

Problem Set 1
10-601 Fall 2012
Due: Friday Sept. 14, at 4 pm

TA: Brendan O'Connor (brenocon@cs.cmu.edu)

Due Date

This is due at **Friday Sept. 14, at 4 pm**. Hand in a hard copy to Sharon Cavlovich, GHC 8215. This document was last updated Tuesday 11th September, 2012, 8:34pm.

Changelog: (9/6) Clarified that graphs need to be printed out and turned in. (9/10) clarified notation on 2.d.1, 2.d.2. (9/11) Added log-scale suggestion for 2.e.2; clarified wording of 2.e.5.

1 Probability Review

Please show all steps in your solution.

1.a Equation of the Reverend

Prove

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

1.b Contingencies

A is a random variable that can take one of two values $\{\diamond, \heartsuit\}$. B is a random variable that can take one of two values $\{\triangle, \square\}$.

There are 117 (A_i, B_i) pairs, with the following “contingency table” of counts: each cell says how many cases there are of that pair type, e.g. 12 cases where $(A, B) = (\diamond, \triangle)$.

	$A = \diamond$	$A = \heartsuit$
$B = \triangle$	12	97
$B = \square$	3	5

Compute the quantities

1. $P(A = \diamond)$
2. $P(A = \diamond \text{ AND } B = \square)$ (this is a notational equivalent of $P(A = \diamond, B = \square)$.)
3. $P(A = \diamond \text{ OR } B = \square)$
4. $P(A = \diamond | B = \square)$
5. Use the law of total probability to rewrite $P(A)$ in terms of conditional probabilities $P(A|B = \triangle)$ and $P(A|B = \square)$. Compute $P(A = \diamond)$ from this equation. (If this is how you did 1b(1), then compute it with a different, more direct, approach.)

1.c Chain rule

Rewrite $P(X, Y, Z)$ as a *product* of several conditional probabilities, and one unconditioned probability involving a single variable. Your conditional probabilities can use only one random variable on the left side of the conditioning bar. For example, $P(A|C)$ and $P(A)$ would be ok, but $P(A, B|C)$ is not.

1.d Total probability and independence

Let X, Y, Z all be binary variables, taking values either 0 or 1.

Assume Y and Z are independent, and $P(Y = 1) = 0.9$ while $P(Z = 1) = 0.8$.

Further, $P(X = 1|Y = 1, Z = 1) = 0.6$, and $P(X = 1|Y = 1, Z = 0) = 0.1$, and $P(X = 1|Y = 0) = 0.2$.

1. Compute $P(X = 1)$. (Hint: use the law of total probability.)
2. Compute the expected value $E[Y]$.
3. Suppose that instead of Y attaining values 0 and 1, it takes one of two values 115 and 20, where $P(Y = 115) = 0.9$. Compute the expected value $E[Y]$.

2 Decision Trees



Untergang der Titanic by Willy Stöwer, 1912

Below is a dataset of the 2201 passengers and crew aboard the RMS Titanic, which disastrously sunk on April 15th, 1912. For every combination of three variables (Class, Gender, Age), we have the counts of how many people survived and did not. We've also included rollups on individual variables for convenience.

Next to the table is a *mosaic plot*, which simply visualizes the counts as proportional areas.¹

2.a Train a decision tree

We are interested in predicting the outcome variable Y , survival, as a function of the input features C, G, A .

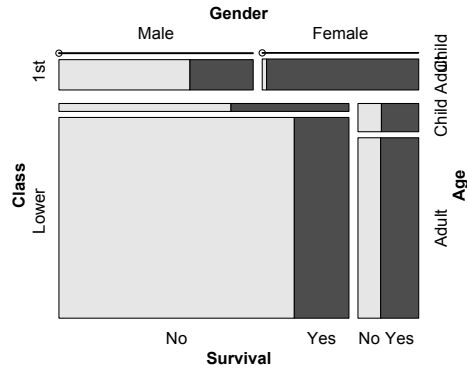
Use the information gain criterion to choose which of the three features C, G or A to use at the root of the decision tree. In fact, your task here is to learn a depth 1 decision tree that uses only this root feature to classify the data (such depth-1 decision trees are often called "decision stumps"). Please show all work, including the information gain calculations for each candidate feature.

Hint: to make information gain easier to calculate, you may wish to use this formula for conditional entropy:

$$-H(Y|X) = \sum_{x,y} p(x,y) \log p(y|x)$$

¹ From R packages *vcd* and *reshape2*, and built-in dataset *Titanic*. The original data has four values for Class; we collapsed 2nd, 3rd, and Crew into "Lower".

Class	Gender	Age	No	Yes	Total
1st	Male	Child	0	5	5
1st	Male	Adult	118	57	175
1st	Female	Child	0	1	1
1st	Female	Adult	4	140	144
Lower	Male	Child	35	24	59
Lower	Male	Adult	1211	281	1492
Lower	Female	Child	17	27	44
Lower	Female	Adult	105	176	281
totals:			1490	711	2201



Class	No	Yes	Total	Gender	No	Yes	Total	Age	No	Yes	Total
1st	122	203	325	Male	1364	367	1731	Child	52	57	109
Lower	1368	508	1876	Female	126	344	470	Adult	1438	654	2092

2.b Evaluation

1. What is the accuracy rate of your decision stump (depth 1 decision tree) on the training data?
2. If you grew a complete decision tree that used all three variables, what would its accuracy be over the training data? [Hint: you don't actually need to grow the tree to figure out the answer.]

2.c Decision Trees and Equivalent Boolean Expressions

The decision tree is a function $h(C, G, A)$ that outputs a binary value. Therefore, it can be represented as a boolean logic formula.

Write a decision tree that is equivalent to the following boolean formula (i.e., a decision tree that outputs 1 when this formula is satisfied, and 0 otherwise).

$$(C \wedge \neg A \wedge \neg G) \vee (C \wedge A) \vee (\neg C \wedge G)$$

2.d Model complexity and data size

Let's think about a situation where there is a true boolean function underlying the data, so we want the decision tree to learn it. We'll use synthetic data generated by the following algorithm. To generate an (\vec{x}, y) pair, first, six binary valued x_1, \dots, x_6 are randomly generated, each independently with probability 0.5. This six-tuple is our \vec{x} . Then, to generate the corresponding y value:

$$f(\vec{x}) = x_1 \vee (\neg x_1 \wedge x_2 \wedge x_6) \quad (1)$$

$$y = f(\vec{x}) \text{ with prob } \theta, \text{ else } (1 - f(\vec{x})) \quad (2)$$

So Y is a possibly corrupted version of $f(X)$, where the parameter θ controls the noisiness. $\theta = 1$ is noise-free. $\theta = 0.51$ is very noisy.

1. What is $P(Y = 1 \mid (X_1 \vee (\neg X_1 \wedge X_2 \wedge X_6)) = 1)$?
2. What is $P(Y = 1 \mid \neg((X_1 \vee (\neg X_1 \wedge X_2 \wedge X_6))) = 1)$?
3. Does $P(Y = 1 \mid X_2 = 1) = P(Y = 1)$? Why?
4. Does $P(Y = 1 \mid X_4 = 1) = P(Y = 1)$? Why?

5. Consider learning a decision tree classifier h . Assume the learning algorithm outputs a decision tree h that exactly matches f (despite the noise in the training data, it has so much data that it still learns f correctly). Assume the training data was generated by the above process. What should h 's accuracy rate be on the training data?
6. Assume new test data is also generated from the same process. What should its accuracy rate be on this new test data (assuming plenty of test data)?
7. Decision trees can overfit, so let's think about controlling the tree's model complexity. Instead of using pruning like we learned in lecture, here we use a maximum depth parameter.
Assuming a very large amount of training data, what's the smallest maximum-depth setting necessary to perfectly learn the generating function f ?

2.e Train/Test Experiments

Now we experimentally investigate the relationships between model complexity, training size, and classifier accuracy. Get code and test data from: http://www.cs.cmu.edu/~tom/10601_fall12012/hw/hw1_code.tgz

We provide a Matlab implementation of ID3, without pruning, but featuring a `maxdepth` parameter: `train_tree(trainX, trainY, maxdepth)`. It returns an object representing the classifier, which can be viewed with `print_tree(tree)`. Classify new data via `classify_with_tree(tree, testX)`. We also provide the simulation function to generate the synthetic data: `generate_data(N, theta)`, that you can use to create training data. Finally, there is a fixed test set for all experiments (generated using $\theta = 0.9$).

See `tt1.m` for sample code to get started.

Include printouts of your code and graphs.

1. For a depth=3 decision tree learner, learn classifiers for training sets size 10 and 100 (generate using $\theta = 0.9$). At each size, report training and test accuracies.
2. Let's track the learning curves for simple versus complex classifiers.

For `maxdepth=1` and `maxdepth=3`, perform the following experiment. For each training set size $\{2^1, 2^2, \dots, 2^{10}\}$, generate a training set, fit a tree, and record the train and test accuracies. For each (depth, trainsize) combination, average the results over 20 different simulated training sets.

Make three learning curve plots, where the horizontal axis is training size, and vertical axis is accuracy. First, plot the two testing accuracy curves, for each `maxdepth` setting, on the same graph. For the second and third graphs, have one for each `maxdepth` setting, and on each plot its training and testing accuracy curves. Place the graphs side-by-side, with identical axis scales. It may be helpful to use a log-scale for data size.

Next, answer several questions with *no more than three sentences* each:

3. When is the simpler model better? When is the more complex model better?
4. When are train and test accuracies different? If you're experimenting in the real world and find that train and test accuracies are substantially different, what should you do?
5. For a particular `maxdepth`, why do train and test accuracies converge to the same place? Comparing different `maxdepth`s, why do test accuracies converge to different places? Why does it take smaller or larger amounts of data to do so?
6. For `maxdepths` 1 and 3, repeat the same vary-the-training-size experiment with $\theta = 0.6$ for the training data. Show the graphs. Compare to the previous ones: what is the effect of noisier data?

3 Maximum Likelihood and MAP Estimation

This question asks you to explore a simple case of maximum likelihood and MAP estimation. The material for this question will not be covered in class until Tuesday, September 11, so you might want to wait until then to attempt it. Please print out all plots and code used to create them.

Our data is a set of n Boolean (0 or 1) values drawn independently from a single Bernoulli probability distribution, for which $P(X = 1) = \theta$, and therefore $P(X = 0) = 1 - \theta$. We define n Boolean-valued random variables, $X_1 \dots X_n$ to represent the outcomes of these n distinct draws. This problem asks you to explore how to estimate the value of θ from the observed values $X_1 \dots X_n$.

Turn in printouts of your graphs.

3.a Maximum Likelihood Estimate

1. Write a formula for $P(X_1 \dots X_n | \theta)$ in terms of θ . This is called the dataset's *likelihood*. We write $L(\theta) = P(X_1 \dots X_n | \theta)$, to indicate that the likelihood of the data $X_1 \dots X_n$ is a function of θ .
2. Assume a dataset size $n = 9$, consisting of 6 heads and then 3 tails:

$$(X_1, \dots, X_n) = (1, 1, 1, 1, 1, 1, 0, 0, 0)$$

Plot the likelihood curve as a function of θ , using a fine-grained grid of θ values, say for $\theta \in \{0, 0.01, 0.02, \dots, 1\}$. For the plot, the x-axis should be θ and the y-axis $L(\theta)$. Scale your y-axis so that you can see some variation in its value. Make sure to turn in both the plot and code that made it (should only be 3 or so lines of code). [Hint: In Matlab, it's easiest to first create the vector of θ values, then compute a corresponding vector of $L(\theta)$ values.]

3. In class we discussed that the maximum likelihood estimate of θ , which we call θ^{MLE} is the value that maximizes the likelihood $L(\theta)$:

$$\theta^{MLE} = \arg \max_{\theta} L(\theta)$$

On your plot, mark the value of θ along the x-axis that maximizes the likelihood. Does your θ^{MLE} agree with the following closed-form maximum likelihood estimator for a binomial distribution, which we learned in class?

$$\theta^{MLE} = \frac{\sum_i X_i}{n}$$

4. Create two more likelihood plots: one for a dataset of 2 heads and 1 tail; and one for a dataset of 40 heads and 20 tails.
5. Describe how the likelihood curves, maximum likelihoods, and maximum likelihood estimates compare?

3.b MAP Estimation

This section asks you to explore Maximum A Posteriori Probability (MAP) estimation of θ , in contrast to Maximum Likelihood estimation. Whereas the maximum likelihood estimate chooses a θ to maximize $P(X_1 \dots X_n | \theta)$, the MAP estimate instead chooses the θ that maximizes $P(\theta | X_1 \dots X_n)$. That is,

$$\theta^{MAP} = \arg \max_{\theta} P(\theta | X_1 \dots X_n)$$

which, by Bayes rule, is the same as

$$\theta^{MAP} = \arg \max_{\theta} \frac{P(X_1 \dots X_n | \theta) P(\theta)}{P(X_1 \dots X_n)}$$

and since the denominator $P(X_1 \dots X_n)$ is independent of θ this is equivalent to the simpler

$$\theta^{MAP} = \arg \max_{\theta} P(X_1 \dots X_n | \theta) P(\theta) \quad (3)$$

Thus, to find θ^{MAP} we just need to find the θ that maximizes $P(X_1 \dots X_n | \theta) P(\theta)$. This requires that we choose some probability distribution $P(\theta)$ that represents our prior assumptions about which values of θ are most probable before we have seen the data. For this, we will use the $Beta(\theta; \beta_H, \beta_T)$ distribution:

$$P(\theta) = \frac{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1}}{B(\beta_H, \beta_T)} = Beta(\theta; \beta_H, \beta_T) \quad (4)$$

where the denominator $B(\beta_H, \beta_T)$ is a normalizing function that does not depend on θ . Therefore, we ignore this denominator when maximizing θ .

1. Let's use a $Beta(\theta; 3, 3)$ distribution as our prior $P(\theta)$. Plot this as a function of θ . [Hint: The value of the normalizing denominator $B(3, 3) = 0.0333$].
2. Now plot the expression in the argmax of Equation 3, versus θ . Use your earlier data set containing 6 heads and 3 tails, and use $P(\theta) = Beta(\theta; 3, 3)$ as your prior. Where is the maximum on this plot? How does your θ^{MAP} from this plot compare with your earlier θ^{MLE} estimate for the same 6 heads, 3 tails data?
3. Above you used a $Beta(\theta; 3, 3)$ prior. Can you pick a different $Beta(\theta; \beta_H, \beta_T)$ distribution that, when used as a prior along with the earlier 2 heads and 1 tail data, will yield a $P(\theta|D)$ that has the same shape as your likelihood plot for 6 heads and 3 tails? If so, explain *in at most two sentences* why you are sure these two curves will have identical shape. If not, explain *in at most two sentences* why this is impossible.