

The Charlatan Problem

CS/DSA 5703: Machine Learning Practice

The Charlatan Problem

What is a good algorithm for empirically choosing a stock broker?

The Charlatan Problem

What is a good algorithm for empirically choosing a stock broker? One possibility:

- We ask the broker to make a judgement on a set of stocks as to whether they will go up or down in value in the next week (each is a binary question)
- At the end of the week, we ask how many the broker got right (also binary questions)

The Charlatan Problem

How do we evaluate this statistically?

The Charlatan Problem

How do we evaluate this statistically?

- Null hypothesis: the broker is a charlatan & doesn't really know how to choose
- Assume that choices are just made with a coin flip ($p = 0.5$)
- We will assume that the true probability of going up is also $p = 0.5$
- How well do we expect the broker to do in this case?

Guessing for a Single Stock

Under the null hypothesis, the charlatan will be correct for any one guess 50% of the time

- How about with N stocks?

Guessing the Outcome of N Stocks

We expect half of the guesses will be correct

- Are there other possible outcomes?

Guessing the Outcome of N Stocks

Are there other possible outcomes?

- Yes! $N/2-1$ and $N/2+1$ are equally likely (with each other)
- The probability drops off as we get further away from $N/2$
- As N gets large, what does this distribution look like?

Guessing the Outcome of N Stocks

- As N gets large, the distribution tends toward a Gaussian with mean $N/2$
- Central Limit Theorem: the sum (or mean) of N samples from any distribution tends towards a Gaussian distribution as N gets large
- $N=30$ is a good place to be ($N=20$ is still very close)

Guessing the Outcome of N Stocks

How do we decide whether to hire the stock broker after we have done this experiment?

Guessing the Outcome of N Stocks

How do we decide whether to hire the stock broker after we have done this experiment?

- If the number correct is large enough, then the probability of guessing correctly is a small probability under the null hypothesis assumption

NOTES_M9_L02b

- Probability distribution under the null hypothesis
- Statistical value: a particular observation
- Integral of likelihood above the critical value: the probability of observing that or greater value

Guessing the Outcome of N Stocks

Some terms:

- p-value: the estimated probability of incorrectly rejecting the null hypothesis
- alpha-value: the largest acceptable probability of incorrectly rejecting the null hypothesis
- critical value: the value of the statistic that corresponds to the alpha-value
 - In our case, this is the number of correctly selected stocks at which we accept that the stock broker is not a charlatan

Choosing Alpha

Choosing alpha depends on context and how well we need to trust the result

- Typical: 5%
- Stretch: 10% (but some will argue)
- Life and death: 1% - .1%

This is just hypothesis testing ...

- Assume that there is no difference in the way two models perform
- Ask how likely it is that we observe samples of performance from each of the two models under this assumption
- If the probability is too small, then we reject the assumption (which is what we want)

Implications

- A model, even if it is not better than a competitor, can look good with some probability
- But, we can control the probability of making a mistake if we present the model with enough tests and have a high enough criterion

The Many Charlatans Problem

CS/DSA 5703: Machine Learning Practice

The Multiple Charlatans Problem

- We have already discussed a test that will reveal whether a stock broker can be hired. This test only makes a mistake with probability α
- How do we make the search process for a broker more efficient?

The Multiple Charlatans Problem

How do we make the search process for a broker more efficient?

- Let's test K brokers in parallel!
- Each broker gets the same stocks to judge
- Choose the broker with the highest accuracy
 - Must also have a low enough p-value

The Multiple Charlatans Problem

The broker with the highest accuracy is selected and this accuracy is above a critical threshold

- What is the probability that we have made a mistake?
- I.E.: assume that all are charlatans. What is the probability that we still accept someone?

NOTES_M8_L03b

- Math for computing aggregate alpha assuming alphas of 0.05
- Math for computing individual alphas given an aggregate alpha of 0.5
- Corrections options: Bonferroni vs Sidak

Solving Multiple Comparisons Problem

Comparing one sampling against more than one other sampling dilutes the power of the individual comparisons.

Options for addressing:

- Correct the alpha
 - This is a very conservative approach, but is effective
- Once we have selected the best, we take new samples to do the final comparisons

Selecting Model Types and Hyper-Parameters

This is really a multi-level question

- For a given model type, we first need to know what the best hyper-parameter set is. This can involve **a lot** of comparisons
- Then, we can begin to compare model types
- Typical approach: use different data sets for these two levels

Selecting Model Types and Hyper-Parameters

Typical approach: use independent data for these two levels

- Validation data: use for selection of hyper-parameters
- Test data: use for comparing models

Comparing Models

CS/DSA 5703: Machine Learning Practice

Holistic Cross-Validation

- When we are selecting model parameters, we potentially have a large number of parameter sets
- Neighboring parameter sets should yield similar performance
 - Be suspicious if this is not the case
- We address the many charlatans problem by using independent data sets:
 - Validation data set: parameter selection
 - Test data set: only use at the very end when you are justifying your use of the model (including comparison to a small number of alternatives)

Holistic Cross-Validation: Comparing Models

For each model type, we have:

- N performance measures (one for each model instance / test set combination)
- Because these performance measures are independent of the data used for hyper-parameter selection, we won't be fooled by hyper-parameter over-fit
- Have the opportunity to perform paired comparisons
 - Pairs are between models for a given test data set

NOTES_M8_L07

- Hyper parameter selection: gives way to best choice for each model
- Paired samples
- Can simply pick the best model: note that this involves $O(M)$ comparisons to find the max
 - Statistically comparing the best against the next best has to take these M comparisons into account
- Comparing M models against a baseline model: if we want to make a statistical argument that one of the M is better than the baseline, then we assume that we have M charlatans that we are comparing against the baseline

Comparing Model Types

Compare two models statistically:

- Often see a student t-test (one tailed)
 - Nominally assumes an underlying normal distribution of the performance metric
- Sample based approaches: computationally involved, but are robust to any underlying distribution
 - For example: bootstrap resampling and bootstrap randomization
 - And, can be constructed to address lack of independence in the samples or even censored samples

Comparing M Model Types

- If we have more than 2 model types in question, then we again have the multiple charlatans problem
- ANOVA:
 - Tells us that at least one model is different from the others (but not which one)
 - But, this gives us permission to perform pairwise tests

Comparing M Model Types

- We can perform all $M * M-1 / 2$ comparisons
 - But this is often overkill
- A couple common alternatives:
 - We want to pick the best of M & make a statistical argument that this is the best one
 - This involves $O(M)$ comparisons
 - We want to show that at least one of M is statistically better than some baseline model
 - Exactly M comparisons are made here
- In any of these cases: we use Bonferroni or Sidèk correction to adjust the p-value cutoff

Comparing M Model Types

- Remember that Bonferroni and Sidèk are very conservative adjustments to the p-value cutoff
- An alternative: use new data
- Cross-validation then involves 4 data sets:
 - Train, validation, test1, test2
 - Use test1 to pick the favorite model of M (and the next best one)
 - Use test2 to confirm that it is the best (a single statistical test!)

Final Thoughts

- Hyper-parameters:
 - If you are doing tests that examine how much training data you need for a given model, this is also a hyper-parameter
 - Make this selection based on validation performance
- Data cutting:
 - Be sure that samples across the folds are independent
 - Time series data can have a lot of autocorrelation
 - Cutting a contiguous block of data (in time) can yield folds that exhibit some autocorrelation

Final Thoughts

Cross-Validation (even Holistic) can be fooled

- The cutting of the data can give you “lucky” results
- Typical approach (not very common, though):
 - Re-cut the data into a new set of folds and repeat the procedure (multiple times)
 - We expect consistent results
 - Beyond our scope: formal methods for deciding how many times to repeat

