**Using Machine Learning Techniques to Uncover What Makes Understanding Spoken Chinese Difficult for Non-native Speakers**

[EXTENDED ABSTRACT REVISION]

John Kowalski
Advisor: Geoff Gordon

**Abstract**

The Pinyin Tutor has been used for the past few years in over thirty classrooms at universities around the world. A large amount of data have been collected from this program on the types of errors students make when trying to spell the pinyin of the Chinese phrase spoken to them. I plan to use this data to help answer the question of what is hard about understanding Chinese. Is it a particular set of consonants, vowels, or tones? Or perhaps do certain difficulties arise in the context in which these sounds are spoken? Since each pinyin phrase can be broken down into features (consonants, vowel sounds, and tones), we can apply machine learning techniques to uncover the most confounding aspects for beginning students of Chinese. We can extend the methods we developed here to create an ML engine that learns on the fly for each student what they find difficult. The items to be presented to the learner can be chosen from a pool based on the predicted probability of being correctly answered by the student. This will allow the Chinese learner to focus on what he or she is having most difficulty and hopefully more quickly understand spoken Chinese than without such focused "intelligent" instruction.

**1. Introduction**

One of the hurdles for students in introductory Chinese courses is understanding what Chinese word or phrase has been spoken to them. On top of learning new vocabulary, translation, and grammar skills, simply answering, "What did you hear?" is often difficult for beginning students. Improving this rudimentary skill for each student in the classroom can be tedious and in general a poor use of class time. Luckily, computer tutors are perfectly suited for such instruction. At the Pittsburgh Science of Learning Center (PSLC), the Pinyin Tutor was developed for this purpose *(figure 1)*. Pinyin is a system of writing Standard Mandarin using Roman characters and is commonly taught to first-semester students of Chinese in the first year of instruction. The Pinyin Tutor's instructional model is quite simple: a Chinese word or phrase is "spoken" through the students' personal computer speaker and the task is to enter the pinyin of the sound they heard. If correct, the tutor congratulates the student and presents the next item in the lesson. If incorrect, the tutor gives feedback on what part of the phrase is incorrect and gives the student an opportunity to try again. The items the student answered incorrectly on the first try are put back into the pool to be presented again; items answered correctly on the first try are eliminated from the pool. The student assignment is to continue until all items in the lesson have been answered correctly on the first try (without feedback from the tutor).

**2. Pinyin Basics**

Each syllable in Pinyin can be thought of as being comprised of three components: an "initial" consonant sound (b, n, zh, …), a "final" vowel sound (ai, ing, uang, …), and a tone marking (1,2,3,4,5) indicating the rising/falling pitch of the syllable. In total there are twenty-three initials, thirty-six finals, and five tones *(figure 2)*.

**3. Identifying What Causes Students the Most Difficulty**

The current version of the Pinyin Tutor keeps track of what Chinese words or phrases are

answered incorrectly by students and re-drills the word or phrase until the student can identify it without help from the tutor.  Beyond knowing what items students having trouble with, the obvious next question to ask is what component(s) of these items are causing the most difficulty.  Are there some initials, finals, or tones particularly problematic?   And then beyond asking this question for individual components, we must consider difficulties that may arise in the context in which these sounds are spoken.  For instance, is final "ao" easy to hear in "hao3", but not in "bao4"?

To answer these questions, we can train a statistical model of student performance.  If we can obtain a reasonably low generalization error with this model, we can use the parameters learned to discover what the easy and difficult components are.

As a first try, we can model a Chinese syllable by constructing a linear prediction with 64 covariates (one for each 23 initials, 36 finals, 5 tones), and output whether this item was answered correctly.  So we have:

$$\hat{y} = \beta_0 + \beta_1 * x_1 + \beta_2 * x_2 + \ldots + \beta_{64} * x_{64}$$

Where:

$y$=1 (correct), $y$=0 (incorrect).
$x_n$=1 if feature n is present, 0 otherwise.
$\beta_n$=coefficient for feature n to be learned.

We can supplement this single-syllable model with interaction terms, adding a term for each interaction of initial-final, initial-tone, and initial-tone.  This will enable us to answer for instance, whether some features are typically easy by themselves, but when combined are difficult.  We can also add terms to indicate the number of times a student was exposed to these features and contexts, enabling us to potentially discover some feature to be initially difficult, but fast to learn once presented.

Another way to train our models is by using L1-regularized logistic regression.  Instead of fitting our model directly to the training data, in logistic regression we fit our data to a logistic curve.

$$\hat{y} = f(\beta_0 + \beta_1 * x_1 + \beta_2 * x_2 + \ldots + \beta_{64} * x_{64})$$

Where:

$y$=1 (correct), $y$=0 (incorrect).
$x_n$=1 if feature n is present, 0 otherwise.
$\beta_n$=coefficient for feature n to be learned.
$f()$ = sigmoid function

Logistic analysis is a natural choice for fitting a model to our categorical data (correct/incorrect).  We use L1-regularized logistic regression for its feature selection properties (Koh, 2007).  This will give us the flexibility similar to the "s" parameter in lasso. We use the method developed by Koh, et al. for its ability to scale well to large sparse problems (like for our 2374 covariate model for all interaction terms).

[Perhaps flesh-out L1-regularized logistic regression info here a bit more…]

*3.1.Training the Models*

Having constructed representations of Chinese syllables with various levels of detail, we now want to train these models and arrive at the coefficients to best predict student correctness.  The method we will first use is the Least-Angle Regression (LARS), modified to simulate the

Lasso (least absolute shrinkage and selection operator). Lars and Lasso are regression algorithms especially suited for high-dimensional data (Efron, et al, 2004).

Using data we collected from the Pinyin Tutor in Fall 2008, we have approximately 250,000 training examples (tuples of Chinese syllable and binary indicator of whether student correctly identified it). The Lasso takes this input and fits our linear model described above using the criterion:

$$\text{Minimize}(\Sigma(y - \hat{y})^2), \text{ under constraint that:}$$
$$\Sigma_j (\text{absoluteValue}(\beta_j)) <= s$$

Where:
- y is the actual indicator of correctness from the training data (0=incorrect, 1=correct)
- $\hat{y}$ is the running predicted output of our model computed by the lasso. We assume if $\hat{y} >$ .5, then predict correct, $\hat{y} <= .5$, then predict incorrect
- s is a bound that can be used as a tuning parameter.


What makes the lasso special is the "s" parameter. If it is sufficiently large, the lasso is basically the ordinary multiple linear least squares regression of y on the $x_n$ covariates. But for smaller positive values, lasso computes a "shrunken" version of the usual least squares estimate. This shrunken version will often lead to some of the coefficients ($\beta$'s) being zero, so choosing a value for s is like choosing the number of predictors in the model (Tibshirani, 1996). For a graph of LARS-Lasso coefficient learning progress for our 64-covariate pinyin model see *(figure 3)*. To see how well the model learned by lasso predicts, see the graph of the K-fold cross-validated mean squared prediction error as the model is learned by LARS-Lasso *(figure 4)*.


**4. Making the Pinyin Tutor More Intelligent**

While the practice the Pinyin Tutor gives students is valuable (Zhang, 2008), it may not be the most efficient way for students to learn since remediation is based on the whole item, not on the parts of the item they have shown to have difficulty. For example, if the tutor presents the two-syllable item "ni3hao3", but the student incorrectly types "ni3ho4", the tutor will put this item back into the pool for a future re-drilling. But the student did not demonstrate difficulty with initials "n" or "h", final "i", or the tone in the first syllable. The student did, however, demonstrate difficulty with the final "ao". Specifically, the student demonstrated difficulty with the final "ao", when in the second syllable, preceded by initial "h", with tone 3. What if we could give students more practice on the parts they have demonstrated difficulty? For our example, what if we could give more practice on items most similar to having a final "ao" in the second syllable, preceded by initial "h", with tone 3? We can in fact augment the Pinyin Tutor so that we can have such "intelligent" behavior by training a hidden Markov model for each possible feature and calculating the probability for each item in the lesson of being in the "Unlearned" state. This technique of estimating the probability a student knows a skill after observing their attempts is known as "knowledge tracing" in the intelligent tutor community (Corbett & Anderson, 1995).

*4.1. Hidden Markov Model Basics*

A hidden Markov model (HMM) is a statistical model that represents a Markov process with "hidden" state. A Markov process model can be thought of as a finite state machine where transitions between states are assigned probabilities and the probability of transition to state B at time t+1 depends only on the state the machine is in at time t. The "hidden" part is that we are not able to directly observe what state the machine is in currently, just the output dependent on the current state. The output symbols per state are assigned a probability distribution.

*4.2. Hidden Markov Model Application*

We can create an HMM for each skill necessary to master correctly spelling the pinyin of a spoken Chinese phrase. The skill set can be the 64 basic features. We can also supplement the basic 64 skill models with models for interactions between each, so a skill model for each interaction between initial-final, initial-tone, and final-tone (totaling $64 + 23*36 + 23*5 + 36*5 = 1187$ HMMs (skill models)).

The HMM for each skill has two states: "Learned" (L) and "Unlearned" (U). In each of these hidden states, we are able to observe whether a skill was answered correctly (C) or incorrectly (I). The transition probabilities between each state are:
- P(learn), the prior probability of transitioning (U)->(L) at each step given that we start at (U)
- P(forget), the prior probability of transitioning (L)->(U) at each step given that we start at (L)

And conditional observational probabilities within each "hidden" state are:
- P(slip), the probability making an error even if the student is in the (L) state
- P(guess), the probability of guessing correctly even if the student is in the (U) state
See *(figure 5)*.

Now that we have the structure (states and transition paths) of the HMMs necessary to represent the skills to tutor, we need to calculate the transition probabilities between the two states. While there is no known way to optimally "train" an HMM (it's an NP-complete problem), fortunately there is a way to approximate transition probabilities via the Baum-Welch algorithm, which makes use of the forward-backward algorithm used frequently in HMM computations (Rabiner, 1989).

Using the data collected from the Pinyin Tutor in the Fall 2008 semester and considering only the student first attempts at each item (and thus not influenced by tutor feedback), we have approximately 250,000 student-tutor interactions we can train our pinyin skill HMMs utilizing the Baum-Welch algorithm.

[*Perhaps include **condensed** details here of how the Baum-Welch algo works. It is explained well in (Rabiner, 1989).*]

*4.3. Using Trained HMMs to Estimate Student Learning State*

Now that we have HMMs for each skill, we can calculate the probability a skill is in the "Learned" or "Unlearned" state after observing correct/incorrect observations for that skill as the student works through the tutor lesson. Rabiner refers to this as "Problem 2" of HMM applications. Namely, given an observation sequence, O, (of correct/incorrect responses) and an HMM, $\lambda$, for the skill, what is the probability we are in a state $S_i$ ((U) or (L)) at time t?

$\gamma_t(i) = P(q_t = S_i \mid O_1, O_2, ..., O_t, \lambda)$

[*Perhaps I include **condensed** details here of how the computation is done. It is explained well in (Rabiner, 1989).*]

Using the probabilities of being in the "Unlearned" state for each skill, we can now average the probabilities of each feature present in an item being in the "Unlearned" state and base our selection of the next item to present on this calculation. Our aim is that by choosing items with highest average probability of being in the "Unlearned" state, the tutor will tailor its instruction for the student precisely on skills he or she is having the most difficulty. Our ultimate goal is to

get the students to a certain level of mastery much faster than the previous tutor model.

## 5. Conclusions and Future Work

With the tutor giving focused and tailored instruction, we expect students will learn more quickly than with the previous non-intelligent model. Preliminary plans are being made for an in-lab study at the Hong Kong International School to test this theory. We expect the full version of the tutor as described in this paper to be online by summer 2010, and used in schools for the Fall 2010 semester. It is currently being advertised on the Pinyin Tutor website and a beta version is currently being tested by at least two sites.

Running analyses on the entire data set has been a recurring challenge, as has been running analyses with large numbers of covariates due to memory constraints and analysis programs crashing. We believe completing these extra-large analyses is achievable (perhaps before May).

Another potential area for future work is to explore if there are any results interesting to the machine learning community. (Perhaps there are methods or insights comparing results of our HMM training with the Lars-lasso or L1-regularized logistic results, other ideas...)

Figure 1.

Figure 2.
**Ordering of Covariates for "Basic" Model (without interaction terms)**
Initials:

```
1       b
2       p
3       m
4       f
5       d
6       t
7       n
8       l
9       g
10      k
11      h
12      j
13      q
14      x
15      z
16      c
17      s
18      r
```

Finals:

```
19      zh
20      ch
21      sh
22      w
23      y
24      a
25      e
26      i
27      o
28      u
29      v
30      ai
31      ao
32      an
33      ei
34      en
35      ia
36      ie
37      iu
38      in
39      ou
40      ua
41      uo
42      ui
43      un
44      ve
45      vn
46      ue
47      ang
48      eng
49      ian
50      ing
51      iao
52      ong
53      uai
54      uan
55      van
56      iang
57      iong
58      uang
59      ueng
```

Tones:

```
60      Tone 1
61      Tone 2
62      Tone 3
63      Tone 4
64      Tone 5
```

Figure 3.
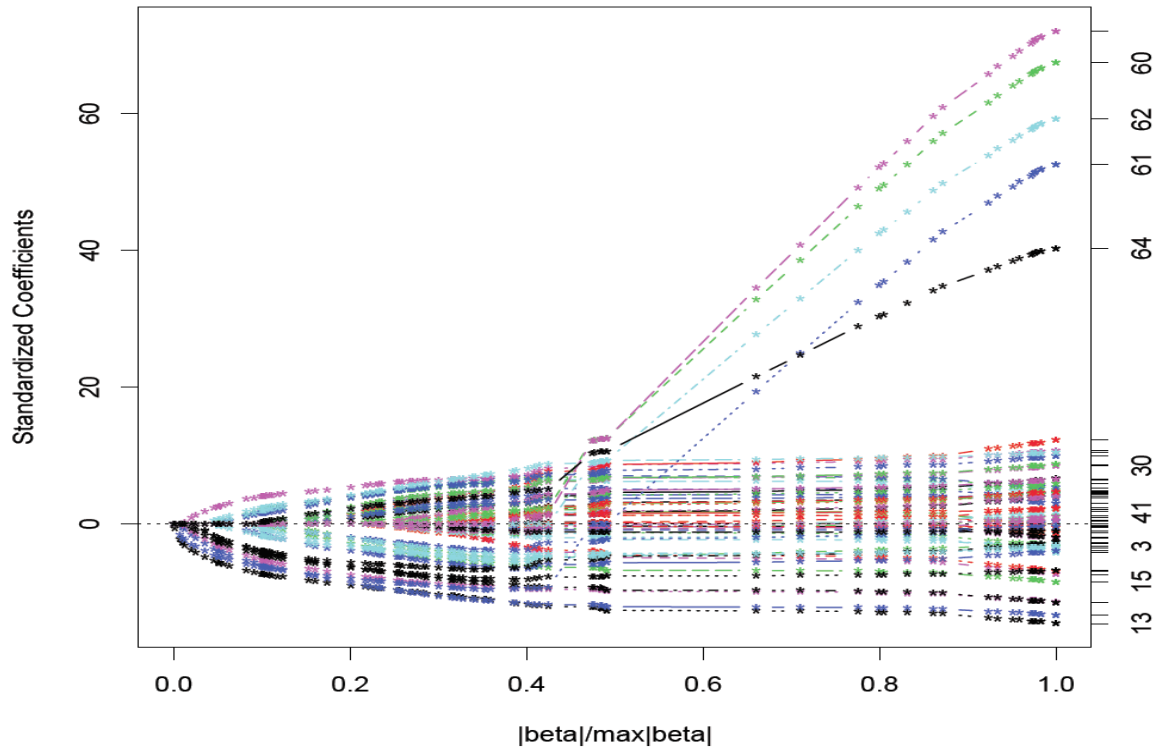Coefficient learning progress of LARS-Lasso on 64 covariate model.



Figure 4.
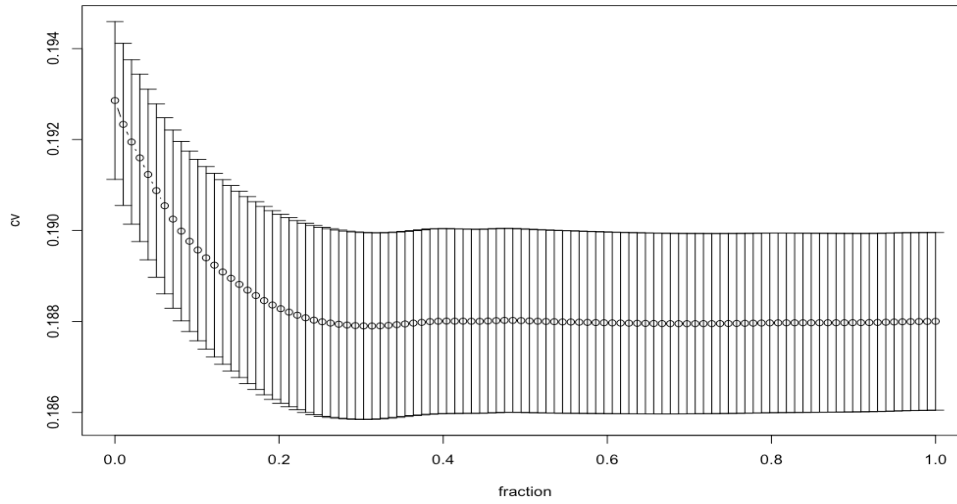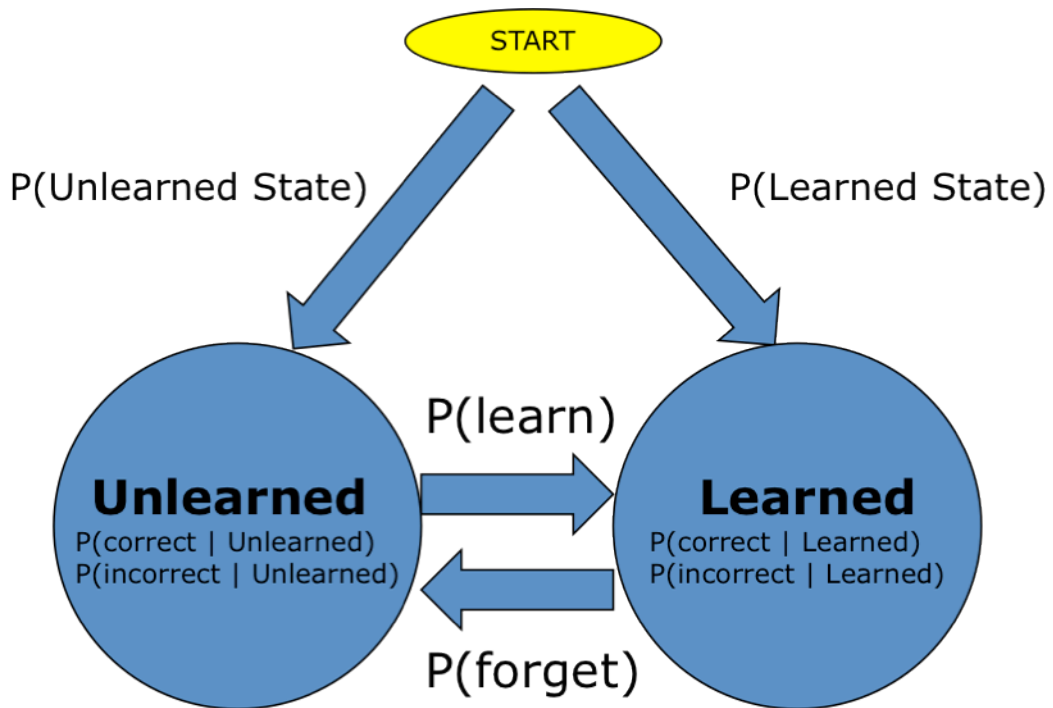K-fold cross-validated mean squared prediction error as modeled by LARS-lasso

Figure 5.
Hidden Markov Model representation (one for each pinyin skill).

**References**

[1] Efron, B., Hastie, T., Johnstone, I., Tibshirani, R. (2004), Least Angle Regression, Annals of
    Statistics, Volume 32, Number 2, 407-499.

[2] L. R. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech
Recognition," Proc. of the IEEE, Vol.77, No.2, pp.257-286, 1989.

[3] Tibshirani, R., (1996), Regression Shrinkage and Selection via the Lasso. Journal of the
Royal Statistical Society, Series B (Methodological), Volume 58, Issue 1, 267-288.

[4] Zhang, Y. (2008), Cue Focusing for Robust Phonological Perception in Chinese.

[5] Corbett, A. and J. Anderson, Knowledge tracing: Modeling the acquisition of procedural
knowledge., 1995.

[6] Casella, George, and Berger, Roger L. (2002). Statistical Inference, Second Edition. Duxbury
Press, Pacific Grove, California.

[7] Hastie, Tibshirani, and J. H. Friedman. Elements of Statistical Learning New York: Springer,
2009

[8] Koh, K., Kim, S., Boyd S. (2007) An Interior-Point Method for Large-Scale L1-Regularized
Logistic Regression, Journal of Machine Learning Research Number 8, 1519-1555.)

[9] Cen, H., Generalized Learning Factors Analysis: Improving Cognitive Models with Machine
Learning.