

# Master's Thesis Proposal

Marcel Gietzmann-Sanders

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Appendices</b>	<b>3</b>

## 1 Introduction

### 1.1 Breaking Down Modeling

One of the functions of science is to build models of the world around us. Specifically to look for functions of the form:

$$v = F(\theta)$$

where  $\theta$  is some vector of information,  $v$  is a particular outcome given that information, and  $F$  is the model in question. This we might term a *deterministic model* because for any specific  $\theta$  there is one single outcome  $v$ .

However, life is rarely so kind to us. From quantum physics to fisheries science the world is replete with examples where the same information does not always lead to the same conclusion. Sometimes that is because we do not possess all of the relevant information but sometimes things are just truly random. As such, the model above will in some sense be lying to us in its deterministic certainty - while we will always predict  $v$ ,  $v$  is not always what we will receive.

Instead suppose that for a specific  $\theta$  there exist a set of possible outcomes  $V(\theta) = \{v\}$ . We can define a *stochastic model* as:

$$P(v|\theta) = F(v, \theta)$$

where now the model is predicting the probability of  $v$  given  $\theta$  instead of just predicting  $v$  itself.

Now what is especially interesting about models of this form is the fact that in some sense for any combination of  $V$  and  $\theta$  they always exist. Whereas to get a deterministic model we need to know precisely what kind of information ( $\theta$ ) we need to make a deterministic prediction, we can always make a stochastic one, even in the presence of no information. The model can always give us a correct answer, it's just a question of how useful that answer is. If you base it off of better data it will become more useful.

Unfortunately, unless we ourselves have built the part of the world we're studying,  $F$  is not known to us. Instead our models are really just proposals  $\hat{F}$  on what  $F$  could be. This presents us with an immediate problem, how do we evaluate any specific  $\hat{F}$  if we don't know  $F$  itself? Well suppose we've captured a whole series of pairs of  $\theta_i$  and  $v_i$  where  $i$  allows us to index the pair in question. Given we know  $\hat{F}$  we can directly ask what the likelihood of the data we have collected is, given  $\hat{F}$ :

$$\mathcal{L} = \prod_i \hat{F}(v_i, \theta_i)$$

or equivalently (and more easily computed)

$$\ln \mathcal{L} = \sum_i \ln \hat{F}(v_i, \theta_i)$$

$\hat{F}$ 's with higher  $\ln \mathcal{L}$  (log likelihood) are better fits to the given data and as we have more and more comprehensive data those data better and better represent  $F$  (see Appendix 1).

Therefore finding the "true"  $F$  can be summarized with two steps:

1. Collect large amounts of comprehensive data.
2. Find the  $\hat{F}$  that maximizes the likelihood of that data.

So is this all there is to modeling? Absolutely not. Everything we've said so far is dependent upon a chosen information set  $\theta$ . For each choice there is a  $F$  and obviously the usefulness of these is a balance between how accessible  $\theta$  is when the model needs to be used and how predictive  $F$  is (i.e. how much variance is represented by it's predicted  $P(v|\theta)$ ). Determining the best  $\theta$  is really what makes any kind of modeling an art form.

## 1.2 Searching for Maximal Functions - Machine Learning

While the first step outlined is obviously highly dependent on what we are trying to model, once the data has been collected

## 2 Appendices

### 2.1 Proof of Log Likelihood Maximization

**Theorem 1.** *Suppose we have a set of possible outcomes  $V = \{v_k\}$  with probabilities  $P_k$ . Such that:*

$$\sum_k P_k = P$$

*Next suppose we concoct a new series of probabilities  $U_k = P_k \alpha_k$  s.t.*

$$\sum_k U_k = P - \epsilon$$

*where  $\epsilon \geq 0$ . There does not exist a set of  $\alpha_k$  s.t.*

$$\prod_k \left( \frac{P_k \alpha_k}{P_k} \right)^{P_k N} > 1$$

*Proof.* We proceed by induction.

First note that from the equation above we get:

$$\sum_k P_k \ln \alpha_k > 0$$

and subtracting  $\sum_k U_k$  we arrive at:

$$\sum_k P_k (\ln \alpha_k - \alpha_k) > -P + \epsilon$$

Now suppose for  $k$  outcomes we know no such  $\alpha_k$  exist as to satisfy the above. Suppose we incorporate a new  $k+1$  term s.t:

$$P_{k+1} + \sum_k P_k = P_{k+1} + P$$

and:

$$P_{k+1} \alpha_{k+1} + \sum_k P_k \alpha_k = P_{k+1} + P - \epsilon = P_{k+1} \alpha_{k+1} + (P - P_{k+1} (\alpha_{k+1} - 1)) - \epsilon$$

Now we need:

$$P_{k+1} (\ln \alpha_{k+1} - \alpha_{k+1}) + \sum_k P_k (\ln \alpha_k - \alpha_k) > -P_{k+1} - P + \epsilon$$

However given our inductive assumption we know that at best,

$$\sum_k P_k (\ln \alpha_k - \alpha_k) = -P$$

therefore at best:

$$P_{k+1} (\ln \alpha_{k+1} - \alpha_{k+1}) - P > -P - P_{k+1} + \epsilon$$

Furthermore the easiest case for us is if  $\epsilon = 0$ . If it cannot be satisfied than this certainly doesn't work for  $\epsilon > 0$ . So let's consider:

$$P_{k+1} (\ln \alpha_{k+1} - \alpha_{k+1}) - P > -P - P_{k+1}$$

or equivalently:

$$\ln \alpha_{k+1} > \alpha_{k+1} - 1$$

Now if  $\alpha_{k+1} = 1$  we have:

$$\ln 1 = 1 - 1$$

Further consider the derivatives of each side:

$$\partial_\alpha \ln \alpha_{k+1} = \frac{1}{\alpha_{k+1}}$$

$$\partial_\alpha (\alpha_{k+1} - 1) = 1$$

If  $\alpha_{k+1} > 1$  then the log component rises more slowly than the constant component. I.e. our left side will be larger than the right. Likewise if  $\alpha_{k+1} < 1$  the log component will shrink faster than the constant component which means it will also be less than the constant component. Therefore given our assumptions:

$$\ln \alpha_{k+1} \not> \alpha_{k+1} - 1$$

and so:

$$P_{k+1} (\ln \alpha_{k+1} - \alpha_{k+1}) - P \not> -P - P_{k+1} + \epsilon$$

So much for the  $k + 1$ th case. What about  $k = 1$ . This is trivial because any  $\alpha_1$  that satisfies:

$$\sum_k U_k = P - \epsilon$$

must be less than or equal to 1 and therefore our product of quotients:

$$\prod_k \left( \frac{P_k \alpha_k}{P_k} \right)^{P_k N} > 1$$

must be less than or equal to one.

□

With this proof in hand we can now show how in the limit as the number of samples taken  $N$  goes to infinity our likelihood is only maximized if  $\hat{F} \rightarrow F$ .

For some given information  $\theta$  and outcomes  $V = \{v_k\}$  we have, in gathering our data, observed  $v_k$   $N_k$  times. Therefore our overall likelihood is:

$$\mathcal{L} = \prod_k \hat{F}(v_k, \theta)^{N_k}$$

Now suppose that we represent:

$$\hat{F}(v_k, \theta) = F(v_k, \theta) \alpha_k$$

That is the ratio of our likelihoods between  $\hat{F}$  and  $F$  is given by:

$$\prod_k \left( \frac{F(v_k, \theta) \alpha_k}{F(v_k, \theta)} \right)^{N_k}$$

but we also know that:

$$\sum_k F(v_k, \theta) = \sum_k F(v_k, \theta) \alpha_k = 1$$

and that  $\lim_{N \rightarrow \infty} N_k = F(v_k, \theta)N$  so that we have:

$$\prod_k \left( \frac{F(v_k, \theta) \alpha_k}{F(v_k, \theta)} \right)^{F(v_k, \theta)N}$$

which is now in the exact same form as our theorem above. And we now know that this ratio is maximized when  $\alpha_k \equiv 1$ . So as  $N \rightarrow \infty$  a higher  $\mathcal{L}$  means we are closer to approximating  $F$ .