propose to build a Hidden Markov Model (HMM). Our choice relies on the fact that HMM models are Markov processes that assume the existence of underlying latent states, as well as the ability of transitioning between these latent states. These characteristics fit perfectly with the nature of our problem. Recall that in our scenario, we assume that a given skill's expertise is dynamic and evolves over time; Hence transitioning between latent levels of expertise is required.

HMM: We first define a set of discrete, unobserved (latent) states of expertise (E). These states emit with different probability distributions observations in the set $X \in \{m_1, m_2, \dots, m_M\}$. If we consider that a completion of each job is a unit of time, then we expect that as a worker continues to complete tasks of the given skill, the worker will be transitioning between these unobserved states and will be emitting different observations.

Let a sequence of observations of a worker on a given skill at a given time t be $X_{\{1:t\}}$. We can estimate the conditional probability of the worker being at state $e_t, e_t \in E$, $\Pr(e_t|X_{\{1:t\}})$, by using the forwards algorithm presented in Algorithm 1 (Murphy 2012). The algorithm takes as input some prior probability distribution across all the available states, a transition matrix between the available states, as well as some local evidence, which is a vector of the conditional probability distribution over the emitted observations. After the completion of the first task, the algorithm computes the conditional probability of being at each available state, given the observation. As a worker continues to complete new tasks, the algorithm updates the worker's probability of being at each state, by taking into account the input transition matrix between states, as well as the set of observations up to that point. At each iteration, the algorithm returns the most probable state of the worker at hand. Simply put, the algorithm predicts the state of a worker by choosing the state that best explains the worker's sequence of emitted observations.

Up to this point, we assumed that the input parameters for Algorithm 1 are known. In practice, we estimate these parameters by using a version of the Baum-Welch algorithm (Baum et al. 1970). The particular pseudocode that we use is presented in Algorithm 2. The algorithm is a variant of Expectation Maximization (EM): it assumes that we have a total of N workers, each of whom completes $T_i, i \in \{1, \dots, N\}$ number of tasks. In addition, it considers K different levels of expertise, i.e. $E = \{e_1, e_2, \dots, e_K\}$. At the E-step, the algorithm estimates the expected log-likelihood of the old parameter vector (θ^{old}). At the M-step, the algorithm finds the new parameters that maximize the previously estimated log-likelihood. Note that the parameter vector $\theta = [\Pr(e_1), A(j, k), \Pr(X_t|e_t)]'$. The algorithm keeps iterating until the parameter vector θ converges (Murphy 2012).

```

1: repeat
2:   E-step
    1:  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) = \sum_{k=1}^K \mathbb{E}[N_k^1] \log(\Pr(e_1 = k)) + \sum_{j=1}^K \sum_{k=1}^K \mathbb{E}[N_{jk}] \log A(j, k)$ 
        $+ \sum_{i=1}^N \sum_{t=1}^{T_i} \sum_{k=1}^K \left[ \Pr(e_{i,t} = k | \mathbf{X}_{i,t}; \boldsymbol{\theta}^{old}) \log \Pr(\mathbf{X}_{i,t} | e_{i,t} = k) \right]$ 
    2:  $\mathbb{E}[N_k^1] = \sum_{i=1}^N \Pr(e_{i,1} = k | \mathbf{X}_{i,1}; \boldsymbol{\theta}^{old})$ 
    3:  $\mathbb{E}[N_{jk}] = \sum_{i=1}^N \sum_{t=2}^T \Pr(e_{i,t-1=j} e_{i,t} = k | \mathbf{X}_{i,t}; \boldsymbol{\theta}^{old})$ 
4:   M-step
    1:  $\Pr(e_1 = k) = \frac{\mathbb{E}[N_k^1]}{N}$ 
    2:  $A(j, k) = \frac{\mathbb{E}[N_{jk}]}{\sum_l \mathbb{E}[N_{jl}]}$ 
    3:  $\Pr(\mathbf{X}_t = m | e_t = k) = \frac{\sum_{i=1}^N \sum_{t=1}^{T_i} \Pr(e_t = k | \mathbf{X}_{i,1:T_i}; \boldsymbol{\theta}) \mathbb{1}_{x_{i,t}=m}}{\mathbb{E}[N_j]}$ 
4: until  $A(i, j)$ ,  $\Pr(X_t|e_t = j)$  and  $\Pr(e_1 = j)$  converge.
5: return  $A(i, j)$ ,  $\Pr(X_t|e_t = j)$ ,  $\Pr(e_1 = j)$ 
```

Algorithm 2: Baum-Welch algorithm: estimates the parameter vector θ

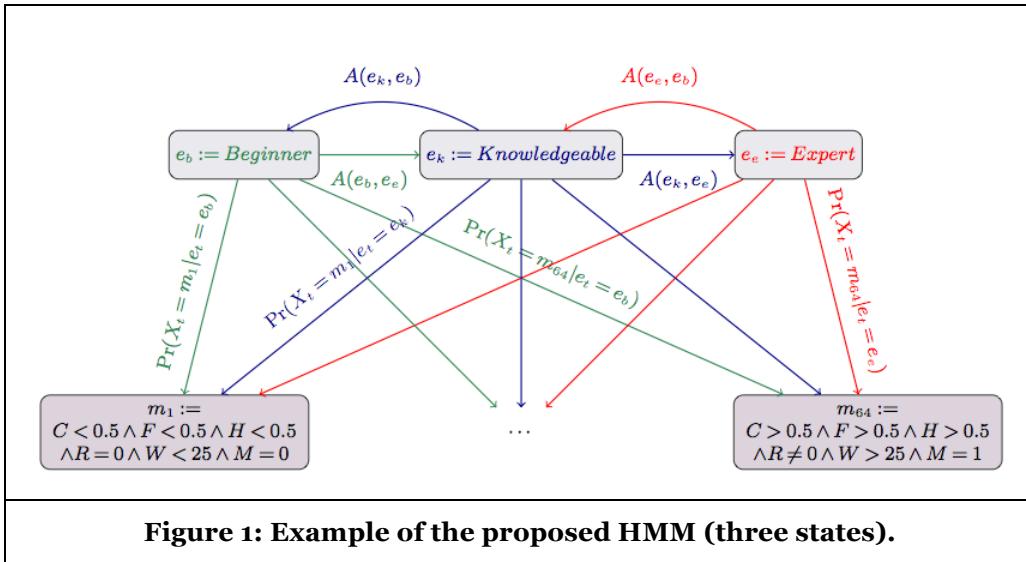


Figure 1: Example of the proposed HMM (three states).

HMM example: Suppose that we consider only three latent states: e_b , which describes an entry-level (beginner) worker; e_k , which describes an intermediate-level worker (knowledgeable); and e_e , which represents an expert worker in a given skill. Recall that we consider six different features to identify the expertise of a worker (see Table 1). We can encode the observations of these features into $M=64$ levels as follows: $m_1 := \{C < 0.5 \wedge F < 0.5 \wedge H < 0.5 \wedge R = 0 \wedge W < 25 \wedge M = 0\}, \dots, m_{64} := \{C > 0.5 \wedge F > 0.5 \wedge H > 0.5 \wedge R = 1 \wedge W > 25 \wedge M = 1\}$. The HMM that describes this scenario is presented in Figure 1.

Skills on oDesk.com

In this section, we provide an empirical analysis of the value of a set of skills that are found on the oDesk.com platform. We start this analysis by describing the data that we use, then we present some implementation details, and finally we evaluate our proposed approaches and compute the estimated conditional utilities of each one of the skills that we study.

Data

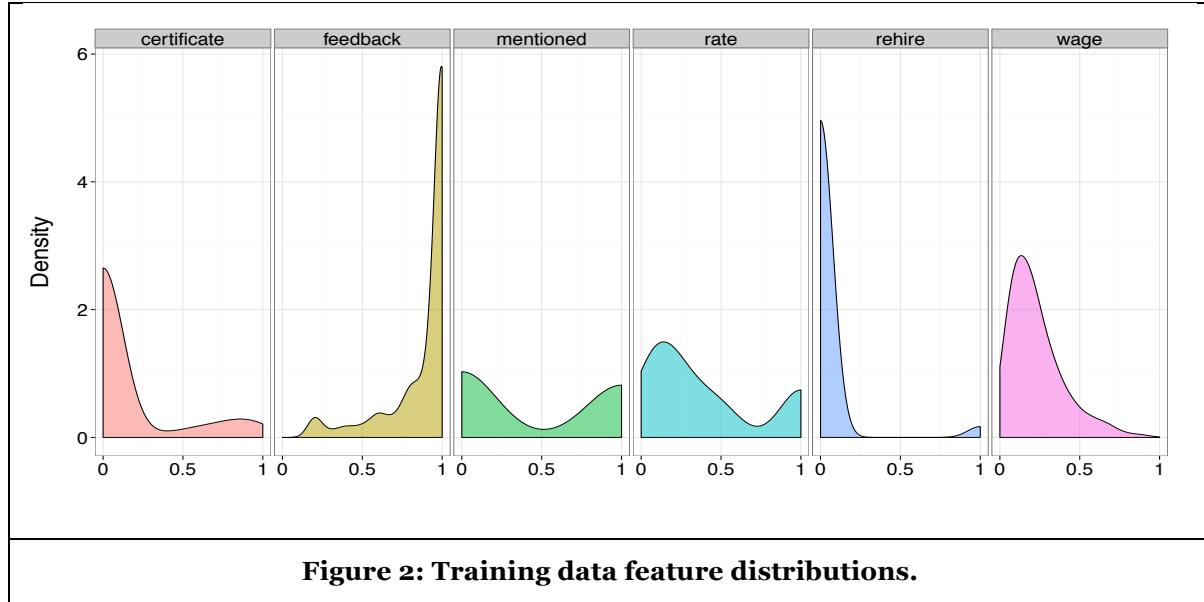
For our study, we use a unique transactional dataset from oDesk.com. oDesk is a global job marketplace with a plethora of tools targeted to businesses that intend to hire and manage remote workers. The company reports more than 500,000 hours of work billed per week, as well as an exponentially growing transaction volume of more than \$300 million USD per year.

The particular dataset we use for this work was collected between September 1st of 2012 and December 31st of 2013, and consists of 1,417,387 applications by 29,309 workers, and 147,555 hiring decisions (completed tasks/received feedback) by 50,516 employers. The analyzed tasks span 5 categories: Software Development, Web Development, Writing & Translation, Sales & Marketing, and Design & Multimedia.

In this work, we use a total of 26 skills from all 5 categories. The selection of the skills was based on two constraints. First, we wanted to limit our analysis to skills for which the oDesk.com platform provides certifications⁴. Second, in order to have sufficient data to train our HMM, we included skills for which we had information about more than 500 workers in our *training set*. Note that this is the number of workers, but not the number of total instances we use to train the HMM (i.e. the number of completed tasks for a given skill).

⁴ In the future, we plan to extend our HMM for skills for which certification tests are not provided by the platform.

We split our data into training (66%) and testing (34%) sets based on workers (i.e. the set of completed tasks of each worker belongs to either the training or the test set). For each opening, we extract the required skills. Then for each hiring decision, we have the snapshot of each worker's profile at the time of application, from which we extract the current state of all (if any) available certifications. In order to compute the current hiring rate of each worker, we keep track of the workers' applications, independent of whether or not they lead to a hiring decision. We also keep track of whether a hiring decision is a rehire. Finally, we include the hourly wage of each task, as well as the feedback score that the worker receives on completion of the task. Once the data preprocessing is complete, we use the training set to build the HMM for each skill. We then use the testing set to evaluate our approach.



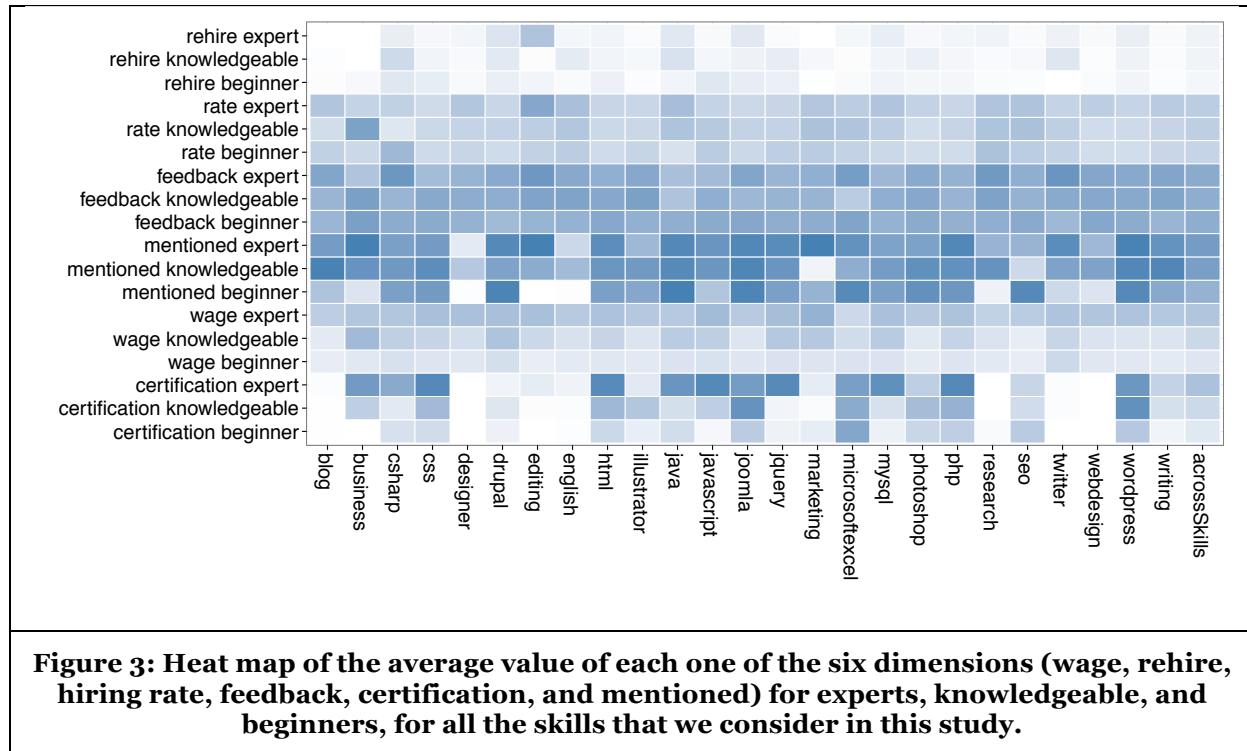
The (smoothed) distributions of all the features we used in our training set are shown in Figure 2. Starting from left to right, we observe that most of our training instances do not have skills certifications. In addition, we observe an increase of the density in high scores (close to one). This increase depicts the tendency of workers to list certifications for which they achieve high scores (as we discussed earlier). The feedback score distribution is as expected: highly skewed towards high scores. The 'mentioned' distribution shows that about 50% of our instances mention the required skills on their profile. From the 'hiring rate' distribution we see that most of the workers accumulate low rates, while a considerable portion of workers have a rate equal to 1. These workers have either been invited to apply or are getting rehired. Next, we see that very few workers are rehired and finally, we observe that the wage distribution (actual wage/50) is skewed to low wages (between \$10 and \$20).

For estimating the level of expertise of a given worker in a given skill, we use an HMM very similar to the one presented in Figure 1. In particular, we use the presented three latent states: beginner, knowledgeable, and expert. For the emitted symbols, instead of splitting every single feature in half, we compute the per feature percentiles; our 64 emitted symbols now become:

$$\begin{aligned}
 m_1 &:= C \in \text{Bottom } 50\% \wedge H \in \text{Bottom } 50\% \wedge F \in \text{Bottom } 50\% \wedge R = 0 \wedge W \in \text{Bottom } 50\% \wedge M = 0 \\
 &\dots \\
 m_{64} &:= C \in \text{Top } 50\% \wedge H \in \text{Top } 50\% \wedge F \in \text{Top } 50\% \wedge R = 1 \wedge W \in \text{Top } 50\% \wedge M = 1
 \end{aligned} \tag{2}$$

Evaluation: Level of expertise

The evaluation of our HMM is not straightforward: our goal is to understand whether the predicted latent states indeed represent the level of expertise of a given worker in a given skill. We propose two evaluating approaches: the first one verifies that our HMM sufficiently separates workers in the three states (beginner, knowledgeable, expert) based on intuition. For the second one, we use the developed HMM as a ‘wage’ predictor, and compare its accuracy with two other baseline predictors.



Evaluation 1, experts/beginners separation: In Figure 3, we present a heat map of the average values of each one of the six dimensions (see Table 1) for experts, knowledgeable, and beginners, for the set of skills that we consider. On the y-axis, we show the Cartesian product of the six dimensions and the three states. Intuitively, if our HMM is correct, we would expect that for every dimension and for every skill, the workers that are predicted to be experts have higher average values (darker shade of blue) than the average values of the workers for which our algorithm predicted to be beginners. On the x-axis is the list of skills we consider in this study. Starting from the bottom horizontal line, we observe that for the majority of the skills, workers that are predicted to be experts have significantly higher *certification* values (darker shade) than those who are predicted to be knowledgeable or beginners. Similarly, the average expert wages appear to be higher than the average knowledgeable and beginner wages, for all skills. For the ‘mentioned’ dimension, the graph is more confusing; still, for the majority of skills, experts have higher values than those that are knowledgeable and beginners. The same applies for the feedback score. As discussed earlier, this should not come as a surprise; the reason is the very skewed feedback score distribution. Furthermore, for the ‘mentioned’ dimension, intuition suggests that workers list their skills independent of how expert they are in these skills. Next, the average hiring rates are either equal between the three states or experts have higher values. Finally, for most of the skills, the average rehire rate is higher for experts than for knowledgeable and beginners.

The bottom line is that the proposed HMM distinguishes experts from the knowledgeable, and knowledgeable workers from beginners, according to intuition. One can further verify this by looking at

the last column of the heat map ('acrossSkills'), which shows the average values of all skills in each one of the proposed dimensions.

Evaluation 2, wage predictor: The second method for evaluating our approach draws on our definition of the conditional utility of a skill (see Equation 1). We use the fact that the utility of each skill is connected to the average wage of each latent state of the proposed HMM, and we naturally propose to evaluate on the average wage of each predicted latent state. In particular, for each one of the available skills, we compute the average wage of each predicted state. We then compute the mean absolute error between the average wage of the predicted state, and the actual wage of the instance at hand. Formally, the mean absolute error of the HMM is given by the following equation:

$$MAE_{HMM} = \sum_{x \in \{\text{beginner, knowledgeable, expert}\}} \frac{1}{|D_{e_x}|} \sum_{d \in D_{e_x}} (\bar{w}_{D_{e_x}} - w_d),$$

where $\bar{w}_{D_{e_x}}$ is the average wage of the set of instances D_{e_x} , for which the HMM predicted that the worker belongs in state e_x .

Baselines: We propose to compare our approach to two different baselines. The first one computes the average wage for each skill. We name this baseline '*Simple Average*' (SA). SA implicitly assumes that the utility of a skill is independent of a worker's level of expertise. The mean absolute error for the simple average is given by the following equation:

$$MAE_{SA} = \frac{1}{|D|} \sum_{d \in D} (\bar{w} - w_d)$$

The second baseline goes one step further. It assumes that a worker is an 'expert' in a certain skill if that skill is listed on the worker's profile. We name this baseline '*Mentioned Average*' (MA), and we compute the mean absolute error as follows:

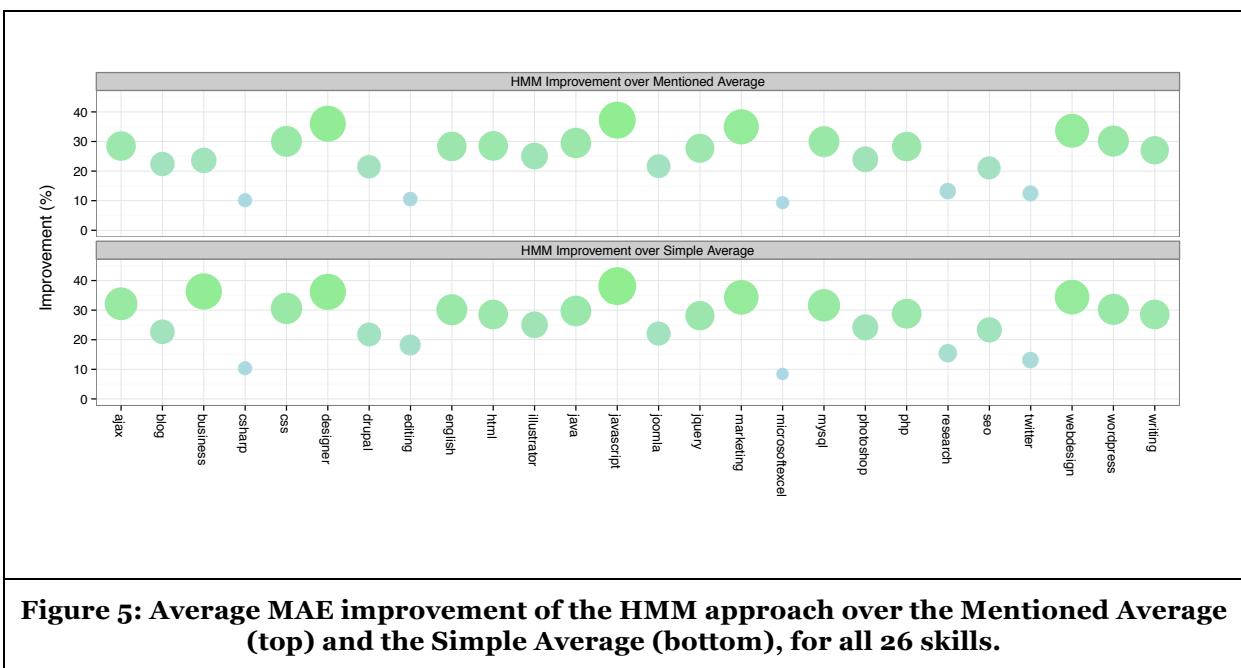
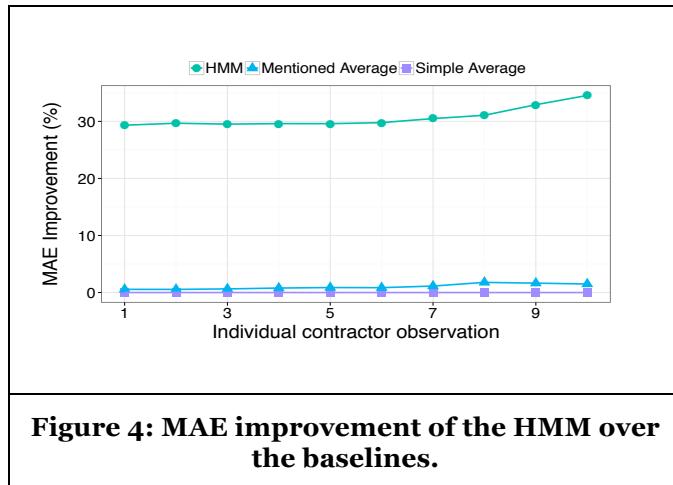
$$MAE_{MA} = \frac{1}{|D_{m_0}|} \sum_{d \in D_{m_0}} (\bar{w}_{D_{m_0}} - w_d) + \frac{1}{|D_{m_1}|} \sum_{d \in D_{m_1}} (\bar{w}_{D_{m_1}} - w_d),$$

where $\bar{w}_{D_{m_0}}$ ($\bar{w}_{D_{m_1}}$) is the average wage of the set of instances D_{m_0} (D_{m_1}), for which the given skill is not mentioned in the worker's profile. As we mentioned earlier, the intuition here is that if a worker feels comfortable enough to list a skill on his/hers profile, he/she must have some level of expertise in that skill. Finally, note that $|D_{m_0}| + |D_{m_1}| = |D_{e_1}| + |D_{e_2}| = |D|$.

To better illustrate the improvement of the proposed HMM over the two baselines, we further define the percentage mean absolute error improvement over the simple average baseline, as follows:

$$MAE \text{ Improvement}_{model} = \frac{MAE_{SA} - MAE_{model}}{MAE_{SA}}, \quad model \in \{HMM, \text{Mentioned Average}\}$$

In Figure 4, we present the average results across all skills. On the x-axis, we show the number of completed tasks needed for each individual worker in our testing dataset to make a wage prediction. On the y-axis, we show the MAE improvement. The Simple Average baseline is at zero. Any positive value on the y-axis shows an improvement over the Simple Average. We can see that the proposed HMM performs up to 35% better than the Simple Average baseline and the Mentioned Average. Additionally, the Mentioned Average baseline performs better (up to 3%) than the Simple Average.



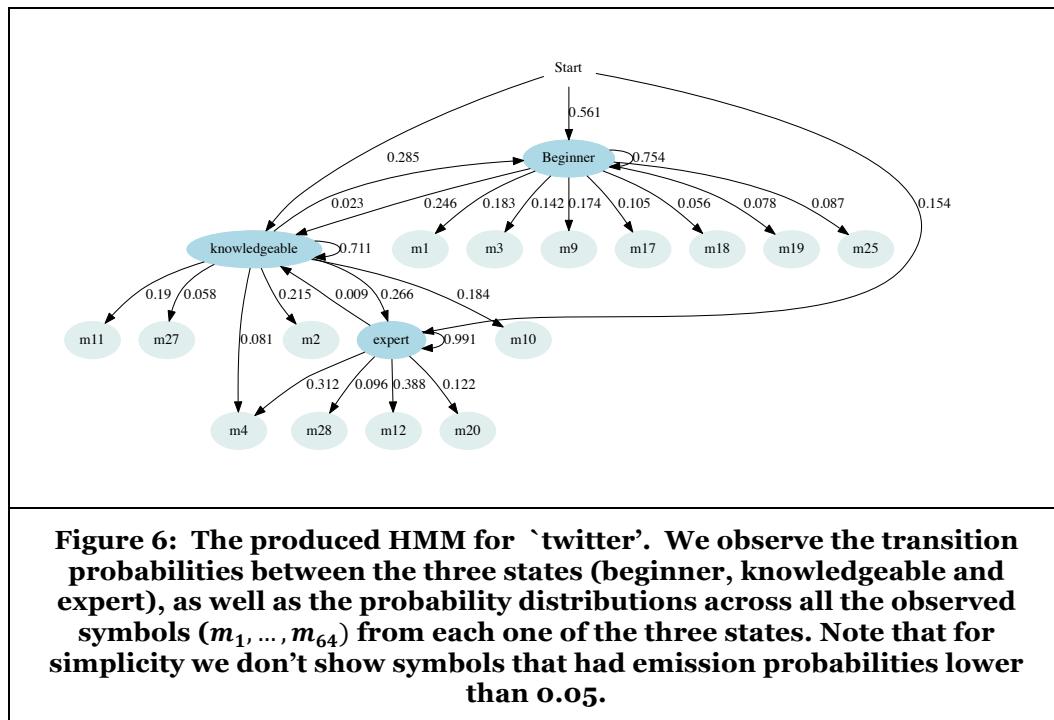
In Figure 5, we show the average mean absolute error improvement of our proposed HMM solution over the Mentioned Average baseline on the top, and over the Simple Average baseline on the bottom, for each one of the 26 skills that we study. One thing to notice is that our approach has a positive improvement across all skills. The highest improvement is observed for ‘javascript’ (40%), while most of the skills appear to have an improvement of around 25%.

In Figure 6, we show the actual computed HMM for one of our skills (twitter⁵). The initial probabilities of being an expert (0.154), knowledgeable (0.285), and a beginner (0.561) represent the distribution of expertise for twitter: 15.4% of workers that complete ‘twitter’ tasks appear to be experts, 28.5% appear to be knowledgeable, and the rest of them beginners. Furthermore, if a person is a beginner, there is a 24.6% probability of transitioning to knowledgeable, and 75.4% of remaining a beginner. Similarly, if one is knowledgeable, there is a very small probability (2.3%) of transitioning to a beginner, and a decent probability (26.6%) of becoming an expert, and so forth. Each expertise level has different distributions

⁵ For additional HMMs for other skills: <http://people.stern.nyu.edu/mk3539/other/hmms.pdf>

across all possible emitted symbols⁶. For example, a beginner is highly likely (24.6%) to emit symbol m1 (all variables in the bottom 50%), while an expert very frequently (38.8%) emits symbol m12 (rehires with feedback scores in the top 50%, and all other features in the bottom 50%). Furthermore, there are certain symbols that are emitted from multiple states with different probabilities; for instance, m4 (rehires=1, everything else in the bottom 50%) is emitted with significant probabilities from both the knowledgeable and the expert states; however, it is almost four times more likely to be emitted from an expert than from a beginner. Finally, note that the graph in Figure 6 is incomplete, since for simplicity we only show emitted symbols with probabilities greater than 0.05.

In general, we observe that an expert is more likely to emit states that include rehires and feedback in the top 50%; a knowledgeable is more likely to emit states in which wages appear in the top 50% and feedback in the top 50%; and a beginner is more likely to emit states in which the hiring rate is in the top 50% and where the workers list ‘twitter’ as a skill on their profile. Also note that there are no states with certification in the top 50%. The reason for that is that in our training set, very few (0.1%) of the workers have taken a certification test on Twitter and decided to list it on their profile.

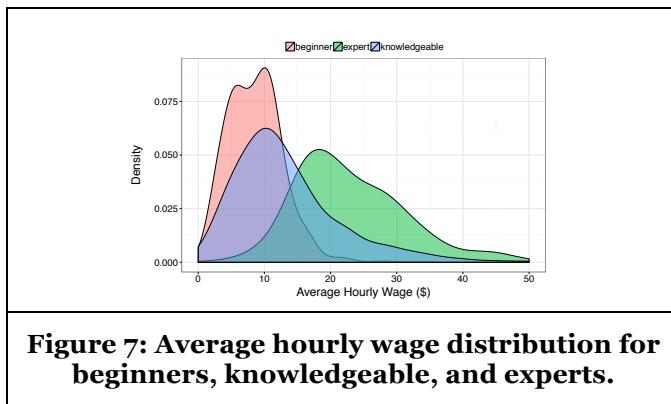


The utility of skills on oDesk.com

After evaluating our proposed approach for estimating the level of expertise for each worker in a given skill, we are ready to move forward and estimate the conditional utilities given by Equation 1 for each one of the skills that we consider in this study.

In Figure 7, we present the average hourly wage distribution across all 26 skills. We observe that the experts’ distribution is significantly shifted towards higher wages, suggesting that workers that are predicted to be experts by our HMM receive significantly higher compensation.

⁶ To interpret the emitted symbols, see Equation (2).



Skill	50 th percentile (\$)	90 th percentile (\$)	Variance
`research'	9.0	15.0	20.14
`java'	14.44	22.22	44.5
`microsoftexcel'	12.34	16.67	27.45
`marketing'	13.0	38.89	136.32
`editing'	16.67	35.0	119.9

Table 2: Wage Percentiles

In Figure 8 we show the utility of each one of the 26 skills, estimated by the use of Equation (1). On the y-axis we have the experts' hourly wage, and on the X-axis, we have the beginners' hourly wage. The utility of each skill is captured by the diameter of each point (the actual utility values are given in the size/color legend). The line is the 'equality line' (45 degrees), indicating that any point above that line has a higher expert wage than beginner wage. First, we note that all skills in our set are above the equality line, verifying intuition. Simply put, this shows that a worker has more value in the marketplace if the worker is an expert on a skill, than a beginner. Second, all skills have positive utility, and greater than \$2.5. Certain skills appear to provide much higher utilities than others: for instance, `marketing' and `editing' have the highest utilities (around \$12.5) while `microsoftexcel' and `java' appear to have the lowest ones (around \$3). Intuitively, this means that being a really good editor, pays a lot more than being let's say a "microsoftexcel" expert. To understand why this is happening, we will use the information provided in Table 2, where we show the 50th and 90th wage percentile, as well as the variance of the wage distribution for five skills. The first three rows, represent low utility skills (`research', `java', `microsoftexcel'), and the last two high utility skills (`marketing', `editing'). We see that the wage differences between the 90th and the 50th percentiles for the low-utility skills are very small, compared to the wage differences of the high-utility skills. This higher spread in distributions is also verified by the higher variance of the high-utility skills. Since our HMM explicitly associates high expertise with high wage, it captures the variance of the wage distributions of each skill by identifying those high-paid workers as experts. This behavior results in higher utilities for skills for which the variance is higher.

Implications, limitations, and future directions

Our study is the first that explicitly quantifies the value of skills in a marketplace. Even though our study has limitations and is constrained by a small set of skills, it clearly communicates a technically sound methodology for analyzing the materialized skills' value.

Implications: There is a direct impact of our work on the online labor marketplace: the platform can strategically suggest to workers to built up their expertise in certain skills. For example, drawing on our earlier discussion regarding `microsoftexcel' and `editing', the marketplace can increase its revenue by acquiring more experts in `editing'. Of course, such a recommendation might change the estimated utility of a skill (e.g., `editing' in our example), and us a result, constant monitoring and system updating is

required for better results. Furthermore, knowing the actual value of being an expert in a given skill is the first step towards building a framework for recommending skills (discussed later in ‘Future directions’). With the deployment of such a recommendation scheme, the marketplace can strategically recommend new skills to workers. If workers start acquiring skills with high utility, both their demand and their income will increase. Since more workers will have skills that have high demand in the marketplace, both the job openings’ closing rate and (as a result) the marketplace’s revenue will increase. Finally, even though the results of this study are constrained within a given online labor market, and a given set of skills, the methodology is widely applicable. For example, TaskRabbit⁷, could potentially use our approach to identify the value of each one of the skills they consider in their marketplace. Furthermore, networks such as LinkedIn⁸ could also use our approach, by simply associating the value of each skill/level of expertise (LinkedIn endorsements) with promotions or new jobs.

Limitations: Besides the assumption regarding the cost of effort that we discussed earlier, in this work we assume that skills are independent from each other. This assumption is critical in defining the utility of one skill, and not the utility of a set of skills. In the future, we intend to relax this assumption and study the associations between skills, in order to provide utility estimations for any given set of skills. Finally, we assume that employers know exactly what they are looking for, and adequately describe the set of required skills for each task that they post.

Future directions: In the future, we intend to include the development of a skills recommendation framework. In particular, consider a scenario where at any given time, workers have two options: (1) to exploit their current skillset and expertise by getting hired and completing a task; or (2) invest their time in improving/expanding their skillset and expect future increased returns. What is the optimal decision for each worker? Exploit or improve? We plan to develop a framework that captures this behavior and recommends the optimal decision for the worker. Such a system could be used as a career development adviser.

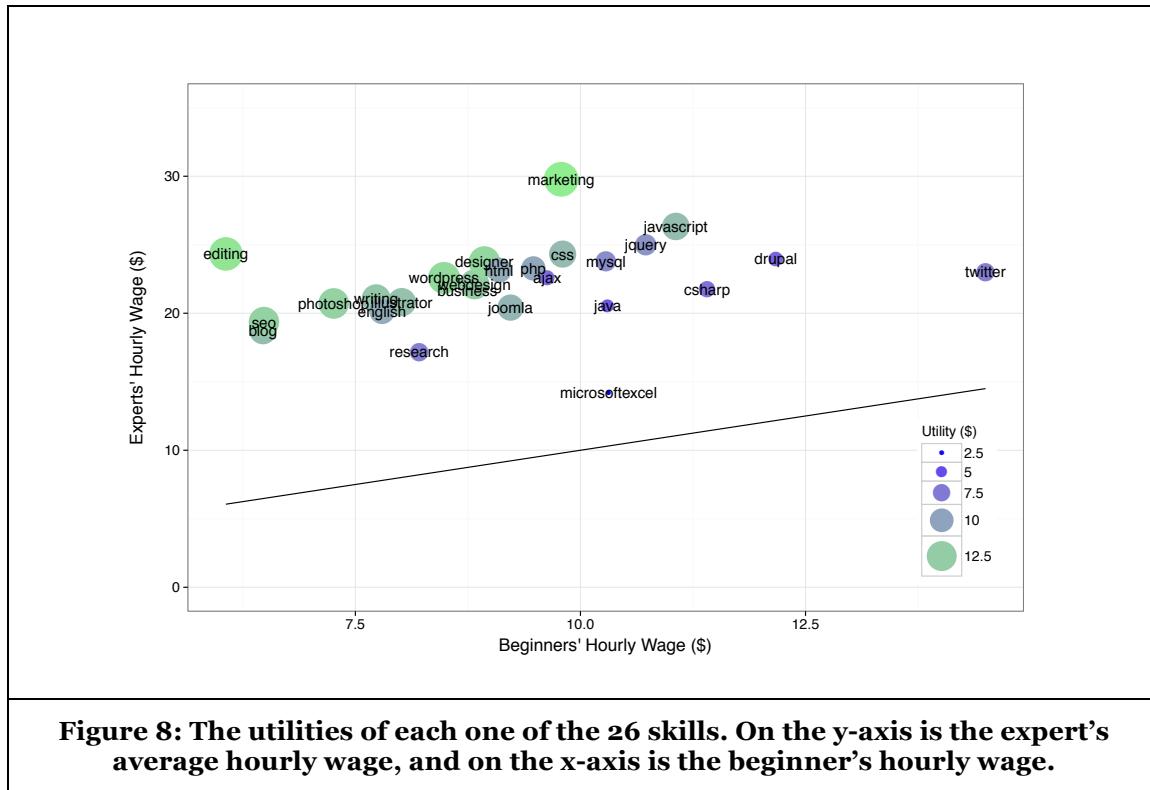


Figure 8: The utilities of each one of the 26 skills. On the y-axis is the expert’s average hourly wage, and on the x-axis is the beginner’s hourly wage.

⁷ <http://www.taskrabbit.com/>

⁸ <http://www.linkedin.com/>

References

- Agrawal, A., Horton, J., Læctera, N., and Lyons, E. 2013. "Digitization and the Contract Labor Market: A Research Agenda," in *Economics of Digitization*, University of Chicago Press.
- Balog, K., and De Rijke, M. Year. "Determining Expert Profiles (With an Application to Expert Finding)," IJCAI2007, pp. 2657-2662.
- Berinsky, A. J., Huber, G. A., and Lenz, G. S. 2012. "Evaluating online labor markets for experimental research: Amazon.com's Mechanical Turk," *Political Analysis* (20:3), pp. 351-368.
- Bouguessa, M., Dumoulin, B., and Wang, S. Year. "Identifying authoritative actors in question answering forums: the case of yahoo! answers," Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM2008, pp. 866-874.
- Brin, S., and Page, L. 1998. "The anatomy of a large scale hypertextual Web search engine," *Computer networks and ISDN systems* (30:1), pp. 107-117.
- Chandler, D., and Horton, J. J. 2011. "Labor Allocation in Paid Crowdsourcing: Experimental Evidence on Positioning, Nudges and Prices," *Human Computation* (11), p. 11.
- Cunha, F., Heckman, J. J., and Schenckach, S. M. 2010. "Estimating the technology of cognitive and noncognitive skill formation," *Econometrica* (78:3), pp. 883-931.
- Danescu Niculescu-Mizil, C., Kossinets, G., Kleinberg, J., and Lee, L. Year. "How opinions are received by online communities: a case study on amazon.com helpfulness votes," Proceedings of the 18th international conference on World wide web, ACM2009, pp. 141-150.
- Dellarocas, C. 2003. "The digitization of word of mouth: Promise and challenges of online feedback mechanisms," *Management science* (49:10), pp. 1407-1424.
- Desmarais, M. C., Maluf, A., and Liu, J. 1995. "User expertise modeling with empirically derived probabilistic implication networks," *User modeling and user adapted interaction* (5:3-4), pp. 283-315.
- Goes, P., and Lin, M. 2012. "Does Information Really "Unravel"? Understanding Factors That Motivate Sellers to Seek Third Party Certifications in an Online Labor Market," *Understanding Factors That Motivate Sellers to Seek Third Party Certifications in an Online Labor Market* (November 1, 2012). NET Institute Working Paper:12-02.
- Hambleton, R. K., Swaminathan, H., and Rogers, H. J. 1991. *Fundamentals of item response theory*, (Sage).
- Handel, M. J. 2003. "Skills mismatch in the labor market," *Annual Review of Sociology*, pp. 135-165.
- Horton, J. J. 2010. *Online labor markets*, (Springer).
- Horton, J. J., Rand, D. G., and Zeckhauser, R. J. 2011. "The online laboratory: Conducting experiments in a real labor market," *Experimental Economics* (14:3), pp. 399-425.
- Hu, N., Zhang, J., and Pavlou, P. A. 2009. "Overcoming the J-shaped distribution of product reviews," *Communications of the ACM* (52:10), pp. 144-147.
- Ipeirotis, P. G., and Horton, J. J. Year. "The need for standardization in crowdsourcing," Proceedings of the CHI2011.
- Ipeirotis, P. G., Provost, F., and Wang, J. Year. "Quality management on amazon mechanical turk," Proceedings of the ACM SIGKDD workshop on human computation, ACM2010, pp. 64-67.
- Jerath, K., Fader, P. S., and Hardie, B. G. 2011. "New Perspectives on Customer "Death" Using a Generalization of the Pareto/NBD Model," *Marketing Science* (30:5), pp. 866-880.
- Jureczyk, P., and Agichtein, E. Year. "Discovering authorities in question answer communities by using link analysis," Proceedings of the sixteenth ACM conference on Conference on information and knowledge management, ACM2007, pp. 919-922.
- Kleinberg, J. M. 1999. "Authoritative sources in a hyperlinked environment," *Journal of the ACM (JACM)* (46:5), pp. 604-632.
- Kokkodis, M., and Ipeirotis, P. G. Year. "Have you done anything like that?: predicting performance using inter-category reputation," Proceedings of the sixth ACM international conference on Web search and data mining, ACM2013, pp. 435-444.
- Lazear, E. P. 2003. "Firm specific human capital: A skill weights approach," National Bureau of Economic Research.
- Liu, Y., Huang, X., An, A., and Yu, X. Year. "Modeling and predicting the helpfulness of online reviews," Data Mining, 2008. ICDM'08. Eighth IEEE International Conference on, IEEE2008, pp. 443-452.
- Lu, Y., Tsaparas, P., Ntoutas, A., and Polanyi, L. Year. "Exploiting social context for review quality prediction," Proceedings of the 19th international conference on World wide web, ACM2010, pp. 691-700.

- Macdonald, C., and Ounis, I. 2008. "Voting techniques for expert search," *Knowledge and information systems* (16:3), pp 259–280.
- Mason, W., and Watts, D. J. 2010. "Financial incentives and the performance of crowds," *ACM SigKDD Explorations Newsletter* (11:2), pp 100–108.
- Murphy, K. P. 2012. *Machine learning: a probabilistic perspective*, (MIT Press).
- Nelson, P. 1970. "Information and consumer behavior," *The Journal of Political Economy*, pp 311–329.
- Pallais, A. 2013. "Inefficient hiring in entry level labor markets," National Bureau of Economic Research.
- Petkova, D., and Croft, W. B. 2008. "Hierarchical language models for expert finding in enterprise corpora," *International Journal on Artificial Intelligence Tools* (17:01), pp 5–18.
- Rand, D. G. 2012. "The promise of Mechanical Turk: How online labor markets can help theorists run behavioral experiments," *Journal of theoretical biology* (299), pp 172–179.
- Saito, S., Watanabe, T., Kobayashi, M., and Takagi, H. 2014. "Skill Development Framework for Micro Tasking," in *Universal Access in Human Computer Interaction. Universal Access to Information and Knowledge*, Springer, pp. 400–409.
- Shaw, A. D., Horton, J. J., and Chen, D. L. Year. "Designing incentives for inexpert human raters," Proceedings of the ACM 2011 conference on Computer supported cooperative work, ACM2011, pp. 275–284.
- Sheng, V. S., Provost, F., and Ipeirotis, P. G. Year. "Get another label? improving data quality and data mining using multiple, noisy labelers," Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM2008, pp. 614–622.
- Smir, E. M., and Hitt, L. M. 2003. "Costly bidding in online markets for IT services," *Management Science* (49:11), pp 1504–1520.
- Spence, M. 1973. "Job market signaling," *The quarterly journal of Economics*, pp 355–374.
- Varshney, K. R., Wang, J., Mojsilovic, A., Fang, D., and Bauer, J. H. Year. "Predicting and recommending skills in the social enterprise," Proc. Int. AAAI Conf. Weblogs Soc. Med. Workshops2013.
- Zhang, J., Ackerman, M. S., and Adamic, L. Year. "Expertise networks in online communities: structure and algorithms," Proceedings of the 16th international conference on World Wide Web, ACM2007, pp. 221–230.