
Instituto Federal do Norte de Minas Gerais

IFNMG – Campus Montes Claros

Bacharelado em Ciência da Computação

Disciplina de Inteligência Artificial

Seminário:

Random Forest

David Walter Jansen (davidwalterjansen@gmail.com)

Iarah Gonçalves de Almeida

Paulo Augusto Borges de Matos

Última atualização: 30 de maio de 2019

Sumário

Resumo	1
1 Introdução: Random Forest Algorithm	2
1.1 História e Motivação	2
1.2 Ideia do Algoritmo	2
1.3 Funcionamento do Algoritmo	3
1.4 Vantagens e Desvantagens	4
2 Exemplo de Aplicação	5

Resumo

O presente trabalho é referente a um seminário sobre modelos de aprendizagem de máquina, proposto pela Professora Dr^a. Luciana Balieiro Cosme na disciplina de Inteligência Artificial do curso de Ciência da Computação - IFNMG (Campus Montes Claros).

O modelo escolhido pela nossa equipe foi a Random Forest (Floresta Aleatória). Esse é um modelo de aprendizado supervisionado usado tanto para regressão quanto classificação [1]. Nosso foco se dará sobre essa último. Falaremos sobre a aplicação do método, como são usadas as árvores de decisão (outro modelo de aprendizagem), bem como suas vantagens e desvantagens. Abordaremos também uma das aplicações mais importantes desse modelo: Seleção de Características.

Capítulo 1

Introdução: Random Forest Algorithm

1.1 História e Motivação

O primeiro algoritmo para florestas de decisão aleatória foi proposto por Tim Kam Ho em 1995. Em seu [artigo](#), Tim Kam Ho menciona que *"as árvores derivadas de métodos tradicionais geralmente não podem ser expandidas até a complexidade arbitrária para possível perda de precisão de generalização em dados não vistos. A limitação da complexidade significaria precisão abaixo do ideal nos dados de treinamento."* [3]. Esse fato serviu de motivação para o desenvolvimento de um novo método, baseado em princípios da modelagem estocástica¹, que utiliza a construção de várias árvores em subespaços selecionados aleatoriamente sobre espaço dos dados totais.

Na trabalho inicial, Ho havia estabelecido as divisões em hiperplanos oblíquos. Posteriormente em 1998 [5], o mesmo autor mostrou que outros métodos de divisão são válidos. O algoritmo de Floresta Aleatória sofreu uma grande evolução com o trabalho de Breiman [2], onde foi introduzida a capacidade de se medir a importância das características.

1.2 Ideia do Algoritmo

Suponha que você está montando sua grade para o próximo período e precisa escolher mais uma disciplina. Usando o paradigma da divisão e conquista, uma das possibilidades é pedir recomendações a seus colegas. Cada colega poderá lhe dar algumas recomendações. Com essas recomendações, você terá uma lista de matérias para completar sua grade. Agora, realize uma votação com base nas recomendações de seus colegas. A disciplina com mais votos será a sua escolha final.

¹Um processo estocástico é formado por variáveis que se comportam, durante o tempo, de uma maneira onde pelo menos parte é considerada randômica.

O processo de escolha de disciplina que descrevemos é uma analogia ao funcionamento da decisão feita pela Floresta Aleatória. Chamamos de floresta a coleção de classificadores por árvores de decisão. Para ficar mais claro, iremos dividir a analogia feita em duas partes: a primeira parte é o recebimento de recomendações dos colegas. Nessa parte, o algoritmo da Floresta de decisão usa árvores de decisão. Cada colega (árvore) faz uma seleção das matérias que ele tem conhecimento. A segunda parte é o procedimento de votação. Após coletar todas as recomendações dos colegas, a votação para encontrar a melhor sugestão é uma dos processos do algoritmo de Florestas Aleatórias.

Na classificação, cada árvore vota e a classe mais popular é escolhida. Para a regressão, a média de todas as saídas de árvore é considerada como resultado final.

1.3 Funcionamento do Algoritmo

O algoritmo constrói o conjunto de árvores de decisão usando divisões aleatórias e independentes sobre o conjunto de dados. No processo de geração é usado um indicador de seleção de atributos, como ganho de informação, taxa de ganho e índice de Gini para cada atributo [1].

O funcionamento da Random Forest pode ser resumido em 4 etapas:

- (1) Selecione amostras aleatórias de um determinado conjunto de dados.
- (2) Construa uma árvore de decisão para cada amostra e obtenha um resultado de previsão de cada árvore de decisão.
- (3) Faça um voto para cada resultado previsto.
- (4) Selecione o resultado da previsão com o maior número de votos como previsão final.

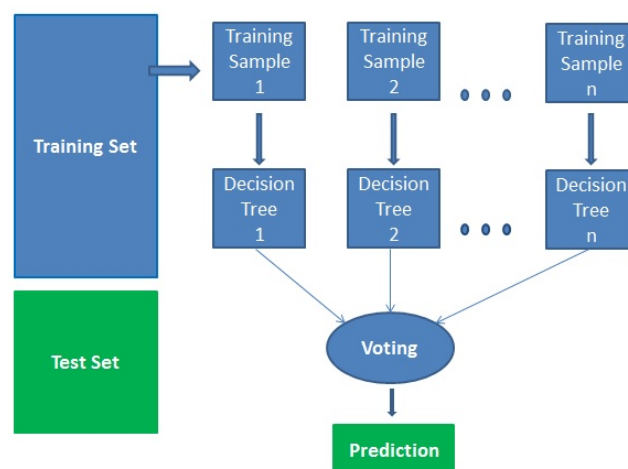


Figura 1.1: Divisão e conquista no algoritmo Random Forest[1]

1.4 Vantagens e Desvantagens

Vantagens

- O uso de muitas árvores de decisão torna a Random Forest um método preciso e robusto.
- Como podemos fazer uso de várias árvores de decisão, esse método não sofre com *overfitting* ².
- Como dissemos antes, podemos usar a Random Forest para classificação e regressão.
- É possível obter a importância relativa da característica, o que ajuda na seleção das melhores para treinar e testar o modelo.

Desvantagens

- Por usar as n -árvores de decisão sempre que se deseja fazer uma previsão e depois realizar a votação, a execução completa do algoritmo pode se tornar lenta.
- Esse modelo requer uma análise mais cuidadosa na interpretação em comparação com uma única árvore de decisão, onde apenas precisamos olhar para o caminho dessa árvore.
- A ordenação das características mais importantes pode alterar até para bases pequenas como a [Iris](#). Isso ocorre devido à natureza aleatória do algoritmo na construção das árvores de decisão.

²Overfitting é um fenômeno em que um modelo de aprendizado de máquina modela os dados de treinamento muito bem, mas não consegue ter um bom desempenho nos dados de teste.

Capítulo 2

Exemplo de Aplicação

A *SciKit Learn* disponibiliza uma implementação do algoritmo [RandomForestClassifier](#). Essa implementação foi usada em um trabalho realizado em 2018 para realizar a predição sobre a base de dados para predição do câncer de mama ([WPBC](#)) [4].



Tópicos em IC: Análise exploratória de dados.

Equipe: David Jansen, Iarah Gonçalves de Almeida, Paulo Borges

1) Introdução

A análise de dados será feita sobre o Breast Cancer Wisconsin (Prognostic) Data Set. (<http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Prognostic%29>) A manipulação dos dados para análise será feita utilizando a linguagem Python no ambiente Jupyter.

2) Informação dos Atributos

- 1) ID number
- 2) Outcome (R = recur, N = nonrecur)
- 3) Time (recurrence time if field 2 = R, disease-free time if field 2 = N)
- 4-33) Ten real-valued features are computed for each cell nucleus:
 1. radius (mean of distances from center to points on the perimeter)
 2. texture (standard deviation of gray-scale values)
 3. perimeter
 4. area
 5. smoothness (local variation in radius lengths)
 6. compactness ($\text{perimeter}^2 / \text{area} - 1.0$)
 7. concavity (severity of concave portions of the contour)
 8. concave points (number of concave portions of the contour)
 9. symmetry
 10. fractal dimension ("coastline approximation" - 1)

Several of the papers listed above contain detailed descriptions of how these features are computed.

The mean, standard error, and "worst" or largest (mean of the three largest values) of these features were computed for each image, resulting in 30 features. For instance, field 4 is Mean Radius, field 14 is Radius SE, field 24 is Worst Radius.

Values for features 4-33 are recoded with four significant digits.

- 34) Tumor size - diameter of the excised tumor in centimeters
- 35) Lymph node status - number of positive axillary lymph nodes observed at time of surgery.

3) Preparação do Ambiente Jupyter

In [29]:

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from urllib.request import urlopen

#URL onde se encontra a base de dados.
UCI_data_URL = 'https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-
-wisconsin/wpbc.data'

#Nomeação das colunas.
names = ['id_number', 'outcome', 'time', 'radius_mean',
         'texture_mean', 'perimeter_mean', 'area_mean',
         'smoothness_mean', 'compactness_mean', 'concavity_mean',
         'concave_points_mean', 'symmetry_mean',
         'fractal_dimension_mean', 'radius_se', 'texture_se',
         'perimeter_se', 'area_se', 'smoothness_se',
         'compactness_se', 'concavity_se', 'concave_points_se',
         'symmetry_se', 'fractal_dimension_se',
         'radius_worst', 'texture_worst', 'perimeter_worst',
         'area_worst', 'smoothness_worst',
         'compactness_worst', 'concavity_worst',
         'concave_points_worst', 'symmetry_worst',
         'fractal_dimension_worst', 'tumor_size', 'lymph_node_status']

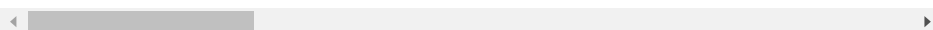
#Leitura do arquivo para o formato csv.
wpbc = pd.read_csv(urlopen(UCI_data_URL), names=names)

#Mostrando resultado em forma de tabela (10 primeiros).
wpbc.head(10)
```

Out[29]:

	id_number	outcome	time	radius_mean	texture_mean	perimeter_mean	area_me
0	119513	N	31	18.02	27.60	117.50	1013.0
1	8423	N	61	17.99	10.38	122.80	1001.0
2	842517	N	116	21.37	17.44	137.50	1373.0
3	843483	N	123	11.42	20.38	77.58	386.1
4	843584	R	27	20.29	14.34	135.10	1297.0
5	843786	R	77	12.75	15.29	84.60	502.7
6	844359	N	60	18.98	19.61	124.40	1112.0
7	844582	R	77	13.71	20.83	90.20	577.9
8	844981	N	119	13.00	21.82	87.50	519.8
9	845010	N	76	12.46	24.04	83.97	475.9

10 rows × 35 columns



4) Tratamento do Conjunto de Dados para análise

O atributo 'outcome' está representado com um objeto string. Faremos um mapeamento do atributo para valores numéricos, onde N = 0.0 e R = 1.0. Isso ajuda na plotagem de gráficos e na aplicação de alguns modelos.

Algumas informações podem ser obtidas através do arquivo wpbc.names

(<https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/wpbc.names>):

- O atributo 'lymph_node_status' está ausente em 4 amostras.
- 151 amostras não apresentam recorrência('outcome' = 0), 47 apresentam recorrência('outcome' = 1).

Dada a primeira informação na lista acima, iremos retirar essas 4 amostras de nossa análise.

As entradas da tabela serão passadas para o formato **float64**.

As colunas **id_number**, **time** serão removidas do nosso conjunto de dados por não se tratarem de informações relevantes em nossa análise.

In [30]:

```
# Mapeia a coluna 'outcome' para valores numéricos, 'R' de recur será 1, 'N' de nonrecu
r será 0
wpbc['outcome'] = wpbc['outcome'].map({'R':1.0, 'N':0.0})

# Retira os 4 valores null da coluna 'lymph_node_status'
wpbc = wpbc[wpbc['lymph_node_status'] != '?']

# Converte todo o dataframe para float64
wpbc = wpbc.astype('float64')

#Removendo colunas "id_number" e "time"
wpbc = wpbc.drop(columns=['id_number', 'time'])

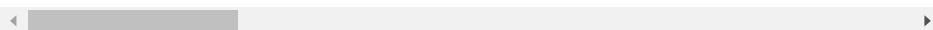
# Reseta os índices para não ter falha dos valores retirados
wpbc.reset_index(inplace=False)

wpbc.head(10)
```

Out[30]:

	outcome	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_
0	0.0	18.02	27.60	117.50	1013.0	0.09489
1	0.0	17.99	10.38	122.80	1001.0	0.11840
2	0.0	21.37	17.44	137.50	1373.0	0.08836
3	0.0	11.42	20.38	77.58	386.1	0.14250
4	1.0	20.29	14.34	135.10	1297.0	0.10030
5	1.0	12.75	15.29	84.60	502.7	0.11890
7	1.0	13.71	20.83	90.20	577.9	0.11890
8	0.0	13.00	21.82	87.50	519.8	0.12730
9	0.0	12.46	24.04	83.97	475.9	0.11860
10	0.0	16.02	23.24	102.70	797.8	0.08206

10 rows × 33 columns



5) Informações da tabela

Podemos ter mais algumas informações estatísticas do conjunto de dados, usando a função **describe**.
Informações dos tipos de dados, podem ser obtidos com a função **info**.

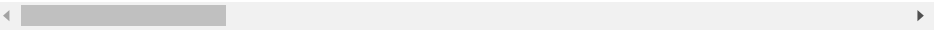
In [31]:

```
wpbc.describe()
```

Out[31]:

	outcome	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness
count	194.000000	194.000000	194.000000	194.000000	194.000000	194.000000
mean	0.237113	17.402320	22.300979	114.781495	969.092268	0.102708
std	0.426413	3.171672	4.335292	21.430694	353.159959	0.012601
min	0.000000	10.950000	10.380000	71.900000	361.600000	0.074951
25%	0.000000	15.052500	19.342500	98.160000	702.525000	0.093601
50%	0.000000	17.290000	21.795000	113.700000	929.100000	0.102200
75%	0.000000	19.580000	24.782500	129.650000	1193.500000	0.111301
max	1.000000	27.220000	39.280000	182.100000	2250.000000	0.144700

8 rows × 33 columns



In [32]:

```
wpbc.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 194 entries, 0 to 197
Data columns (total 33 columns):
outcome                194 non-null float64
radius_mean            194 non-null float64
texture_mean           194 non-null float64
perimeter_mean         194 non-null float64
area_mean              194 non-null float64
smoothness_mean        194 non-null float64
compactness_mean        194 non-null float64
concavity_mean          194 non-null float64
concave_points_mean     194 non-null float64
symmetry_mean           194 non-null float64
fractal_dimension_mean  194 non-null float64
radius_se              194 non-null float64
texture_se              194 non-null float64
perimeter_se            194 non-null float64
area_se                 194 non-null float64
smoothness_se           194 non-null float64
compactness_se           194 non-null float64
concavity_se            194 non-null float64
concave_points_se       194 non-null float64
symmetry_se             194 non-null float64
fractal_dimension_se    194 non-null float64
radius_worst            194 non-null float64
texture_worst           194 non-null float64
perimeter_worst         194 non-null float64
area_worst              194 non-null float64
smoothness_worst        194 non-null float64
compactness_worst       194 non-null float64
concavity_worst         194 non-null float64
concave_points_worst    194 non-null float64
symmetry_worst          194 non-null float64
fractal_dimension_worst  194 non-null float64
tumor_size              194 non-null float64
lymph_node_status       194 non-null float64
dtypes: float64(33)
memory usage: 51.5 KB
```

6) Exploração de Dados

Em nosso [trabalho intermediário](#)

(<https://github.com/DWalterJansen/topicosIC/blob/master/Equipe/Trabalho%20Intermedi%C3%A1rio/trabalho>) analisamos apenas as características de valor médio(_mean), tumor_size, e lymph_node_status para avaliar a influência na recorrência ou não recorrência do câncer de mama. Os resultados lá obtidos, nos motivaram a aprofundar na análise. Nessa versão final do trabalho para a disciplina, iremos aumentar o número de características observadas e aplicá-las em alguns modelos de classificação.

6.1) Floresta Aleatória com conjunto de treinamento e teste

Iniciamos usando uma **RandomForestClassifier** aplicada ao nosso dataframe **wdbc**, onde usaremos as 32 características para fazer a predição do valor na coluna **outcome**.

6.1.1) Procedimentos para classificação

Primeiramente, vamos separar nossas colunas em dependentes e independentes. A divisão entre o conjunto de treinamento e teste será feito usando o algoritmo **K-Fold** que realiza o *split* no conjunto de dados.

6.1.2) Treinamento do Modelo e Resultado do treinamento

Depois de separarmos os conjuntos, podemos aplicar o treinamento. Em seguida, verificamos a precisão do modelo através da média dos resultados obtidos para cada *split*, o que nos diz quão frequente nosso classificador estará correto.

Durante a execução do algoritmo, guardamos sempre o melhor resultado de predição. Usaremos esse resultado para medir a importância das características para a predição.

6.1.3) Importância das Variáveis:

Uma outra importante funcionalidade da **RandomForestClassifier** é a possibilidade de exibirmos a importância das características para o modelo.



In [33]:

```
#####
#6.1)

#Separa as colunas em dependentes e independentes

namesrfc = ['radius_mean',
            'texture_mean', 'perimeter_mean', 'area_mean',
            'smoothness_mean', 'compactness_mean', 'concavity_mean',
            'concave_points_mean', 'symmetry_mean',
            'fractal_dimension_mean', 'radius_se', 'texture_se',
            'perimeter_se', 'area_se', 'smoothness_se',
            'compactness_se', 'concavity_se', 'concave_points_se',
            'symmetry_se', 'fractal_dimension_se',
            'radius_worst', 'texture_worst', 'perimeter_worst',
            'area_worst', 'smoothness_worst',
            'compactness_worst', 'concavity_worst',
            'concave_points_worst', 'symmetry_worst',
            'fractal_dimension_worst', 'tumor_size', 'lymph_node_status']

X = wpbc[namesrfc] # Características
y = wpbc['outcome'] # Rótulo

# Treinamento do modelo
# Importa 'RandomForestClassifier' de 'sklearn.ensemble'
from sklearn.ensemble import RandomForestClassifier
# Importa 'metrics' de 'sklearn'
from sklearn import metrics

# Cria um classificador gaussiano
clf = RandomForestClassifier(n_estimators=100)

# Importa o K-fold para separar entre treino e teste
from sklearn.model_selection import KFold

# Valor do K (número de divisões na base de dados)
kf = KFold(n_splits = 10)

best_precision = 0.0000
sum_precisions = 0.0000

#####
# 6.1.1)
i = 0
for train_index, test_index in kf.split(X):
    i = i + 1 #contado de split
    X_train, X_test = X.iloc[train_index], X.iloc[test_index]
    y_train, y_test = y.iloc[train_index], y.iloc[test_index]

    # Treina o modelo usando o conjunto de treinamento
    clf.fit(X_train, y_train)
    y_pred = clf.predict(X_test)

    # Modelo de precisão, que diz quão frequente o classificador está correto
    precision = metrics.accuracy_score(y_test, y_pred)
    print("Precisão encontrada para %dº split: " %(i) ,end="")
    print(precision)
    sum_precisions = sum_precisions + precision
    if best_precision < precision:
```

```

        best_precision = precision
        clf_best = clf
        y_test_best = y_test
        y_pred_best = y_pred

print("Resultados Atuais..")
print("\nPrecisão média para as %d predições: " %i, end="")
print(sum_precisions/i);
print("Melhor resultado: ", end="")
print(best_precision)

#####
# 6.1.3

# Ordenação das as características importantes pela pontuação usando a iteração com a melhor precisão
print("\nCaracterísticas ordenadas pela pontuação:")
# O número dentro de round é o limite de casas decimais
result = sorted(zip(map(lambda x: round(x, 6), clf_best.feature_importances_), namesrfc), reverse=True)

# Criando dataframe com o resultado
df = pd.DataFrame(data=result)
df = df.rename(index=str, columns={0: "importance", 1: "feature"})

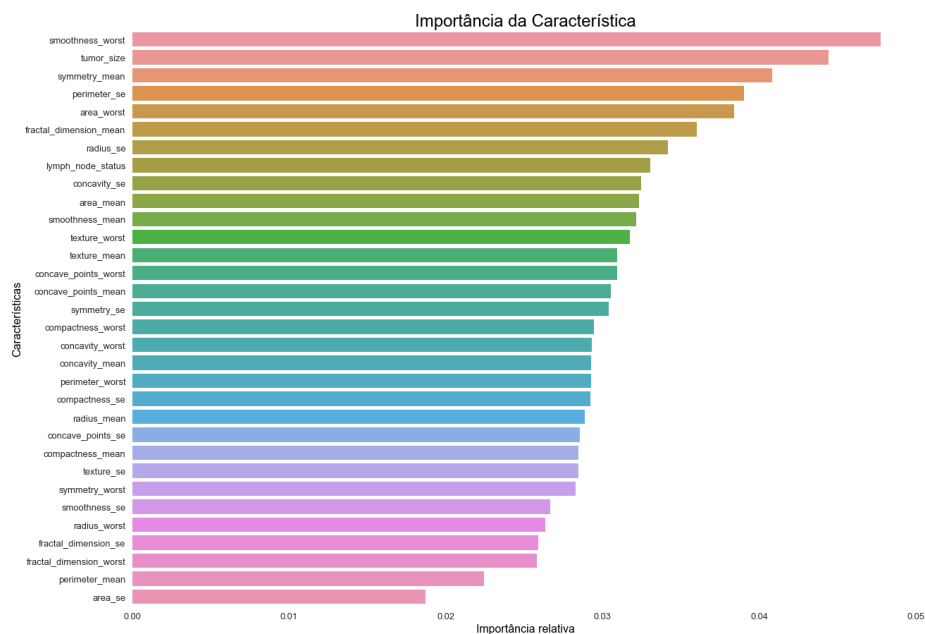
import matplotlib.pyplot as plt
#Configurações para exibição do gráfico
sns.set_color_codes("dark")
sns.set_style("white")
sns.set_context("talk")
plt.figure(figsize=(20,15))
g = sns.barplot(x="importance", y="feature", data=df)
g.axes.set_title('Importância da Característica', fontsize=24,color="black",alpha=2)
g.set_xlabel("Importância relativa", size = 16,color="black")
g.set_ylabel("Características", size = 16,color="black")
sns.despine(left=True, bottom=True)

```


Precisão encontrada para 1º split: 0.7
 Precisão encontrada para 2º split: 0.8
 Precisão encontrada para 3º split: 0.65
 Precisão encontrada para 4º split: 0.75
 Precisão encontrada para 5º split: 0.9473684210526315
 Precisão encontrada para 6º split: 0.631578947368421
 Precisão encontrada para 7º split: 0.631578947368421
 Precisão encontrada para 8º split: 0.631578947368421
 Precisão encontrada para 9º split: 0.7894736842105263
 Precisão encontrada para 10º split: 1.0
 Resultados Atuais..

Precisão média para as 10 predições: 0.7531578947368421
 Melhor resultado: 1.0

Características ordenadas pela pontuação:



6.2) Correlação

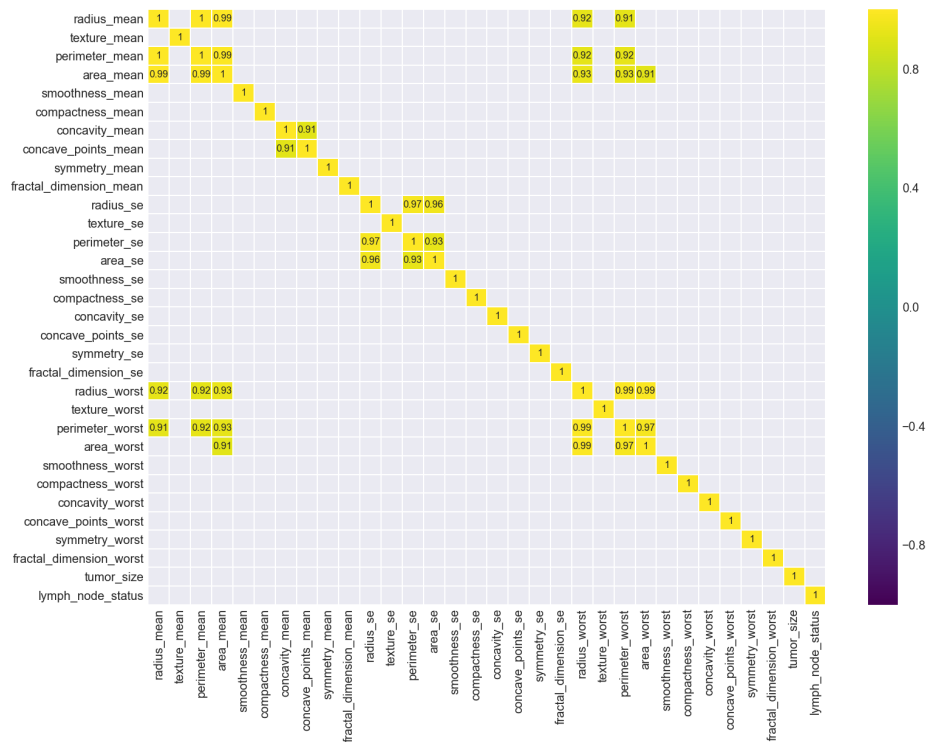
Nossa base de dados possui muitas características, portando iremos tentar reduzir essa quantidade olhando para os atributos correlacionados.

Uma ferramenta que nos permite aprofundar um pouco mais na exploração da correlação de características é mapa de calor:

```
# Gerando a correlação das características
corr = wpbc.drop('outcome', axis=1).corr() # We already examined SalePrice correlations
plt.figure(figsize=(25, 18))
sns.set_context("poster")
sns.set(font_scale=2.0)

# Exibindo as correlações maiores que 0.75 e menor que -0.75
c = sns.heatmap(corr[(corr >= 0.9) | (corr <= -0.9)],
                cmap='viridis', vmax=1.0, vmin=-1.0, linewidths=1,
                annot=True, annot_kws={"size": 16}, square=False);

size = len(namesrfc)
cor_features = list();
data = (corr[(corr >= 0.9) | (corr <= -0.9)])
for i in range(0,size):
    for j in range(i):
        if (data.iloc[i,j] > 0):
            cor_features.append([namesrfc[i], namesrfc[j]])
```



Referências Bibliográficas

- [1] Understanding random forests classifiers in python. <https://www.datacamp.com/community/tutorials/random-forests-classifier-python#algorithm>. Acessado: 2019-05-25.
- [2] Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, Oct 2001.
- [3] Tin Kam Ho. Random decision forests. In *Proceedings of the Third International Conference on Document Analysis and Recognition (Volume 1) - Volume 1*, ICDAR '95, pages 278–, Washington, DC, USA, 1995. IEEE Computer Society.
- [4] A. Iarah Gonçalves J. David Walter, B. M. Paulo Augusto. *Random Forest para predição sobre a Breast Cancer Wisconsin (Prognostic) Data Set*. 2018.
- [5] Tin Kam Ho. The random subspace method for constructing decision forests. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(8):832–844, Aug 1998.