

**MSiA 421, Data Mining**  
**Assignment 3: CLV and Churn Modeling**

Due: Thursday, February 6 by 1pm

1. (Based on Gupta and Lehmann (2003)) In Aug 1994, Jason and Matthew Olim launched CDNow in the basement of their parents' house in Ambler, Pennsylvania. Within a year, revenues reached \$2 million. Like most Web-based startup companies, CDNow focused heavily on acquiring new customers. By 1998, after acquiring a rival company, CDNow had a customer base of more than 3 million customers. The company highlighted the number of new customers in its reports to financial analysts. Was their emphasis on acquisition misplaced? The answer to this question depends on whether CLV exceeds acquisition costs.

From 1998–2000 the average customer acquisition cost ranged from \$30–\$55, according to company reports. During the same time, the annual gross margin per customer ranged from \$10–\$20 and averaged about \$15. The retention rate ranged from 51–68%. Assume an annual discount rate of 12%.

- (a) Estimate CLV (not counting acquisition costs) for the following scenarios assuming payments at the end of the year (ordinary):
    - i. Pessimistic:  $m = \$10$ ,  $r = .51$
    - ii. Average:  $m = \$15$ ,  $r = .595$
    - iii. Optimistic:  $m = \$20$ ,  $r = .68$
  - (b) Compare your answers with the reported acquisition costs (\$30–\$50). Was the acquisition strategy profitable? Continue to assume that payments come at the end of the period.
2. At the end of 1999 CDNow reported a loss of over \$100 million and in Mar 2000 it publicly announced that it had only enough cash to sustain operations for 6 months. Soon after, the media giant Bertelsmann entered negotiations to acquire CDNow. Bertelsmann bought CDNow for \$117 million in Jul 2000. Did Bertelsmann overpay? This is usually a complex question depending on many different factors. CDNow, however, had few physical assets and its major asset was its customers. Assuming the total number of customers (in Jun 2000) was 3.29 million customers, an average annual margin of \$15 and an annual discount rate of 12%, find customer equity (assuming payments at the end of the year) under the following scenarios:
  - (a) Average:  $r = .595$
  - (b) Optimistic:  $r = .68$
  - (c) Really optimistic:  $r = .8$  (due to Bertelsmann having better management, market position and more cash)

3. A company that offers cable television service acquires customers who pay a monthly fee of \$60 until they cancel the service. Customers may cancel at any time. You may assume the event that a customer cancels during any specific time period is statistically independent of the event that the customer cancels during any other time period. Assume an immediate annuity, where the first payment occurs at time 0, i.e., the customer pays for the service at the beginning of the month and if the customer cancels before the end of the month the customer does not receive any refund. Suppose that the company retains 82% of its customers each month and that this retention rate is constant across all months and customers.

- (a) What is the expected and median time of attrition in months? That is, if  $T$  is a discrete random variable that represents the time in years that a customer cancels, what is the mean and median of  $T$ ?
- (b) What is the expected life-time revenue of a customer assuming a monthly discount rate of 1%?
- (c) *Customer equity* is the sum of CLV across customers. Suppose that the company recently acquired 1000 customers as described above. What is the customer equity of these 1000 customers? (total CLV across 1000 customers)
- (d) The remaining parts discuss the effects of *unobserved heterogeneity* — .

Suppose that the retention rate is not constant across customers, i.e., the retention rate is heterogeneous. Assume that customers are from two segments: half the customers acquired during a particular month have a monthly retention rate of 92% (the “loyalists”) while the other half have a monthly retention rate of 72% (non-loyalists). Note that the aggregate retention rate during the first month is 82% (averaged across all customers), as in the original problem. Compute CLV for a single loyalist and separately for a single non-loyalist.

- (e) Suppose the company acquires 500 loyalists and 500 non-loyalists during a given month. Find the customer equity of the 1000 customers. Compare this with the customer equity obtained in part (c).
4. The Comcast “Digital Starter” plan costs \$29.99 for the first six months and \$60.48 thereafter. Payments come at the beginning of the month, i.e., the first payment is made when the technician connects the household. There are no cancellation fees. Suppose that the retention rate is 95% during the four months ( $r_1 = \dots = r_4 = .95$ ),  $r_5 = 70\%$  during the month prior to the first higher payment,  $r_6 = 60\%$  during the next month, and  $r_t = 90\%$  thereafter ( $t \geq 7$ ). Let random variable  $T$  be the month of cancelation, i.e., customers are retained for  $T - 1$  months and cancel during month  $T$ . For example, if  $T = 1$  then the customer makes one payment at the time of installation and cancels before the second payment, receiving no refund for the unused portion of the first month. Assume Comcast does not spend any money on marketing to retain customers.

- (a) Find the probability that a customer defaults during month 8.
  - (b) Find the probability that a customer will cancel after the first year (i.e., has made 12 or more payments)?
  - (c) Find the probability that a customer is retained for the first two years. (i.e., makes at least 24 payments)
  - (d) Find the probability that a customer cancels before the end of the first year (i.e.,  $T < 12$ ).
  - (e) Find the expected month of attrition (i.e., the mean of  $T$ ).
  - (f) Graph the hazard, PDF and Survival functions.
  - (g) Find expected CLV, assuming an annual discount rate of 12%.
5. (39 points) The `cell.csv` data set has a sample of customers from a leading US cell phone provider.
- (a) (2 points) The variable `billmonth` is a string variable (\$6.) that tells the “current” month. It ranges from Jun 2007 through May 2008. The `churn` variable indicates whether a customer churned during the “current” month (1 = churn, 0 = not churn). Run a crosstab of these two variables. Describe the apparent sampling plan.
  - (b) (4 points) How many people canceled? How many opportunities were there to cancel? What is the retention rate? Use the `t2` variable for this part, which gives the difference in months between `billmonth` and `contractdt`. The `contractdt` variable gives the start date of the current contract, so `t2` will tell us the time until churn or censoring.
  - (c) (2 points) Briefly discuss the effect of the sampling plan on your estimates of the retention rate?
  - (d) (4 points) Compute average monthly revenue using the `TOTAL_REV_AMT` variables. Assume payments come at the beginning of the month. Compute expected lifetime revenue using your retention rate estimate, assuming a monthly discount rate of 1%.
  - (e) (3 points) Run a life table analysis and generate hazard and survival function plots for the `T2` variable. Submit only the hazard or survival plot and briefly describe what it tells you about churning.
  - (f) (3 points) Generate life tables stratifying on accounts of type “B” (business) and “I” (individual). Use `account_type` variable. Submit only a survival or hazard plot and indicate which group is more likely to churn?
  - (g) (3 points) Generate separate life tables (or KM) and survival curves for different numbers of lines (`line_count`), after capping `line_count` at 5 (set all values greater than 5 equal to 5). What is the story? How does the number of lines affect retention rates? Submit only the survival or hazard plot.

- (h) Which of the following factors affect the likelihood of churn: account type, line count, time out of contract? You will build a discrete-time survival model using logistic regression to answer this question. First, we will study the effect of time out of contract after controlling for some other variables.
- i. (5 points) Create a variable `timeout` giving the number of months that someone is out of his/her contract, e.g.,  $-1 = 1$  month until the end,  $0 =$  contract just ended,  $+1 =$  one month out of contract, etc. Make a new variable that equals `timeout`, and then cap it at 7 and  $-7$  (any values greater than 7 should be set equal to 7 and values less than  $-7$  should be set to  $-7$ ). Using logistic regression, predict the variable “event” ( $=1$  if customer cancels during the current month, 0 otherwise) on your capped `timeout` variable, treating it as categorical, `account_type` and `line_count` (capped at 5).
  - ii. (2 points) Interpret the effects for `account_type` and `line_count`. Do they tell you the same story as the life tables?
  - iii. (4 points) Describe how the probability of canceling depends on `timeout`. Create a graph showing how the log-odds of canceling depends on `timeout`.
  - iv. (2 points) Which variables are most important in explaining cancelations?
- (i) (5 points) Make a new copy of `timeout` and cap it at 4 and  $-1$ . Create a variable that indicates whether or not there was an overage charge during the previous month (`lagover`). Use logistic regression to predict `event` from the new capped version of `timeout`, `account_type`, `line_count`, and lagged overage. Submit the estimated coefficients. Which variables seem to be good predictors of churning?