

# Machine Learning Engineer Nanodegree

## Capstone Proposal

Hamdy Tawfeek

January 17th, 2019

## Proposal

---

### Domain Background

In a one-click shopping world with on-demand everything, the life insurance application process is antiquated. Customers provide extensive information to identify risk classification and eligibility, including scheduling medical exams, a process that takes an average of 30 days.

The result? People are turned off. That's why only 40% of U.S. households own individual life insurance. Prudential wants to make it quicker and less labor intensive for new and existing customers to get a quote while maintaining privacy boundaries.

### Problem Statement

Can you make buying life insurance easier?

By developing a predictive model that accurately classifies risk using a more automated approach, you can greatly impact public perception of the industry.

## Datasets and Inputs

The dataset provided by Prudential on [Kaggle](#).

In this dataset, a hundred variables describing attributes of life insurance applicants. The task is to predict the “Response” variable for each Id in the test set. “Response” is an ordinal measure of risk that has 8 levels.

### File descriptions

- train.csv - the training set, contains the Response values
- test.csv - the test set, you must predict the Response variable for all rows in this file

### Data fields

Variable	Description
Id	A unique identifier associated with an application.
Product_Info_1-7	A set of normalized variables relating to the product applied for.
Ins_Age	Normalized age of applicant.
Ht	Normalized height of applicant
Wt	Normalized weight of applicant

BMI	Normalized BMI of applicant
Employment_Info_1-6	A set of normalized variables relating to the employment history of the applicant.
InsuredInfo_1-6	A set of normalized variables providing information about the applicant.
Insurance_History_1-9	A set of normalized variables relating to the insurance history of the applicant.
Family_Hist_1-5	A set of normalized variables relating to the family history of the applicant.
Medical_History_1-41	A set of normalized variables relating to the medical history of the applicant.
Response	This is the target variable, an ordinal variable relating to the final decision associated with an application

The following variables are all categorical (nominal):

Product\_Info\_1, Product\_Info\_2, Product\_Info\_3, Product\_Info\_5, Product\_Info\_6, Product\_Info\_7, Employment\_Info\_2, Employment\_Info\_3, Employment\_Info\_5, InsuredInfo\_1, InsuredInfo\_2, InsuredInfo\_3, InsuredInfo\_4, InsuredInfo\_5, InsuredInfo\_6, InsuredInfo\_7, Insurance\_History\_1, Insurance\_History\_2, Insurance\_History\_3, Insurance\_History\_4, Insurance\_History\_7, Insurance\_History\_8, Insurance\_History\_9, Family\_Hist\_1, Medical\_History\_2, Medical\_History\_3, Medical\_History\_4, Medical\_History\_5, Medical\_History\_6, Medical\_History\_7, Medical\_History\_8, Medical\_History\_9, Medical\_History\_11, Medical\_History\_12, Medical\_History\_13, Medical\_History\_14,

Medical\_History\_16, Medical\_History\_17, Medical\_History\_18,  
Medical\_History\_19, Medical\_History\_20, Medical\_History\_21,  
Medical\_History\_22, Medical\_History\_23, Medical\_History\_25,  
Medical\_History\_26, Medical\_History\_27, Medical\_History\_28,  
Medical\_History\_29, Medical\_History\_30, Medical\_History\_31,  
Medical\_History\_33, Medical\_History\_34, Medical\_History\_35,  
Medical\_History\_36, Medical\_History\_37, Medical\_History\_38,  
Medical\_History\_39, Medical\_History\_40, Medical\_History\_41

The following variables are continuous:

Product\_Info\_4, Ins\_Age, Ht, Wt, BMI, Employment\_Info\_1,  
Employment\_Info\_4, Employment\_Info\_6, Insurance\_History\_5,  
Family\_Hist\_2, Family\_Hist\_3, Family\_Hist\_4, Family\_Hist\_5

The following variables are discrete:

Medical\_History\_1, Medical\_History\_10, Medical\_History\_15,  
Medical\_History\_24, Medical\_History\_32  
Medical\_Keyword\_1-48 are dummy variables.

## **Solution Statement**

In this project, I would like to use compare different predictive classification models and see which is better. I will implement logistic regression, XGBoost, Neural Network.

## **Benchmark Model**

Prudential Kaggle competition used a **benchmark score** of 0. So, I will

take a score of 0 as my benchmark for this project.

## Evaluation Metrics

Submissions are scored based on the quadratic weighted kappa, which measures the agreement between two ratings. This metric typically varies from 0 (random agreement) to 1 (complete agreement). In the event that there is less agreement between the raters than expected by chance, this metric may go below 0.

The response variable has 8 possible ratings. Each application is characterized by a tuple (ea,eb), which corresponds to its scores by Rater A (actual risk) and Rater B (predicted risk). The quadratic weighted kappa is calculated as follows.

First, an  $N \times N$  histogram matrix  $O$  is constructed, such that  $O_{i,j}$  corresponds to the number of applications that received a rating  $i$  by A and a rating  $j$  by B. An  $N$ -by- $N$  matrix of weights,  $w$ , is calculated based on the difference between raters' scores:

$$w_{i,j} = \frac{(i-j)^2}{(N-1)^2}$$

An  $N$ -by- $N$  histogram matrix of expected ratings,  $E$ , is calculated, assuming that there is no correlation between rating scores. This is calculated as the outer product between each rater's histogram vector of ratings, normalized such that  $E$  and  $O$  have the same sum.

From these three matrices, the quadratic weighted kappa is calculated

as:

$$\kappa = 1 - \frac{\sum_{i,j} w_{i,j} O_{i,j}}{\sum_{i,j} w_{i,j} E_{i,j}}.$$

## Project Design

The project will roughly consist of the following steps:

1. Understanding the data.
2. Data Preprocessing
3. Model Building and evaluation.
4. Compare different results from different models and choose the model with the highest predictive power.