# Machine Learning Engineer Nanodegree

## Projeto Final

Gleber A. Baptistella
Fevereiro de 2018

# I. Definição

## Visão geral do projeto

A área de Reconhecimento de Atividades Humanas (RAH) está em franca expansão, com dispositivos e sensores cada vez mais fazendo parte do dia-a-dia das pessoas.

Suas pesquisas podem ser aplicadas em diversas áreas, entre elas da saúde (melhorando a detecção e diagnósticos de doenças) e de práticas esportivas, visando a melhoria do desempenho do indivíduo.

## Descrição do problema

Este projeto tem como objetivo identificar uma determinada postura de um indivíduo após a coleta de dados através de sensores que foram previamente dispostos no corpo da pessoa. O embasamento deste projeto encontra-se em **http://groupware.les.inf.puc-rio.br/har**, bem como o dataset utilizado.

O dataset disponibilizado está em formato CSV e está bem estruturado, facilitando a leitura.

A saída do processamento é um valor categórico que representa a posição do indivíduo e suas possibilidades são: sitting, sittingdown, standing, standingup, walking. Em português seria algo como: sentado, sentando, de pé, levantando e andando.

Dessa forma, trata-se de um problema de aprendizado supervisionado de classificação.

Irei avaliar outros algoritmos de classificação diferentes do usado no artigo da PUC Rio, buscando um desempenho superior ao obtido.

## Métricas

As métricas para avaliação dos modelos serão:

1. Recall
2. Precision
3. F1
4. ROC AUC

Todas as métricas são referentes aos problemas de classificação.

É comum utilizar a métrica *Accuracy* nos problemas de classificação, porém em alguns casos ela não é uma boa métrica. Tomemos como exemplo um sistema de classificação de fraudes em transações de cartão de crédito. A grande maioria das transações, digamos 99%, são legítimas. Apenas cerca de 1% das transações são fraudulentas. Se fizermos um modelo de predição onde considero que todas as transações são legítimas, minha taxa de acerto será de 99%, ou seja, a métrica acurácia será de 99%. Contudo este não é um modelo confiável apesar da alta taxa de acerto, já que todas as transações fraudulentas serão classificadas como legítimas. Para situações como esta temos as métricas referidas: Recall, Precision, F1 e ROC AUC.

As métricas levam em consideração as taxas de falsos positivos, verdadeiros positivos, verdadeiros negativos e falsos negativos para analisar a qualidade do modelo. Abaixo uma *confusion matrix* representando os resultados possíveis num problema de classificação:

Dessa forma definimos as métricas da seguinte forma aplicando no exemplo do cartão de crédito:
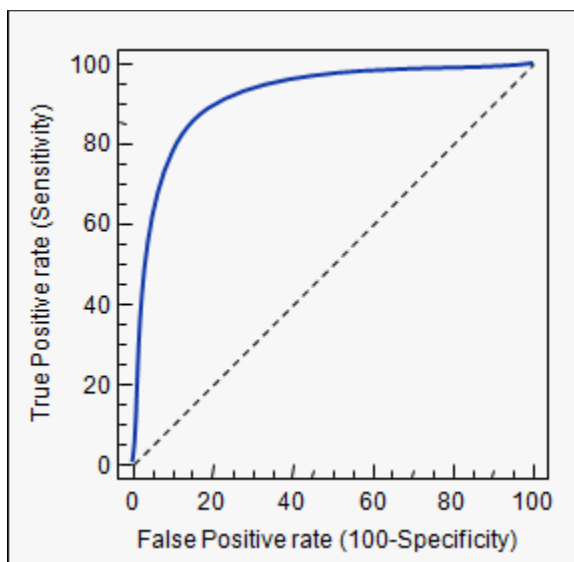
**Recall:** é a capacidade do modelo predizer corretamente as operações legítimas. Seguindo pela tabela a fórmula seria: `VP/(VP+FN)`

**Precision**: é a capacidade do modelo predizer corretamente as operações fraudulentas. Seguindo pela tabela a fórmula seria: `VN/(VN+FP)`

**F1**: é a média harmônica entre Recall e Precision. A fórmula é: `2 * (precision * recall)/(precision + recall)`. O resultado é um número entre 0 e 1, onde 0 é o pior modelo e 1 é o modelo mais preciso.

**ROC AUC:** a curva ROC (*Receiver Operating Characteristics*) é um gráfico bidimensional que utiliza a taxa de verdadeiros positivos no eixo Y e a taxa de falsos positivos no eixo X.



A linha tracejada no meio do gráfico seria um modelo aleatório. A métrica AUC (*area under curve*) é a área abaixo da curva ROC e pode variar de 0 a 1, sendo 0 o pior valor e 1 o melhor valor.
Como veremos na sessão "**II. Análise**", o dataset para o problema proposto também encontra-se desbalanceado. Dessa forma as métricas propostas são as mais adequadas para verificarmos a qualidade do modelo.

# II. Análise

*(approx. 2-4 pages)*

## Exploração de dados

O conjunto de dados selecionado encontra-se disponível em http://groupware.les.inf.puc-rio.br/static/har/dataset-har-PUC-Rio-ugulino.zip e seus atributos estão descritos abaixo:

Dados Pessoais:

- user: nome da pessoa
- gender: sexo da pessoa
- how_tall_in_meters: altura da pessoa
- weight: peso da pessoa
- body_mass_index: Índice de massa corpórea da pessoa

Dados dos sensores:
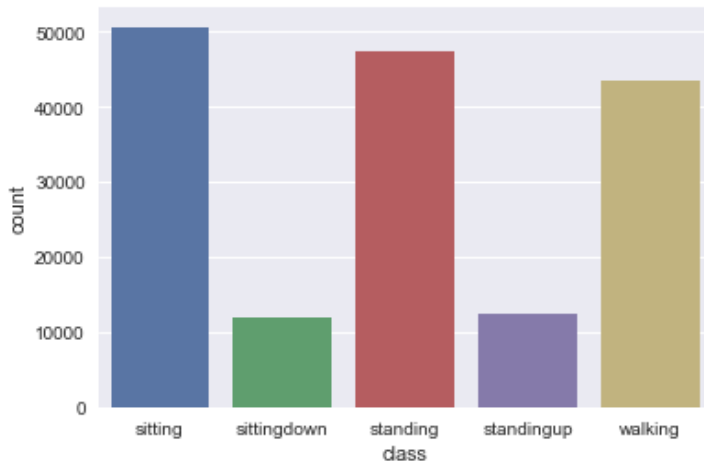
São colocados 4 sensores no corpo do indivíduo.

Para cada sensor são definidos os eixos x, y e z e no dataset temos para o sensor 1 os dados, x1, y1 e z1. Para o sensor 2, os dados x2, y2 e z2, e assim por diante. A posição de cada sensor é:

- Sensor 1: cintura
- Sensor 2: coxa esquerda
- Sensor 3: canela direita
- Sensor 4: braço direito.

Por fim, temos a *feature* "class". Trata-se da variável que queremos prever e que pode assumir os seguintes valores: 'sitting', 'sittingdown', 'standing', 'standingup', 'walking'

Durante a coleta de dados, que consistiu em 8 horas de atividades de quatro indivíduos distintos, foram geradas 165362 amostras dos sensores. Trata-se de um bom volume de dados para aplicação dos algoritmos.

A distribuição das classes de classificação do dataset estão distribuídas da seguinte forma:

- sittingdown = 7,15%
- standingup = 7,50%
- walking = 26,23%
- standing = 28,64%
- sitting = 30,45%

Aparentement as duas primeiras classes estão sub-representadas, e pode ser que tenhamos que balancear o dataset para termos uma distribuição melhor entre as classes.

described and, if possible, have basic statistics and information presented (such as discussion of input features or defining characteristics about the input or environment). Any abnormalities or interesting qualities about the data that may need to be addressed have been identified (such as features that need to be transformed or the possibility of outliers). Questions to ask yourself when writing this section:

- *If a dataset is present for this problem, have you thoroughly discussed certain features about the dataset? Has a data sample been provided to the reader?*
- *If a dataset is present for this problem, are statistics about the dataset calculated and reported? Have any relevant results from this calculation been discussed?*
- *If a dataset is **not** present for this problem, has discussion been made about the input space or input data for your problem?*

- *Are there any abnormalities or characteristics about the input space or dataset that need to be addressed? (categorical variables, missing values, outliers, etc.)*

## Exploratory Visualization

In this section, you will need to provide some form of visualization that summarizes or extracts a relevant characteristic or feature about the data. The visualization should adequately support the data being used. Discuss why this visualization was chosen and how it is relevant. Questions to ask yourself when writing this section:

- *Have you visualized a relevant characteristic or feature about the dataset or input data?*
- *Is the visualization thoroughly analyzed and discussed?*
- *If a plot is provided, are the axes, title, and datum clearly defined?*

## Algorithms and Techniques

In this section, you will need to discuss the algorithms and techniques you intend to use for solving the problem. You should justify the use of each one based on the characteristics of the problem and the problem domain. Questions to ask yourself when writing this section:

- *Are the algorithms you will use, including any default variables/parameters in the project clearly defined?*
- *Are the techniques to be used thoroughly discussed and justified?*
- *Is it made clear how the input data or datasets will be handled by the algorithms and techniques chosen?*

## Benchmark

In this section, you will need to provide a clearly defined benchmark result or threshold for comparing across performances obtained by your solution. The reasoning behind the benchmark (in the case where it is not an established result) should be discussed. Questions to ask yourself when writing this section:

- *Has some result or value been provided that acts as a benchmark for measuring performance?*

- *Is it clear how this result or value was obtained (whether by data or by hypothesis)?*

# III. Methodology

*(approx. 3-5 pages)*

## Data Preprocessing

In this section, all of your preprocessing steps will need to be clearly documented, if any were necessary. From the previous section, any of the abnormalities or characteristics that you identified about the dataset will be addressed and corrected here. Questions to ask yourself when writing this section:

- *If the algorithms chosen require preprocessing steps like feature selection or feature transformations, have they been properly documented?*
- *Based on the **Data Exploration** section, if there were abnormalities or characteristics that needed to be addressed, have they been properly corrected?*
- *If no preprocessing is needed, has it been made clear why?*

## Implementation

In this section, the process for which metrics, algorithms, and techniques that you implemented for the given data will need to be clearly documented. It should be abundantly clear how the implementation was carried out, and discussion should be made regarding any complications that occurred during this process. Questions to ask yourself when writing this section:

- *Is it made clear how the algorithms and techniques were implemented with the given datasets or input data?*
- *Were there any complications with the original metrics or techniques that required changing prior to acquiring a solution?*
- *Was there any part of the coding process (e.g., writing complicated functions) that should be documented?*

## Refinement

In this section, you will need to discuss the process of improvement you made upon the algorithms and techniques you used in your implementation. For example, adjusting parameters for certain models to acquire improved solutions would fall under the refinement category. Your initial and final solutions should be reported, as well as any significant intermediate results as necessary. Questions to ask yourself when writing this section:

- *Has an initial solution been found and clearly reported?*
- *Is the process of improvement clearly documented, such as what techniques were used?*
- *Are intermediate and final solutions clearly reported as the process is improved?*

# IV. Results

*(approx. 2-3 pages)*

## Model Evaluation and Validation

In this section, the final model and any supporting qualities should be evaluated in detail. It should be clear how the final model was derived and why this model was chosen. In addition, some type of analysis should be used to validate the robustness of this model and its solution, such as manipulating the input data or environment to see how the model's solution is affected (this is called sensitivity analysis). Questions to ask yourself when writing this section:

- *Is the final model reasonable and aligning with solution expectations? Are the final parameters of the model appropriate?*
- *Has the final model been tested with various inputs to evaluate whether the model generalizes well to unseen data?*
- *Is the model robust enough for the problem? Do small perturbations (changes) in training data or the input space greatly affect the results?*
- *Can results found from the model be trusted?*

## Justification

In this section, your model's final solution and its results should be compared to the benchmark you established earlier in the project using some type of statistical

analysis. You should also justify whether these results and the solution are significant enough to have solved the problem posed in the project. Questions to ask yourself when writing this section:

- *Are the final results found stronger than the benchmark result reported earlier?*
- *Have you thoroughly analyzed and discussed the final solution?*
- *Is the final solution significant enough to have solved the problem?*

# V. Conclusion

*(approx. 1-2 pages)*

## Free-Form Visualization

In this section, you will need to provide some form of visualization that emphasizes an important quality about the project. It is much more free-form, but should reasonably support a significant result or characteristic about the problem that you want to discuss. Questions to ask yourself when writing this section:

- *Have you visualized a relevant or important quality about the problem, dataset, input data, or results?*
- *Is the visualization thoroughly analyzed and discussed?*
- *If a plot is provided, are the axes, title, and datum clearly defined?*

## Reflection

In this section, you will summarize the entire end-to-end problem solution and discuss one or two particular aspects of the project you found interesting or difficult. You are expected to reflect on the project as a whole to show that you have a firm understanding of the entire process employed in your work. Questions to ask yourself when writing this section:

- *Have you thoroughly summarized the entire process you used for this project?*
- *Were there any interesting aspects of the project?*
- *Were there any difficult aspects of the project?*

- *Does the final model and solution fit your expectations for the problem, and should it be used in a general setting to solve these types of problems?*

## Improvement

In this section, you will need to provide discussion as to how one aspect of the implementation you designed could be improved. As an example, consider ways your implementation can be made more general, and what would need to be modified. You do not need to make this improvement, but the potential solutions resulting from these changes are considered and compared/contrasted to your current solution. Questions to ask yourself when writing this section:

- *Are there further improvements that could be made on the algorithms or techniques you used in this project?*
- *Were there algorithms or techniques you researched that you did not know how to implement, but would consider using if you knew how?*
- *If you used your final solution as the new benchmark, do you think an even better solution exists?*

---

**Before submitting, ask yourself. . .**

- Does the project report you've written follow a well-organized structure similar to that of the project template?
- Is each section (particularly **Analysis** and **Methodology**) written in a clear, concise and specific fashion? Are there any ambiguous terms or phrases that need clarification?
- Would the intended audience of your project be able to understand your analysis, methods, and results?
- Have you properly proof-read your project report to assure there are minimal grammatical and spelling mistakes?
- Are all the resources used for this project correctly cited and referenced?
- Is the code that implements your solution easily readable and properly commented?
- Does the code execute without error and produce results similar to those reported?