

# Álgebra Linear aplicada à análise de crédito

O principal tópico abordado nesse trabalho é a utilização de regressão logística para classificação de credores, classificando-os como bons ou mals.

In [1]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import seaborn as sns
5 %matplotlib inline
6
7 from collections import OrderedDict
8
9 # Bibliotecas de treinamento
10 from sklearn.model_selection import train_test_split
11 from sklearn.metrics import accuracy_score, precision_score, recall_score, confusion_ma
12 from sklearn.preprocessing import RobustScaler
13
14 # PCA
15 from sklearn.decomposition import PCA
```

In [2]:

```
1 default = pd.read_csv("UCI_Credit_Card.csv", index_col="ID")
2 pd.set_option('display.max_columns', None)
3 default['EDUCATION'] = np.where(default['EDUCATION'] == 5, 4, default['EDUCATION'])
4 default['EDUCATION'] = np.where(default['EDUCATION'] == 6, 4, default['EDUCATION'])
5 default['EDUCATION'] = np.where(default['EDUCATION'] == 0, 4, default['EDUCATION'])
6 default['MARRIAGE'] = np.where(default['MARRIAGE'] == 0, 3, default['MARRIAGE'])
7 default.rename(columns = {'default.payment.next.month': "default"}, inplace = True)
8 default.rename(columns = lambda x: x.lower(), inplace = True)
9 default.head()
10
```

Out[2]:

|    | limit_bal | sex | education | marriage | age | pay_0 | pay_2 | pay_3 | pay_4 | pay_5 | pay_6 | bill_ |
|----|-----------|-----|-----------|----------|-----|-------|-------|-------|-------|-------|-------|-------|
| ID |           |     |           |          |     |       |       |       |       |       |       |       |
| 1  | 20000.0   | 2   | 2         | 1        | 24  | 2     | 2     | -1    | -1    | -2    | -2    | 38    |
| 2  | 120000.0  | 2   | 2         | 2        | 26  | -1    | 2     | 0     | 0     | 0     | 2     | 21    |
| 3  | 90000.0   | 2   | 2         | 2        | 34  | 0     | 0     | 0     | 0     | 0     | 0     | 29    |
| 4  | 50000.0   | 2   | 2         | 1        | 37  | 0     | 0     | 0     | 0     | 0     | 0     | 46    |
| 5  | 50000.0   | 1   | 2         | 1        | 57  | -1    | 0     | -1    | 0     | 0     | 0     | 81    |

## Descrição dos Dados(Features)

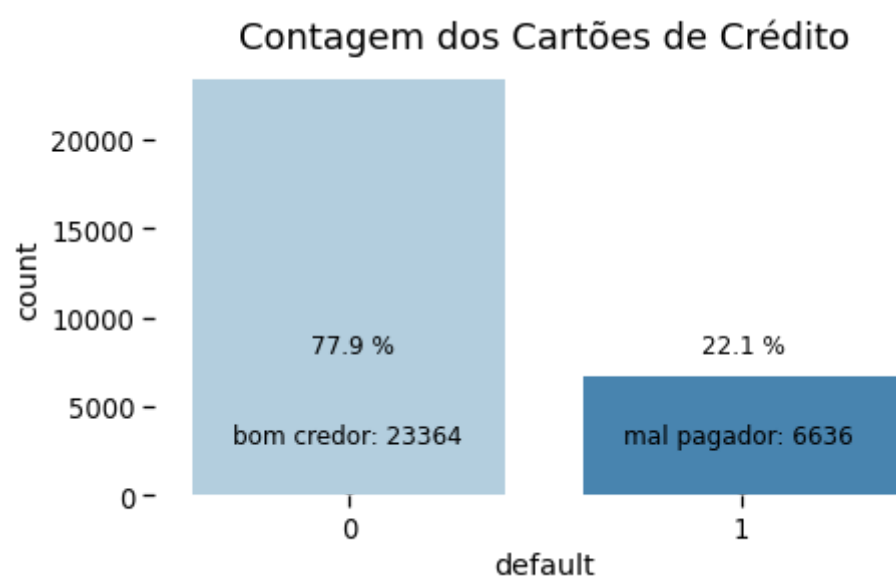
- **limit\_bal**: Limite de crédito considerando o limite individual e familiar.
- **sex**: Gênero

- 1 = Homem
- 2 = Mulher
- **education:** Nível de escolaridade
  - 1 = Pós-Graduação
  - 2 = Universidade
  - 3 = Ensino Médio
  - 4 = Outros
- **marriage:** Estado Civil
  - 1 = Casado(a)
  - 2 = Solteiro(a)
  - 3 = Outros
- **age:** Idade(Anos)
- **pay\_0** até **pay\_6:** Histórico de pagamentos entre os meses de Abril e Setembro de 2005.
  - Números Negativos = Pagou Adiantado
  - 0 = Pagou na data correta
  - 1 = Com um mês de atraso
  - 2 = Com dois meses de atraso, ...
- **bill\_amt1** até **bill\_amt6:** Dívida acumulada nos meses de Abril e Setembro de 2005.
- **pay\_amt1** até **pay\_amt6:** Montante pago em antecipado
- **default:**
  - 1 = Não pagou
  - 0 = Pagou

## Entendendo melhor os dados

In [3]:

```
1 sim = default.default.sum()
2 nao = len(default)-sim
3
4 sim_perc = round(sim/len(default)*100, 1)
5 nao_perc = round(nao/len(default)*100, 1)
6 import sys
7 plt.figure(figsize=(7,4))
8 sns.set_context('notebook', font_scale=1.2)
9 sns.countplot('default',data=default, palette="Blues")
10 plt.annotate('bom credor: {}'.format(nao), xy=(-0.3, 15000), xytext=(-0.3, 3000), size=12)
11 plt.annotate('mal pagador: {}'.format(sim), xy=(0.7, 15000), xytext=(0.7, 3000), size=12)
12 plt.annotate(str(nao_perc)+" %", xy=(-0.3, 15000), xytext=(-0.1, 8000), size=12)
13 plt.annotate(str(sim_perc)+" %", xy=(0.7, 15000), xytext=(0.9, 8000), size=12)
14 plt.title('Contagem dos Cartões de Crédito', size=18)
15
16 plt.box(False);
```



Algumas renderizações dos dados.

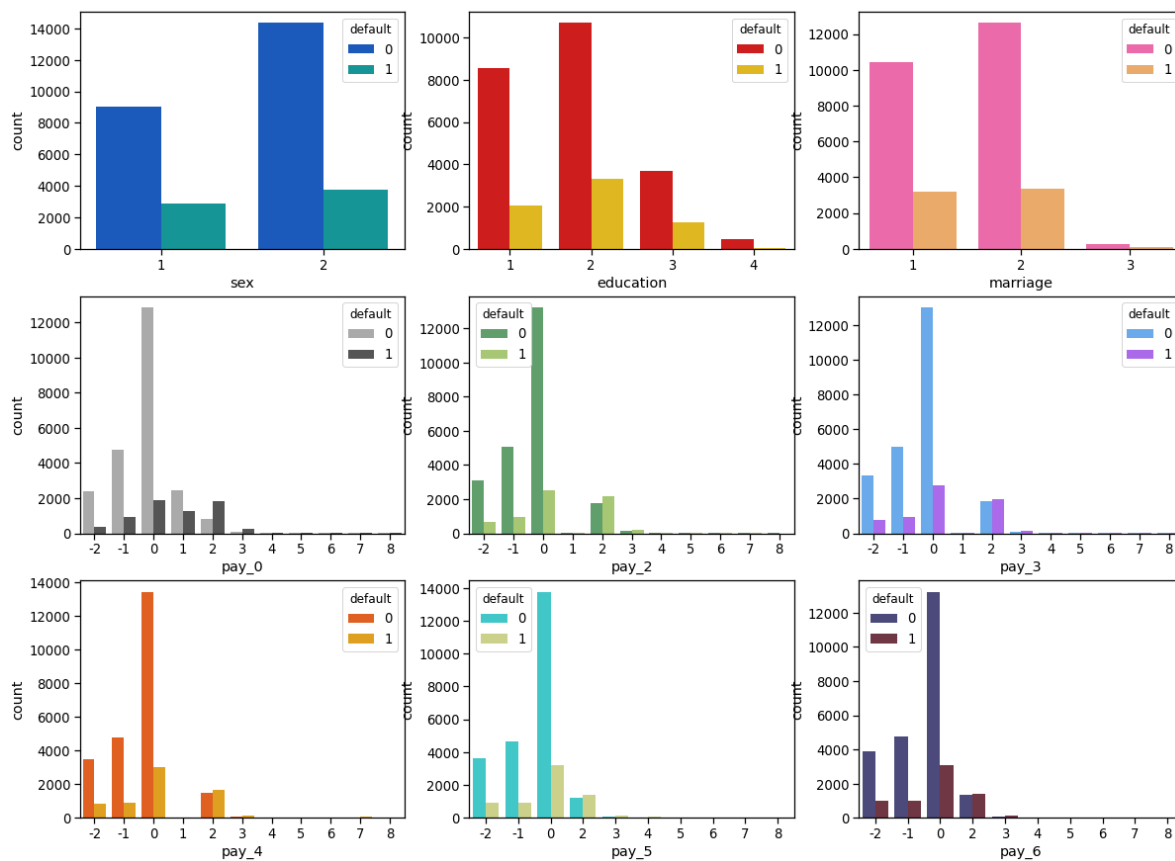
In [4]:

```

1 # Creating a new dataframe with categorical variables
2 subset = default[['sex', 'education', 'marriage', 'pay_0', 'pay_2', 'pay_3', 'pay_4', 'pay_5', 'pay_6']]
3
4
5 f, axes = plt.subplots(3, 3, figsize=(20, 15), facecolor='white')
6 f.suptitle(' Frequência das Variáveis Categóricas')
7 ax1 = sns.countplot(x="sex", hue="default", data=subset, palette="winter", ax=axes[0,0])
8 ax2 = sns.countplot(x="education", hue="default", data=subset, palette="hot", ax=axes[0,1])
9 ax3 = sns.countplot(x="marriage", hue="default", data=subset, palette="spring", ax=axes[0,2])
10 ax4 = sns.countplot(x="pay_0", hue="default", data=subset, palette="binary", ax=axes[1,0])
11 ax5 = sns.countplot(x="pay_2", hue="default", data=subset, palette="summer", ax=axes[1,1])
12 ax6 = sns.countplot(x="pay_3", hue="default", data=subset, palette="cool", ax=axes[1,2])
13 ax7 = sns.countplot(x="pay_4", hue="default", data=subset, palette="autumn", ax=axes[2,0])
14 ax8 = sns.countplot(x="pay_5", hue="default", data=subset, palette="rainbow", ax=axes[2,1])
15 ax9 = sns.countplot(x="pay_6", hue="default", data=subset, palette="icefire", ax=axes[2,2])

```

Frequência das Variáveis Categóricas



A partir da análise visual, é notável que mulheres(2), graduados(2), solteiros(2) e todos os que costumam pagar na data prevista tendem a ser bons credores.

## Preparação dos dados

Antes de realizar qualquer operação com os dados, é preciso que as informações estejam formatadas adequadamente. Diante disso, faz-se necessário dispor os dados qualitativos no formato binário, fazendo alterações semelhantes ao do exemplo abaixo:

Tabela exemplo:

| Sex:   |
|--------|
| Male   |
| Female |

Tabela exemplo no formato binário:

| id | Male | Female |
|----|------|--------|
| 1  | 1    | 0      |
| 2  | 1    | 0      |
| 3  | 0    | 1      |

Os dados da segunda tabela, nessa formatação, são denominados dummy features

In [5]:

```
1 #Formatando os dados
2 default['grad_school'] = (default['education']==1).astype('int')
3 default['university'] = (default['education']==2).astype('int')
4 default['high_school'] = (default['education']==3).astype('int')
5 default.drop('education', axis=1, inplace=True)
6
7 default['male'] = (default['sex']==1).astype('int')
8 default.drop('sex', axis=1, inplace=True)
9
10 default['married'] = (default['marriage'] == 1).astype('int')
11 default.drop('marriage', axis=1, inplace=True)
12 default.sample(2)
```

Out[5]:

|       | limit_bal | age | pay_0 | pay_2 | pay_3 | pay_4 | pay_5 | pay_6 | bill_amt1 | bill_amt2 | bill_ar |
|-------|-----------|-----|-------|-------|-------|-------|-------|-------|-----------|-----------|---------|
| ID    |           |     |       |       |       |       |       |       |           |           |         |
| 6972  | 20000.0   | 37  | 0     | 0     | -2    | -1    | 0     | 0     | 15960.0   | 0.0       | -1953   |
| 14075 | 90000.0   | 28  | 1     | 2     | 2     | 2     | 2     | 0     | 56889.0   | 55407.0   | 64428   |

Os indivíduos que antecipam os pagamentos ou pagam na data prevista são bons credores. Diante disso, é útil resumir-los a variável 0 nas colunas pay\_0 a pay\_6, uma vez que, a predição ficará menos complexa.

In [6]:

```
1 pay_features = ['pay_0', 'pay_2', 'pay_3', 'pay_4', 'pay_5', 'pay_6']
2
3 #Substituindo todos os valores negativos por 0
4 for p in pay_features:
5     default.loc[default[p]<=0,p] = 0
6
```

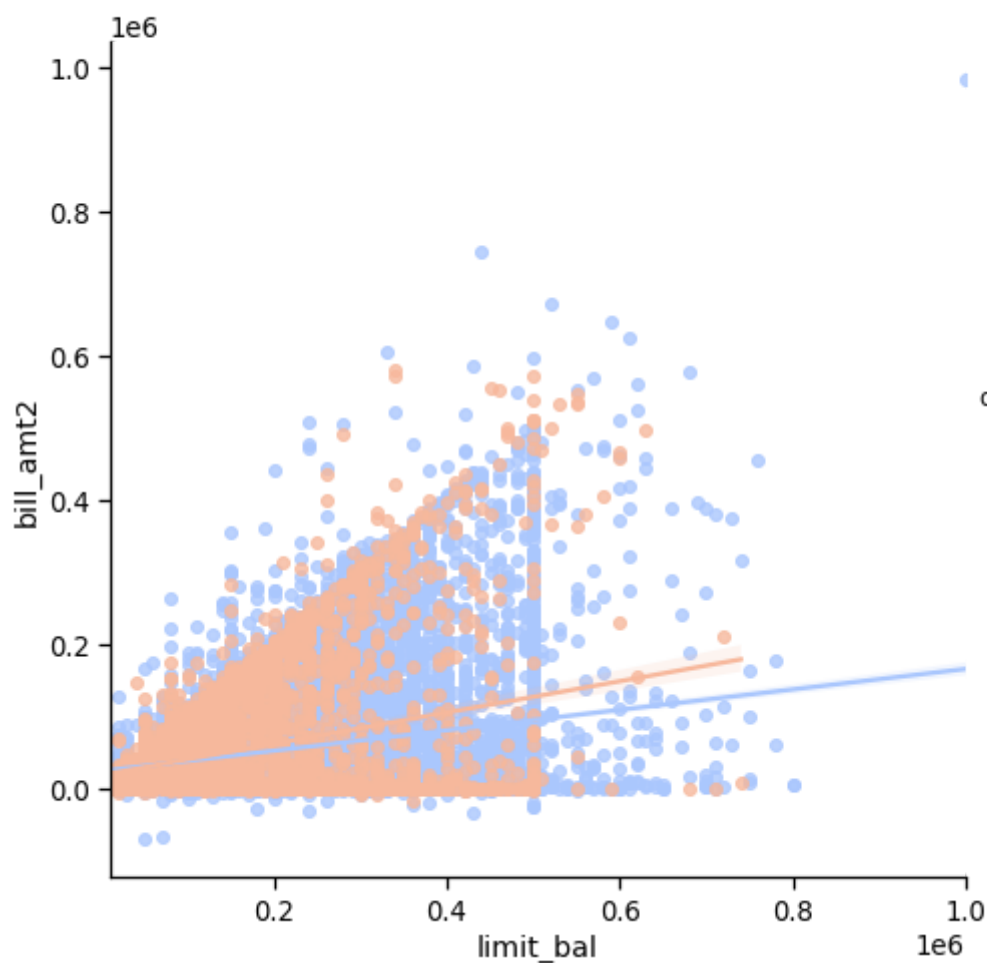
## Regressão Linear

Vamos observar um exemplo com Regressão Linear e o porquê ela não é ideal para esse tipo de análise.

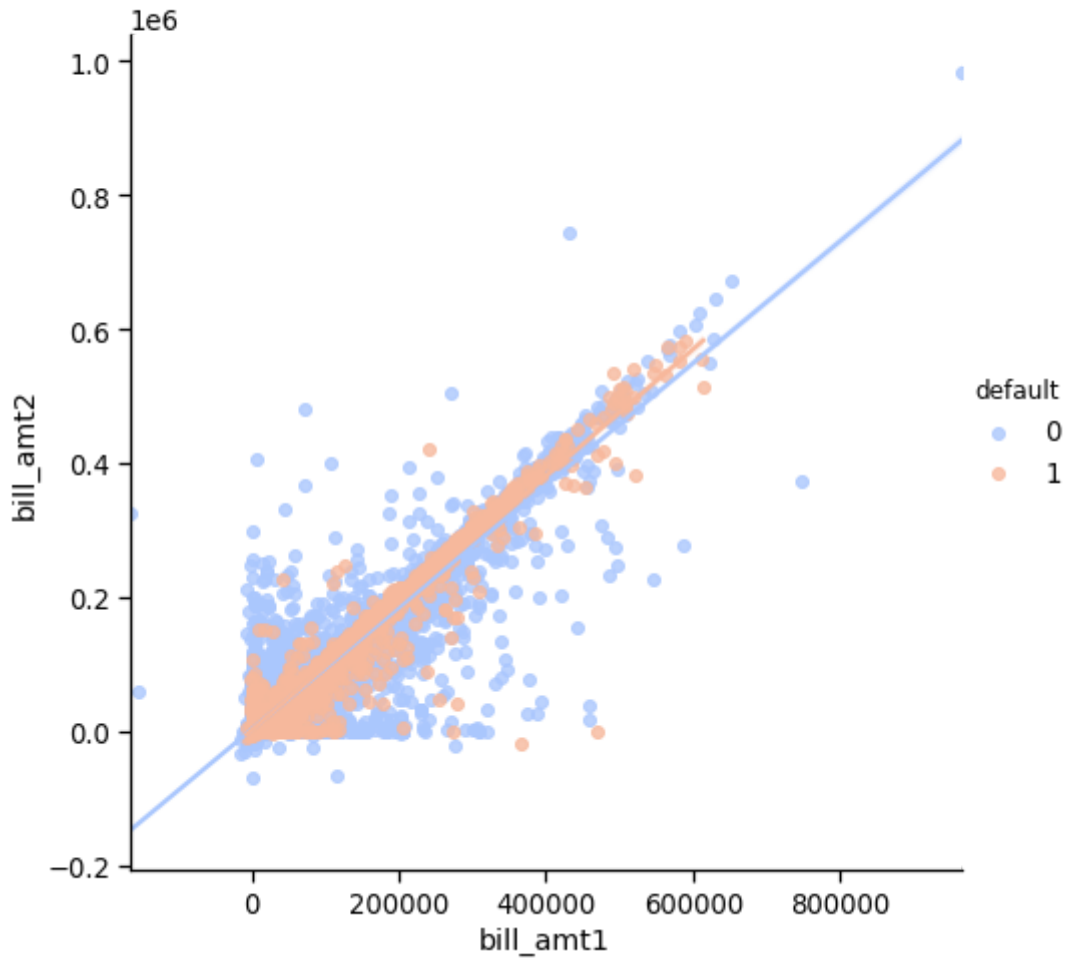
In [7]:

```
1 sns.lmplot(x='limit_bal', y='bill_amt2', data = default, hue = 'default',
2           palette='coolwarm', height=7)
3 plt.title('Regressão Linear: Distinguindo pagadores e não pagadores\n\n', size=18)
4
5 sns.lmplot(x='bill_amt1', y='bill_amt2', data = default, hue = 'default',
6           palette='coolwarm', height=7)
7 plt.title('Regressão Linear:\n Não pode distinguir entre pagadores e não pagadores\n',
8
```

## Regressão Linear: Distinguindo pagadores e não pagadores



## Regressão Linear: Não pode distinguir entre pagadores e não pagadores



Regressão Linear prediz com mais precisão modelos em que a saída é um valor contínuo, enquanto que, Regressão Logística modela melhor casos em que é necessário uma resposta binária como True or False.

## Regressão Logística

A Regressão Logística é uma técnica estatística cujo o objetivo é prever a ocorrência de um evento dado um conjunto de dados.



Nesse notebook, a Regressão Logística será usada para avaliar se um indivíduo é um bom credor, baseado em seus dados pessoais.

## Criando o Primeiro Modelo usando todos os dados

In [8]:

```
1 #Bibliotecas úteis do sklearn
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import accuracy_score, precision_score, recall_score, confusion_ma
4 from sklearn.preprocessing import RobustScaler
```

Como a Regressão Logística é um modelo de Machine Learning supervisionado, é preciso separar do DataFrame as labels e os dados. Dessa forma, o modelo modificará os parâmetros de uma função para que as informações inseridas, de cada usuário, nessa função, resultem nos valores 0 ou 1. De modo que, esses valores binários corresponderão a classificação 0 (bom credor) ou 1 (mal pagador).

$$f : R^n \rightarrow k \in [0, 1]$$

In [9]:

```
1 #Armazenando as labels
2 y = default.copy()['default']
3 #Armazenando os dados
4 X = default.copy().drop('default', axis=1)
```

### Normalizando os dados

Em tese, Regressão Logística não precisa de normalização de Dados, pois esta trabalha com probabilidades(veremos sobre isso adiante), porém a existência de outliers podem prejudicar a performance do modelo. Nesse sentido, a normalização faz-se necessária antes de implementar o modelo e para este fim, a função `RobustScaler`, do `SKLEARN`, é ideal, pois remove os valores que não estão entre o 1º e 2º quartil e normaliza os dados.

In [10]:

```
1 #Removendo os outliers
2 robust_scaler = RobustScaler()
3 X = robust_scaler.fit_transform(X)
4 X = pd.DataFrame(X, columns = default.drop("default", axis = 1).columns)
```

Em seguida, é necessário extrair alguns dados e labels de test, para que ao fim da modelagem sejam usados na avaliação da performance do modelo.

In [11]:

```
1 #Separando os dados de treino e test
2 X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.15, random_state =
```

Após o tratamento dos dados, aparentemente, já é viável executar o modelo. Mas antes, vale entender a Matemática que compõe o modelo de Regressão Logística.

## A Matemática da Regressão Logística

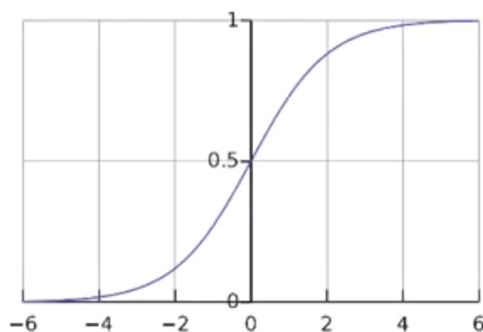
Como dito anteriormente, este é um método de aprendizado supervisionado, ou seja, o modelo realiza a tarefa de aprender uma função que mapeia uma entrada para obter uma saída, com base em pares de entrada-saída.

O termo Logística deriva do termo log odds, no qual, odds refere-se a razão entre a probabilidade de um evento acontecer e a probabilidade de não acontecer. O modelo recebeu este nome devido a utilização de log odds no cálculo da Regressão.

$$Odds = \frac{P(event)}{1 - P(event)}$$

O gráfico da Regressão Logística é gerado a partir de uma função Sigmoide e uma das razões para usar essa curva é por se adequar melhor a modelos de classificação binária, uma vez que, os valores de saída podem ser interpretados como a probabilidade de um evento acontecer, já que estes variam entre 0 e 1.

$$p(X) = \frac{1}{1 + e^{-BX}}; p(x) \in [0, 1]; X, B \in R^n$$



A Regressão Logística estimará o parâmetro B de modo que função Sigmoide atinja a maior similaridade.

### Estimando o parâmetro

- Suponha que exista  $N$  amostras com labels 0 ou 1.
- Algumas amostras correspondem as Labels "1".
- Já outras amostras correspondem as Labels "0"
- O objetivo é estimar um B tal que,  $p(x)$  esteja o mais próximo de 1 nas amostras "1" e  $(1-p(x))$  o mais próximo de 1 nas amostras "0"

$$\prod_{s \text{ em } y_i=1} p(x_i) = p(x_1) * p(x_2) * \dots * p(x_n)$$

$$\prod_{s \text{ em } y_i=0} 1 - p(x_i) = (1 - p(x_1)) * (1 - p(x_2)) * \dots * (1 - p(x_n))$$

$$L(B) = \prod_{s \text{ em } y_i=1} p(x_i) * \prod_{s \text{ em } y_i=0} (1 - p(x_i))$$

$$L(B) = \prod_s p(x_i)^{y_i} * (1 - p(x_i))^{1-y_i}$$

Aplicando log odds para isolar o parâmetro  $BX$  e transformar o produtório em somatório

$$\begin{aligned}
\log(L(B)) &= \log\left(\prod_s p(x_i)^{y_i} * (1 - p(x_i))^{1-y_i}\right) = l(B) \\
l(B) &= \sum_{i=1} y_i \log\left(\frac{1}{1 + e^{-Bx_i}}\right) + (1 - y_i) \log(1 - p(x_i)) \\
l(B) &= \sum_{i=1} y_i Bx_i + \log\left(\frac{1}{1 + e^{Bx_i}}\right) \\
l(B) &= \sum_{i=1} Bx_i - \log 1 + e^{Bx_i} \\
B &= \arg \max_B l(B)
\end{aligned}$$

\$\$

Como essa equação é transcendental, não existe um B exato, entretanto, o Método de Newton oferece uma boa aproximação para este valor.

$$B^{t+1} = B^t - \frac{\nabla_B l(B^t)}{\nabla_{BB} l(B^t)}$$

Calculando o gradiente  $\nabla_B l$

$$\begin{aligned}
\nabla_B l &= \sum_{i=1} \nabla_B [y_i Bx_i] - \nabla_B [\log(1 + e^{Bx_i})] \\
\nabla_B l &= \sum_{i=1} y_i x_i - \left[ \frac{1}{1 + e^{Bx_i}} e^{Bx_i} x_i \right] \\
\nabla_B l &= \sum_{i=1} y_i x_i - \left[ \frac{1}{1 + e^{-Bx_i}} x_i \right] \\
\nabla_B l &= \sum_{i=1} [y_i - p(x_i)] x_i
\end{aligned}$$

Convertendo para a representação matricial

$$\nabla_B = X^T (Y - \hat{Y})$$

Calculando o gradiente  $\nabla_{BB}$ , que é uma matrix de Hessian

$$\begin{aligned}
\nabla_{BB} l &= \sum_{i=1} \nabla_B [y_i - p(x_i)] x_i \\
\nabla_{BB} l &= \sum_{i=1} \nabla_B - p(x_i) x_i \\
\nabla_{BB} l &= \sum_{i=1} \nabla_B - \left[ \frac{1}{1 + e^{-Bx_i}} \right] x_i \\
\nabla_{BB} l &= \sum_{i=1} \left[ \frac{1}{1 + e^{-Bx_i}} \right]^2 e^{-Bx_i} (-x_i) x_i \\
\nabla_{BB} l &= - \sum_{i=1} \left[ \frac{e^{-Bx_i}}{1 + e^{-Bx_i}} \right] \left[ \frac{1}{1 + e^{-Bx_i}} \right] x_i^T x_i \\
\nabla_{BB} l &= - \sum_{i=1} p(x_i) (1 - p(x_i)) x_i^T x_i
\end{aligned}$$

Convertendo para forma matricial

$$\nabla_{BB}l = -X^T P(1 - P)X$$

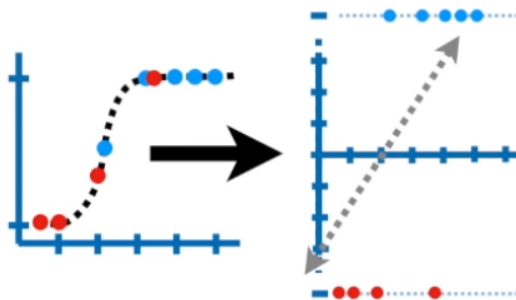
$$\nabla_{BB}l = -X^T W X$$

Desse modo, iterando a fórmula abaixo obtemos o parâmetro B.

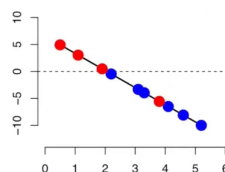
$$B^{t+1} = B^t + (X^T W X)^{-1} X^T (Y - \hat{Y}^t)$$

## Perspectiva Geométrica da Regressão Logística

Quando é aplicado "log odds" na Função Sigmoid, obtém-se a reta  $BX \in \mathbb{R}^n$



Ao calcular o vetor B com Método de Newton, geometricamente, a reta  $BX$  rotacionará até que os coeficientes de B convirjam em um vetor que torna a função Sigmoide o mais similar possível.



## Implementando o modelo

In [12]:

```
1 # 1. Importando o modelo
2 from sklearn.linear_model import LogisticRegression
3
4 #2. Criando uma instância para estimar
5 logistic_regression = LogisticRegression(n_jobs=-1, random_state=15)
6
7 #3. Usando os dados de treino para treinar o modelo
8 logistic_regression.fit(X_train,y_train);
```

## Avaliando o Modelo

In [13]:

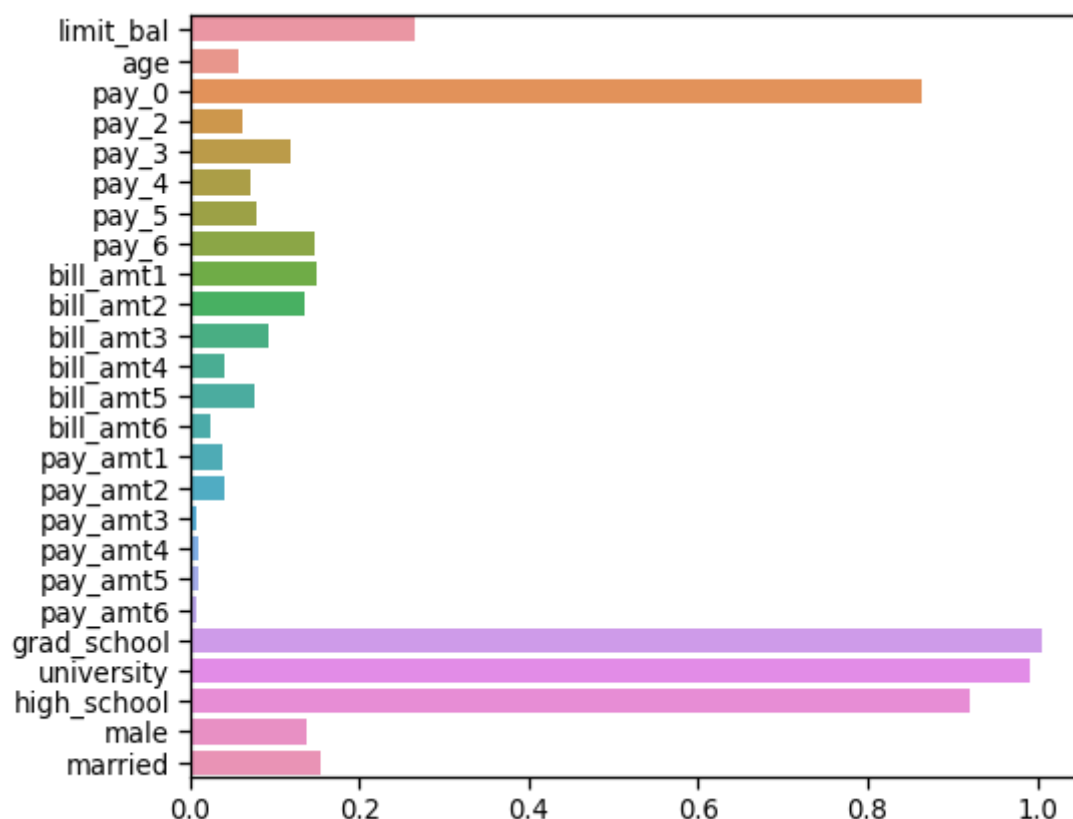
```
1 def CMatrix(y_pred_test,y_true):
2     from pandas import DataFrame
3     CM = confusion_matrix(y_pred_test, y_test)
4     df = pd.DataFrame(
5         data={'Positive(Predicted)':CM[0,:],'Negative(Predicted)':CM[1,:]},
6         index=['Positive','Negative'])
7     return df
8
9 def model_metrics(X_test,y_test, logistic_regression):
10     metrics = pd.DataFrame(index=['accuracy', 'precision', 'recall'],columns=['LogisticReg'],
11                               data=[accuracy_score(y_pred_test,y_true_test),
12                                     precision_score(y_pred_test,y_true_test),
13                                     recall_score(y_pred_test,y_true_test)])
14     print()
15     metrics.loc['accuracy', 'LogisticReg'] = accuracy_score(y_pred=y_pred_test, y_true=y_true_test)
16     metrics.loc['precision', 'LogisticReg'] = precision_score(y_pred=y_pred_test, y_true=y_true_test)
17     metrics.loc['recall', 'LogisticReg'] = recall_score(y_pred=y_pred_test, y_true=y_true_test)
18     metrics['metric_type'] = metrics.index
19     #print(metrics)
20     #print(CMatrix(y_pred_test=y_pred_test, y_true=y_test))
21     return metrics
22
23 metrics_model_1=model_metrics(X_test,y_test,logistic_regression)
```

## Analizando os coeficientes da função

O vetor parâmetro B possui um coeficiente para cada feature e os valores absolutos desses coeficientes representam a importância de cada feature no modelo de predição. Abaixo, há um gráfico com todos os valores de B atrelado as features correspondentes.

In [14]:

```
1 logistic_regression.coef_[0]
2 default.columns
3 f,ax = plt.subplots(figsize=(8, 7))
4 coef_plot = sns.barplot(y=default.columns.drop('default'),x=abs(logistic_regression.coef
```



Ao observar o gráfico, fica visível que os dados de pay\_amt1 a pay\_amt6 oferecem pouca informação para o modelo.

## Analizando a correlação entre as features

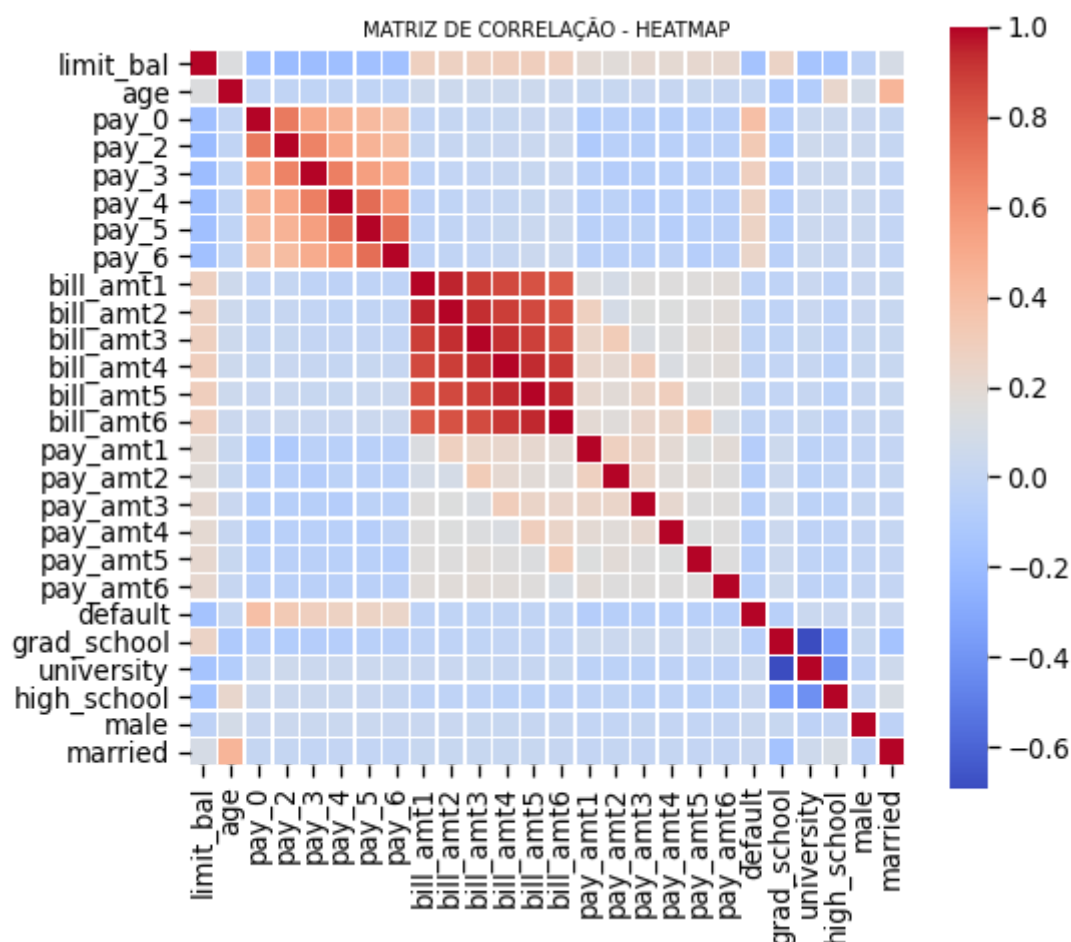
Outro fator importante a ser analisado é a correlação entre os dados, posto que, quando há grande correlação, os dados não oferecem tanta informação para o modelo. Diante disso, é válido observar o heatmap que demonstra a correlação entre as features utilizadas.

In [15]:

```
1 corr = default.corr()
2 f,ax = plt.subplots(figsize=(8, 7))
3 sns.heatmap(corr, cbar = True, square = True, annot = False, fmt= '.1f',
4             xticklabels= True, yticklabels= True
5             ,cmap="coolwarm", linewidths=.5, ax=ax)
6
7 plt.title('MATRIZ DE CORRELAÇÃO - HEATMAP', size=10)
```

Out[15]:

Text(0.5, 1.0, 'MATRIZ DE CORRELAÇÃO - HEATMAP')



Com esse plot, fica visível a correlação entre os bill\_amt. Dessa forma, é possível concluir que a bill\_amt de

um mês é suficiente para o modelo, tendo em vista que as bill\_amts de outros meses apenas repetem a mesma informação.

Após essa análise, conclui-se que é recomendável remover do modelo pay\_amt1, pay\_amt2,..., pay\_amt6 e bill\_amt1,bill\_amt2,...,bill\_amt5, visto que, acrescentam pouca informação ao modelo. Desse modo, pensando na implementação em larga escala, a remoção dessas 11 features pode diminuir bastante os custos de processamento, armazenamento e coleta de dados de uma empresa.

## Re-implementando o modelo

Diante das alterações do conjunto, faz-se necessário re-implementar o modelo e avaliar as alterações nas métricas de predição.

In [16]:

```
1  #Removendo as features
2  default2 = default.copy().drop(['pay_amt1', 'pay_amt2', 'pay_amt3', 'pay_amt4', 'pay_amt5',
3  y_2 = default2['default']
4  X_2 = default2.copy().drop('default',axis=1)
5
6  #Removendo os outliers e normalizando
7  robust_scaler_2 = RobustScaler()
8  X_2 = robust_scaler_2.fit_transform(X_2)
9
10 #Separando os Dados
11 X_train_2, X_test_2, y_train_2, y_test_2 = train_test_split(X_2,y_2, test_size=0.15, ra
12
13 #Reimplementado o modelo
14 logistic_regression_2 = LogisticRegression(n_jobs=-1, random_state=15)
15
16 logistic_regression_2.fit(X_train_2,y_train_2)
17
18 metrics_model_2 = model_metrics(X_test_2,y_test_2,logistic_regression_2)
19
```

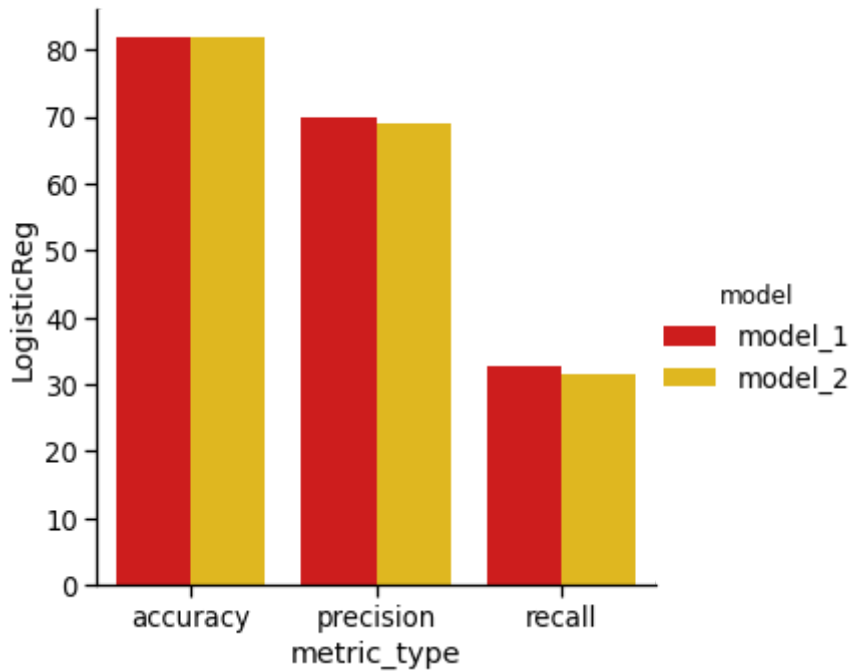
## Comparando resultados

Com os resultados da nova implementação, é possível concluir que, realmente, os dados removidos eram irrelevantes. O gráfico demonstra a precisão, acurácia e recall das duas implementações



In [17]:

```
1 metrics_model_2['model'] = 'model_2'
2 metrics_model_1['model'] = 'model_1'
3 g = sns.catplot(data=pd.concat([metrics_model_1,metrics_model_2]), kind='bar', y='Logis
```



## Utilizando PCA

O Principal Component Analysys (PCA), é uma operação linear que utiliza Singular Value Decomposition(SVD) para reduzir um conjunto de features possivelmente correlacionadas num conjunto de valores de variáveis linearmente não correlacionadas chamadas de componentes principais.

### Matemática do PCA(exemplo explicativo)

Seja o conjunto de dados representado por  $X$ . No cálculo do PCA, primeiramente, obtem-se a média dos  $n$  valores de cada coluna de  $X$ .

$$\frac{\sum_j X_{ji}}{n}$$

Agora, é necessário determinar a matriz de coavariância, para tal, considere uma matriz auxiliar  $M$  obtida da seguinte forma:

$$M = X - \frac{\sum_j X_{ji}}{n}$$

Logo, a matriz de covariância  $C$  é:

$$C = \frac{M^T M}{n - 1}$$

Diante disso, é preciso calcular os auto-valores e auto-vetores da matriz  $C$ . Os auto-valores  $\lambda$  de  $C$  são raízes do seguinte polinômio característico:

$$p_C(Y) = \det(C - YI)$$

E os auto-vetores  $q$  pertencem ao conjunto:

$$(C - \lambda I)q = 0$$

Faz-se necessário calcular o quanto cada auto-vetor(componentes) representa na taxa de variação dos dados, para isso vamos analisar a variância de cada auto-valor:

$$var(\lambda) = \frac{\lambda}{\sum_j \lambda_j}$$

Logo, com o conjunto das variâncias dos auto-valores pode-se determinar as componentes principais do dado conjunto. A importância de cada componente é dada pelo valor da variância do auto-valor correspondente.

## Reduzindo o conjunto de dados às 13 principais componentes

In [18]:

```
1 from sklearn.decomposition import PCA
2 #Definindo a dimensão dos dados
3 pca = PCA(13)
4 #PCA no conjunto de treino
5 pca_X = pca.fit(X)
6
7 X_pca = pd.DataFrame(pca_X.transform(X))
8 X_pca.head()
```

Out[18]:

|   | 0         | 1        | 2         | 3        | 4        | 5         | 6         | 7         |          |
|---|-----------|----------|-----------|----------|----------|-----------|-----------|-----------|----------|
| 0 | -3.339702 | 1.209778 | -0.426773 | 0.309149 | 0.165696 | -0.348627 | 1.198143  | 0.979059  | -0.52719 |
| 1 | -2.915311 | 0.797967 | -0.547978 | 0.430450 | 0.294873 | -0.554742 | 1.188430  | 0.905603  | -0.86573 |
| 2 | -2.246350 | 0.234234 | -0.976335 | 0.378132 | 0.169322 | -0.271923 | 0.824479  | -0.652101 | -0.42399 |
| 3 | -2.267245 | 0.837629 | -0.186583 | 0.242701 | 0.055321 | -0.024043 | -0.120276 | -0.699375 | 0.12211  |
| 4 | 4.589981  | 5.525321 | -0.296593 | 0.641106 | 1.701635 | -2.372474 | 1.211485  | -0.488318 | 1.48947  |

In [19]:

```
1 #Fazendo o split no conjunto reduzido
2 X_train_pca, X_test_pca, y_train, y_test = train_test_split(X_pca, y, test_size = 0.15,
```

In [20]:

```
1 from sklearn.linear_model import LogisticRegression
2
3 #2. Criando uma instância para estimar
4 logistic_regression_pca = LogisticRegression(n_jobs=-1, random_state=15)
5
6 #3. Usando os dados de treino para treinar o modelo
7 logistic_regression_pca.fit(X_train_pca,y_train);
8
9
```

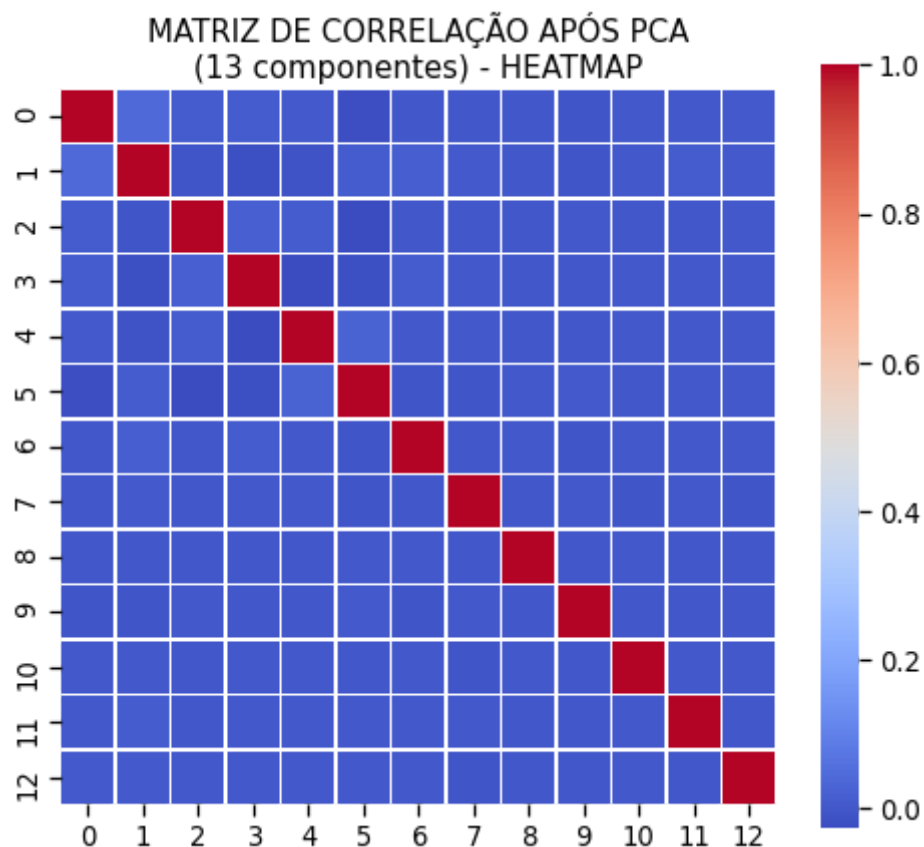
In [21]:

```
1 metrics_model_pca = model_metrics(X_test_pca,y_test,logistic_regression_pca)
```

## Nova correlação entre as variáveis

In [22]:

```
1 corr = X_train_pca.corr()
2 f,ax = plt.subplots(figsize=(8, 7))
3 sns.heatmap(corr, cbar = True, square = True, annot = False, fmt= '.1f',
4             xticklabels= True, yticklabels= True
5             ,cmap="coolwarm", linewidths=.5, ax=ax)
6 plt.title(f'MATRIZ DE CORRELAÇÃO APÓS PCA\n({len(X_test_pca.columns)} componentes) - HEATMAP')
```



Reduzindo o conjunto de dados às 9 principais componentes

## Apenas a título de comparação das métricas

In [23]:

```
1  #Definindo a dimensão dos dados
2  pca = PCA(9)
3  #PCA no conjunto de treino
4  pca_X = pca.fit(X)
5
6  X_pca = pd.DataFrame(pca_X.transform(X))
7
8  #Fazendo o split no conjunto reduzido
9  X_train_pca, X_test_pca, y_train, y_test = train_test_split(X_pca, y, test_size = 0.15,
10
11  #2. Criando uma instância para estimar
12  logistic_regression_pca = LogisticRegression(n_jobs=-1, random_state=15)
13
14  #3. Usando os dados de treino para treinar o modelo
15  logistic_regression_pca.fit(X_train_pca,y_train)
16
17  metrics_model_pca9 = model_metrics(X_test_pca,y_test,logistic_regression_pca)
18
19
```

## Avaliando as métricas

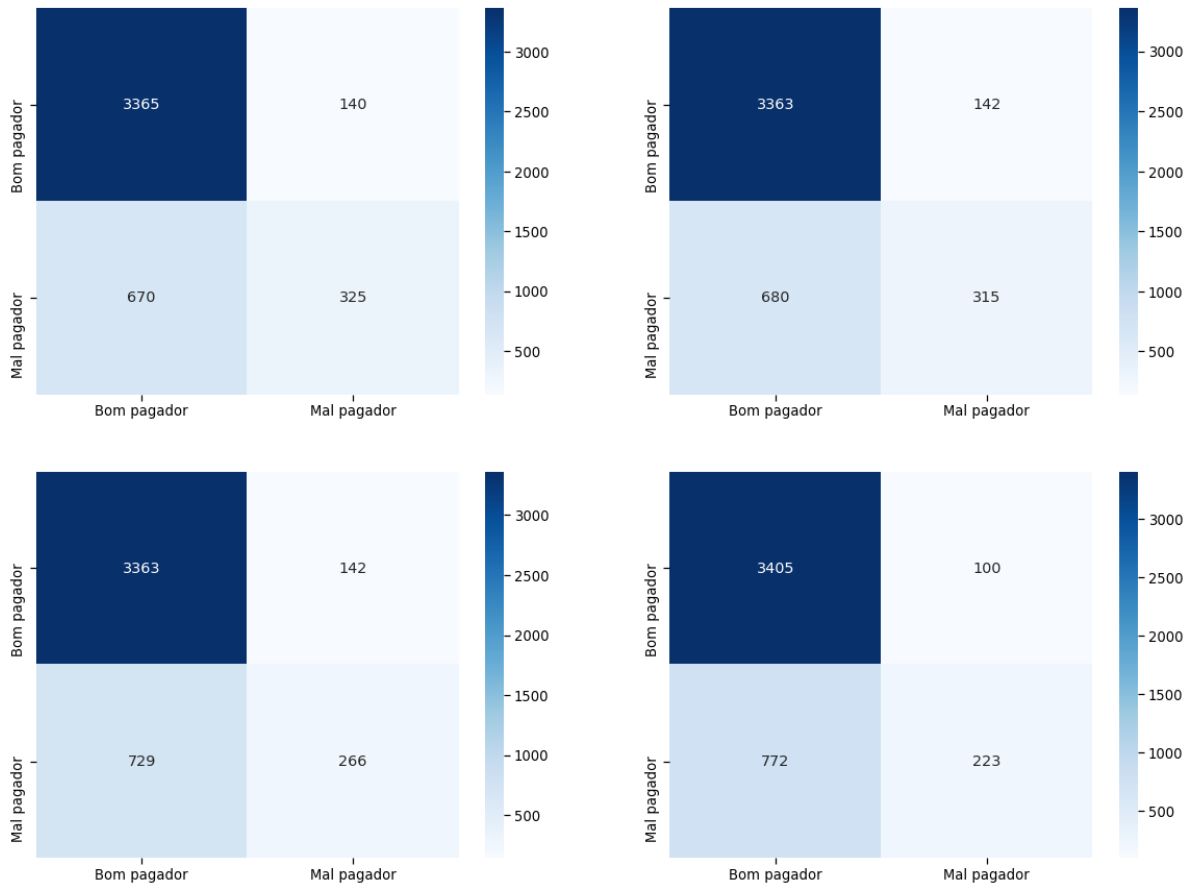
In [24]:

```
1  y_pred_test1 = logistic_regression.predict(X_test)
2  CM1 = confusion_matrix(y_pred = y_pred_test1, y_true = y_test)
3
4  y_pred_test2 = logistic_regression_2.predict(X_test_2)
5  CM2 = confusion_matrix(y_pred = y_pred_test2, y_true = y_test)
6
7  pca = PCA(13)
8  pca_X = pca.fit(X)
9  X_pca = pd.DataFrame(pca_X.transform(X))
10 X_train_pca, X_test_pca, y_train, y_test = train_test_split(X_pca, y, test_size = 0.15,
11 logistic_regression_pca = LogisticRegression(n_jobs=-1, random_state=15)
12 logistic_regression_pca.fit(X_train_pca,y_train)
13
14 y_pred_test_pca13 = logistic_regression_pca.predict(X_test_pca)
15 CM13 = confusion_matrix(y_pred = y_pred_test_pca13, y_true = y_test)
16
17 pca = PCA(9)
18 pca_X = pca.fit(X)
19 X_pca = pd.DataFrame(pca_X.transform(X))
20 X_train_pca, X_test_pca, y_train, y_test = train_test_split(X_pca, y, test_size = 0.15,
21 logistic_regression_pca = LogisticRegression(n_jobs=-1, random_state=15)
22 logistic_regression_pca.fit(X_train_pca,y_train)
23
24 y_pred_test_pca9 = logistic_regression_pca.predict(X_test_pca)
25 CM9 = confusion_matrix(y_pred = y_pred_test_pca9, y_true = y_test)
```

In [25]:

```
1
2 f, axes = plt.subplots(2, 2, figsize=(20, 15), facecolor='white')
3 f.suptitle(' Matrizes de Confusão')
4
5 plt.figure(figsize=(1, 1))
6 ax1 = sns.heatmap(CM1,annot=True, cmap="Blues", fmt="d",
7                   xticklabels = ['Bom pagador', 'Mal pagador'],
8                   yticklabels = ['Bom pagador', 'Mal pagador'], ax=axes[0,0])
9
10 plt.figure(figsize=(1, 1))
11 ax2 = sns.heatmap(CM2,annot=True, cmap="Blues", fmt="d",
12                  xticklabels = ['Bom pagador', 'Mal pagador'],
13                  yticklabels = ['Bom pagador', 'Mal pagador'], ax=axes[0,1])
14
15 plt.figure(figsize=(1, 1))
16 ax3 = sns.heatmap(CM13,annot=True, cmap="Blues", fmt="d",
17                  xticklabels = ['Bom pagador', 'Mal pagador'],
18                  yticklabels = ['Bom pagador', 'Mal pagador'], ax=axes[1,0])
19
20 plt.figure(figsize=(1, 1))
21 ax4 = sns.heatmap(CM9,annot=True, cmap="Blues", fmt="d",
22                  xticklabels = ['Bom pagador', 'Mal pagador'],
23                  yticklabels = ['Bom pagador', 'Mal pagador'], ax=axes[1,1]);
24
```

Matrizes de Confusão



<Figure size 72x72 with 0 Axes>

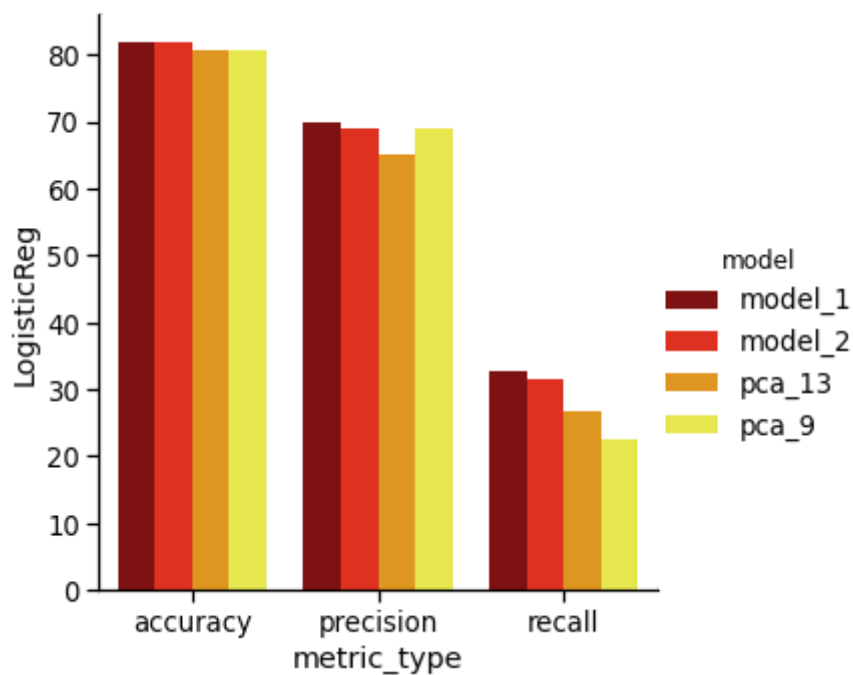
<Figure size 72x72 with 0 Axes>

<Figure size 72x72 with 0 Axes>

<Figure size 72x72 with 0 Axes>

In [26]:

```
1 metrics_model_2['model'] = 'model_2'
2 metrics_model_1['model'] = 'model_1'
3 metrics_model_pca['model'] = 'pca_13'
4 metrics_model_pca9['model'] = 'pca_9'
5 g = sns.catplot(data=pd.concat([metrics_model_1,metrics_model_2,metrics_model_pca,metrics_model_pca9]),
```



## Simulação

### Função para receber dados externos

In [27]:

```
1 def cadastro(cliente):
2
3     #transformando o dict em série
4     cliente = pd.Series(cliente)
5
6     #dimensionando e normalizando
7     info = cliente.values.reshape(1, -1)
8     info = robust_scaler.transform(info)
9
10    #como "info" e "X" são np array, vamos converter para dataframe
11    info = pd.DataFrame(info, columns = default.drop("default", axis = 1).columns)
12
13    #vamos concatenar essa linha aos dados originais
14    novos_dados = pd.concat([info, X])
15
16    return novos_dados
```

In [28]:

```
1 def teste(novos_dados):
2
3     #Definindo a dimensão do dado
4
5     pca_function = PCA(9)
6
7     try:
8         pca_function_teste = pca_function.fit(novos_dados)
9         dado_reduzido = pd.DataFrame(pca_function_teste.transform(novos_dados))
10
11
12         prob = logistic_regression_pca.predict_proba(dado_reduzido[:1])[0][1]
13
14         if prob >= 0.25:
15             print(f"Esse vai dever\n\nScore: {round(100/prob, 3)}\nProb de dever: {roun
16         else:
17             print(f"Camarada que paga!\n\nScore: {round(100/prob, 3)}\nProb de dever: {
18
19     except ValueError:
20         print("Dado com dimensão grande demais")
21
```

In [29]:

```
1 yuri = OrderedDict([("limit_bal", 4000), ("age", 30), ("pay_0", -5), ("pay_2", 0), ("pay
2
3 cliente = cadastro(yuri)
4 teste(cliente)
```

Camarada que paga!

Score: 1494.585

Prob de dever: 6.69%

In [30]:

```
1 luiz = OrderedDict([("limit_bal", 1000), ("age", 19), ("pay_0", 0), ("pay_2", 0), ("pay_
2
3 cliente2 = cadastro(luiz)
4 teste(cliente2)
```

Camarada que paga!

Score: 640.948

Prob de dever: 15.6%

In [31]:

```
1 gustavo = OrderedDict([("limit_bal", 30000), ("age", 19), ("pay_0", 3), ("pay_2", 4), ('
2
3 cliente3 = cadastro(gustavo)
4 teste(cliente3)
```

Esse vai dever

Score: 102.178

Prob de dever: 97.87%

In [32]:

```
1 ari = OrderedDict([("limit_bal", 3000), ("age", 19), ("pay_0", 0), ("pay_2", -1), ("pay_
2
3
4 cliente4 = cadastro(ari)
5 teste(cliente4)
```

Camarada que paga!

Score: 1908.339

Prob de dever: 5.24%

## Referencias e dados

Logistic Regression - THE MATH YOU SHOULD KNOW <https://www.youtube.com/watch?v=YMJtsYlp4kg>  
(<https://www.youtube.com/watch?v=YMJtsYlp4kg>)

<https://www.kaggle.com/selener/prediction-of-credit-card-default/notebook>  
(<https://www.kaggle.com/selener/prediction-of-credit-card-default/notebook>)

Logistic Regression Details Pt1: Coefficients <https://www.youtube.com/watch?v=BfKanl1aSG0>  
(<https://www.youtube.com/watch?v=BfKanl1aSG0>)

Logistic Regression Details Pt 2: Maximum Likelihood <https://www.youtube.com/watch?v=BfKanl1aSG0&t=383s>  
(<https://www.youtube.com/watch?v=BfKanl1aSG0&t=383s>)

Making Predictions with Data and Python : Predicting Credit Card Default | packtpub.com  
<https://www.youtube.com/watch?v=zUqa6KcwRhs&app=desktop> (<https://www.youtube.com/watch?v=zUqa6KcwRhs&app=desktop>)



