

# Predição dos danos sofridos pelas construções do Nepal causados pelo terremoto de 2015

Ana Gabriela Faria da Silva  
Bruno de Assis Silva  
João Phellip de Mello da Rocha  
Rodrigo Matheus Rocha de Medeiros

Universidade de São Paulo  
Instituto de Matemática e Estatística

27 de Novembro de 2020

# Introdução

- Em 2015 ocorreu um terremoto com magnitude de 7,8 na escala Richter que atingiu o Nepal, Índia, Bangladesh, Paquistão e China. O Nepal foi o país mais atingido, sendo este o terremoto mais violento a atingir o país em 81 anos.
- Após o terremoto, o governo local realizou uma grande pesquisa domiciliar para avaliar os danos às construções nos distritos afetados pelo terremoto.
- Os dados coletados são formados principalmente por informações sobre a estrutura das construções da região atingida, e do grau de dano sofrido por elas.

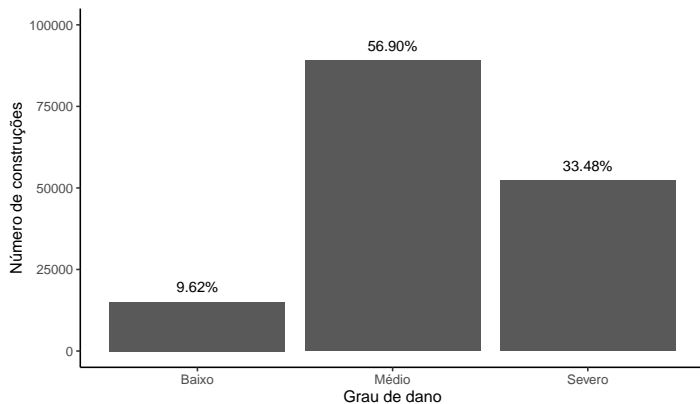
## Descrição dos dados

- Cada linha representa uma construção específica na região que foi atingida.
- A base de dados possui 260.601 observações com 38 variáveis preditoras, em que as variáveis preditoras se dividem entre 5 quantitativas e 33 qualitativas.
- Optamos por utilizar cerca de 60% da base de dados fornecida para treinar os modelos, e os outros 40% como conjunto de validação.
- O conjunto de treinamento resultante possui 156361 observações.

A variável resposta representa o nível de dano à construção causado pelo terremoto. Existem 3 graus de dano:

- 1 — representa grau baixo de dano;
- 2 — representa grau médio de dano;
- 3 — representa grau severo de dano.

# Grau de dano



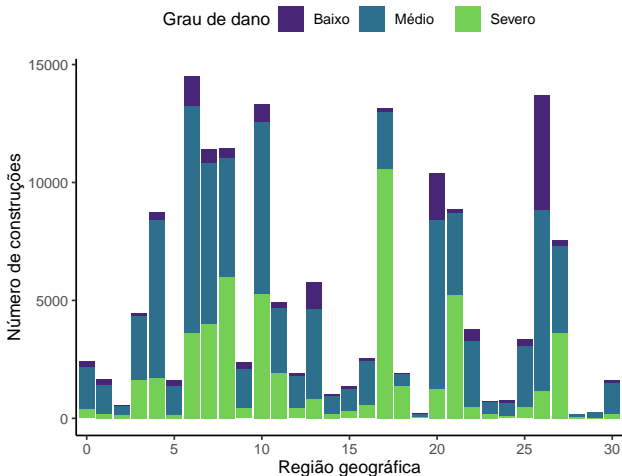
## Variáveis geográficas

- **geo\_level\_1\_id**
- **geo\_level\_2\_id**
- **geo\_level\_3\_id**

Região geográfica em que existem construções, do maior agrupamento (nível 1) à sub-região mais específica (nível 3). Possíveis valores: nível 1: 0–30, nível 2: 0–1427, nível 3: 0–12567;



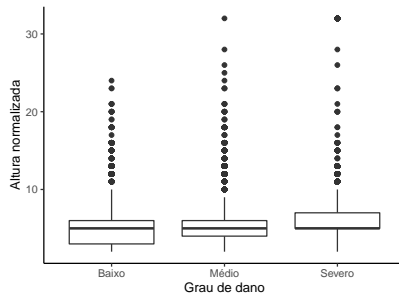
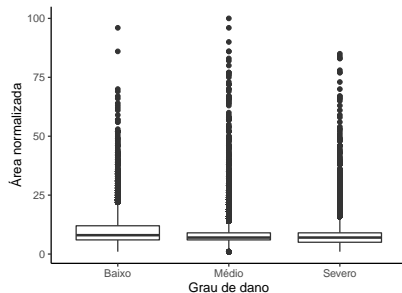
# Descrição dos dados



## Características da construção

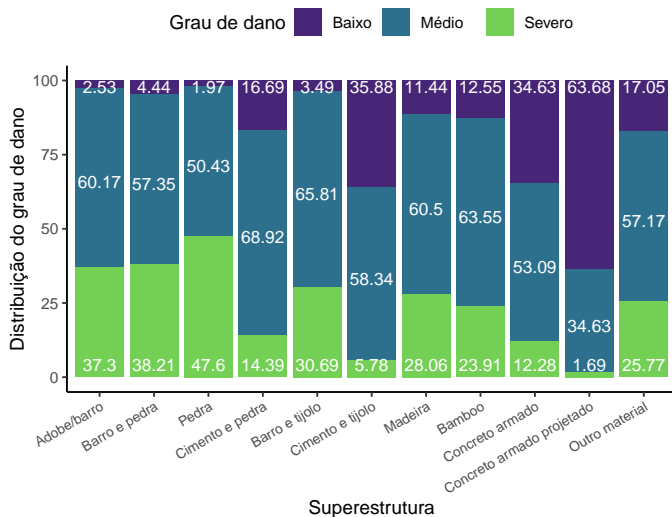
- **land\_surface\_condition** – Condição da superfície do terreno,
- **area\_percentage** – Área (normalizada),
- **height\_percentage** – Altura (normalizada),
- **age** – Idade,
- **count\_floors\_pre\_eq** – Número de andares antes do terremoto,
- **foundation\_type** – Tipo de fundação,
- **ground\_floor\_type** – Tipo de andar térreo,
- **other\_floor\_type** – Tipo de pisos,
- **roof\_type** – Tipo de telhado,
- **position** – Posição,
- **plan\_configuration** – Configuração do plano de construção,
- **legal\_ownership\_status** – Status legal,
- **count\_families** – Número de famílias.

# Descrição dos dados



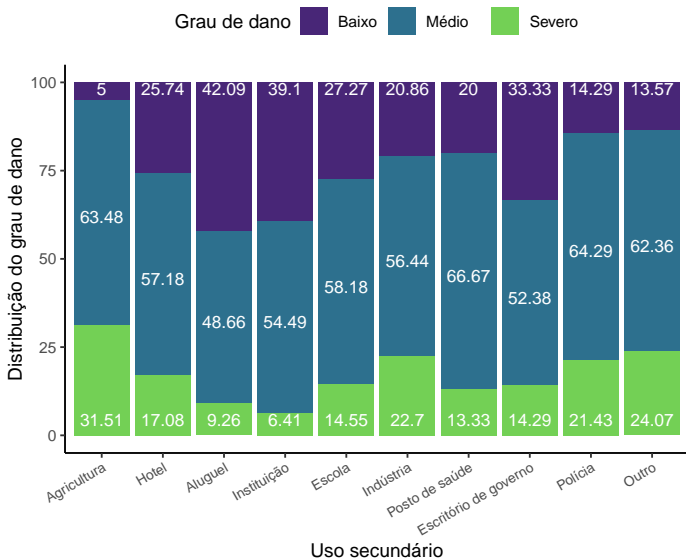
## Superestrutura

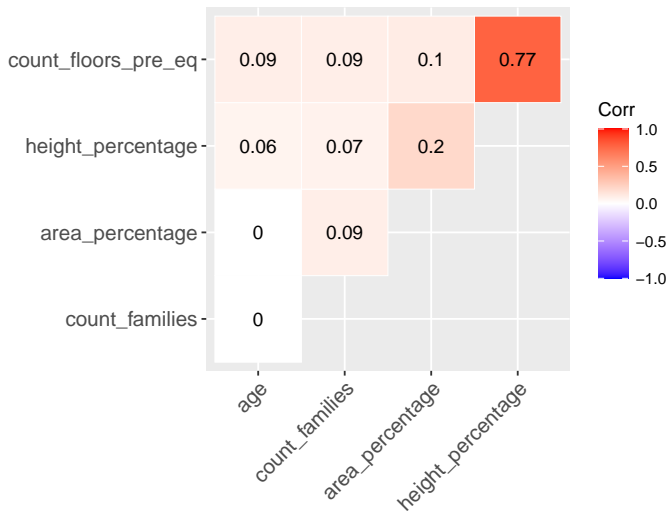
- **has\_superstructure\_adobe\_mud** – Feita de adobe/barro,
- **has\_superstructure\_mud\_mortar\_stone** barro e pedra,
- **has\_superstructure\_stone\_flag** – Feita de pedra,
- **has\_superstructure\_cement\_mortar\_stone** – Feita de cimento e pedra,
- **has\_superstructure\_mud\_mortar\_brick** – Feita de barro e tijolo,
- **has\_superstructure\_cement\_mortar\_brick** – Feita de cimento e tijolo,
- **has\_superstructure\_timber** – Feita de madeira,
- **has\_superstructure\_bamboo** – Feita de bamboo,
- **has\_superstructure\_rc\_non\_engineered** – Feita de concreto armado sem engenharia,
- **has\_superstructure\_rc\_engineered** – Feita de concreto armado projetado,
- **has\_superstructure\_other** – Feita de outro material.



## Uso secundário

- **has\_secondary\_use\_agriculture** – Utilizada para fins agrícolas.
- **has\_secondary\_use\_hotel** – Utilizada como um hotel.
- **has\_secondary\_use\_rental** – Utilizada para aluguel.
- **has\_secondary\_use\_institution** – Utilizada como uma instituição.
- **has\_secondary\_use\_school** – Utilizada como uma escola.
- **has\_secondary\_use\_industry** – Utilizada para fins industriais.
- **has\_secondary\_use\_health\_post** – Utilizada como um posto de saúde.
- **has\_secondary\_use\_gov\_office** – Utilizada como um escritório do governo.
- **has\_secondary\_use\_use\_police** – Utilizada como uma delegacia de polícia.
- **has\_secondary\_use\_other** – Utilizada secundariamente para outros fins.







# Modelos

- Regressão logística multinomial
- Regressão logística ordinal
- Classificação por árvores
- Máquinas de vetores de suporte (SVM)

# Regressão logística multinomial

Para a  $i$ -ésima construção do conjunto de treinamento, assuma que

$$\log \left\{ \frac{\Pr(G_i = k | x = x_i)}{\Pr(G_i = 1 | x = x_i)} \right\} = \beta_{k0} + \beta_k^\top x_i, \quad k = 2, 3,$$

em que  $G_i$  é o grau de dano sofrido pela  $i$ -ésima construção em sua codificação original,  $\beta_{20}, \beta_{30} \in \mathbb{R}$ , e  $\beta_1 = (\beta_{11}, \beta_{12}, \dots, \beta_{189})^\top \in \mathbb{R}^{89}$  e  $\beta_2 = (\beta_{21}, \beta_{22}, \dots, \beta_{289})^\top \in \mathbb{R}^{89}$  são os coeficientes associados as variáveis explicativas.

# Regressão logística multinomial

O modelo é especificado em termos das transformações *logit* em relação a classe 1, de menor dano. Com esta especificação as probabilidades a posteriori são dadas por

$$\Pr(G_i = k | x = x_i) = \frac{\exp\{\beta_{k0} + \beta_k^\top x_i\}}{1 + \sum_{l=2}^3 \exp\{\beta_{l0} + \beta_l^\top x_i\}}, \quad k = 2, 3$$

e

$$\Pr(G_i = 1 | x = x_i) = \frac{1}{1 + \sum_{l=2}^3 \exp\{\beta_{l0} + \beta_l^\top x_i\}}$$

# Regressão logística multinomial

Utilizamos a função `multinom()` do pacote do R `nnet` para ajustar o modelo aos dados do conjunto de treinamento.

```
fit <- multinom(damage_grade ~ ., data = damage)
```

# Regressão logística multinomial

A tabela a seguir mostra a matriz de confusão resultante da classificação no conjunto de validação.

		Predição		
		Baixo	Médio	Severo
Real	Baixo	3637	6248	199
	Médio	2376	47963	8943
	Severo	173	16651	18050

**Taxa de acertos:** 66,82%.

A tabela a seguir mostra o percentual de erros por classe no conjunto de validação.

**Table:** Percentual de erros por classe (%)

Baixo	Médio	Severo
63,93	19,09	48,24



## Regressão logística ordinal

Aplicamos o modelo ordinal quando o número de categorias da variável é maior que dois e elas são ordenadas. O modelo de logito cumulativo é definido como:

$$\text{logito}[P(Y_i \leq j|x_i)] = \log \frac{P(Y_i \leq j|x_i)}{1 - P(Y_i \leq j|x_i)}, \quad j = 1, \dots, c - 1$$

Supondo que as variáveis explicativas tenha diferentes efeitos temos o modelo de logito cumulativo sem chances proporcionais:

$$\text{logito}[P(Y_i \leq j|x_i)] = \alpha_j + \beta'_j \mathbf{x}_i$$

Utilizamos a função `vglm()` do pacote do R VGAM para ajustar o modelo aos dados do conjunto de treinamento.

```
fit <- vglm(damage_grade ~ ., data = damage)
```

Para o conjunto de validação foi obtida a seguinte matriz de confusão:

		Predição		
		Baixo	Médio	Severo
Real	Baixo	5261	4693	130
	Médio	4357	46771	8154
	Severo	338	17501	17035

**Taxa de acertos:** 66,17%.

O percentual de erros por classe no conjunto de validação.

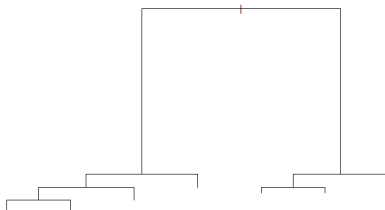
Table: Percentual de erros por classe (%)

Baixo	Médio	Severo
47,82	21,10	51,15

# Classificação por árvores

# Classificação por árvores

Modelos baseados em árvores têm como ideia central a segmentação do espaço de predição em regiões mais homogêneas, de acordo com a resposta.



# Classificação por árvores

## Vantagens:

- Podem ser aplicadas em problemas de regressão/classificação;
- Não precisam de variáveis dummy para lidar com preditores qualitativos;
- São fáceis de interpretar; entre outras.

## Desvantagens:

- Árvores não são muito robustas, ou seja, uma pequena mudança nos dados pode causar uma grande mudança na estimativa final da árvore.

## Possíveis soluções:

- *Bagging*
- *Random Forests*
- *Boosting*



# Classificação por árvores - Random Forests

Utilizamos a função `randomForest()` do pacote do R `RandomForest` para ajustar o modelo aos dados.

```
fit.flor <- randomForest(damage_grade ~ .,  
  data = damage_RF, ntree = 200,  
  mtry = 8, importance = TRUE)
```

Para tratar o desbalanceamento dos dados: `strata` e `sampsize`.

# Classificação por árvores - Random Forests

A tabela a seguir mostra a matriz de confusão resultante da classificação no conjunto de validação.

**Table:** Classificação predita versus valores reais.

		Valores reais					
		Desbalanceado			Balanceado		
		Baixo	Médio	Severo	Baixo	Médio	Severo
Preditos	Baixo	4447	1740	157	7983	8122	1078
	Médio	5494	50792	13727	1921	38991	9103
	Severo	143	6750	20990	180	12169	24693

**Taxa de acertos:** Desbalanceado = 73,13% e Balanceado = 68,75%.

# Classificação por árvores - Random Forests

A tabela a seguir mostra o percentual de erros por classe no conjunto de validação.

**Table:** Percentual de erros por classe (%)

	Baixo	Médio	Severo
RF: Desbalanceado	55,90	14,32	39,81
RF: Balanceado	20,83	34,23	29,19

# Classificação por árvores - Boosting

Utilizamos a função `xgboost()` do pacote do R `XGBoost` para ajustar o modelo aos dados.

XGBoost, que significa Extreme Gradient Boosting, é uma implementação específica do método Gradient Boosting para encontrar o melhor modelo de árvore. Ele tem se mostrado eficiente na solução de diversos problemas e é muito utilizado em competições.

```
xgb.fit <- xgboost(data = features_train,  
label = response_train,  
eta=0.2, max_depth = 10, min_child_weight = 5, gamma = 1,  
subsample = 0.8, colsample_bytree = 0.6,  
nrounds = 149, objective = "multi:softmax",  
num_class = 4, verbose = 0,)
```

Artigo dos autores originais: <https://arxiv.org/abs/1603.02754>

# Classificação por árvores - Boosting

A tabela a seguir mostra a matriz de confusão resultante da classificação no conjunto de validação.

**Table:** Classificação predita versus valores reais.

		Predição		
		Baixo	Médio	Severo
Real	Baixo	4563	5388	133
	Médio	1834	50309	7139
	Severo	153	13528	21193

**Taxa de acertos:**72,97%.

A tabela a seguir mostra o percentual de erros por classe no conjunto de validação.

Table: Percentual de erros por classe (%)

Baixo	Médio	Severo
54,75	15,14	39,23

# Máquinas de vetores de suporte (SVM)

Utilizamos a implementação do pacote `e1071` para treinar uma máquina de vetores de suporte sobre os dados.

Durante a etapa de pré processamento, as variáveis de entrada numéricas foram normalizadas linearmente no intervalo  $[0,1]$  e as variáveis preditoras do tipo categóricas foram transformadas em dummy; para o treinamento, utilizamos um Kernel do tipo `radial basis function` e uma estratégia de multiclassificação por votos: treinamos simultaneamente 3 problemas de classificação binários e a classificação final foi feita para a classe com maior incidência na resposta.

Em função da complexidade temporal do algoritmo, a primeira etapa de grid-search dos hiperparâmetros foi feita sobre um subconjunto de 10K; em seguida, uma busca refinada foi feita sobre um conjunto maior de 20K amostras. Utilizou-se uma estratégia de validação cruzada de 5-fold para ambos.



A tabela a seguir mostra a matriz de confusão resultante da classificação no conjunto de validação.

		Predição		
		Baixo	Médio	Severo
Real	Baixo	4024	5742	162
	Médio	2295	48172	8922
	Severo	184	15049	19690

**Taxa de acertos:** 68,96%.

## Desempenho dos modelos na competição

# Desempenho dos modelos na competição

A tabela a seguir mostra o desempenho de cada modelo na competição.

**Table:** Resultados por método aplicado.

Modelos	Conjunto de Teste - Competição(%)
Regressão Logística Multinomial	66,90
Regressão Logística Ordinal	66,22
Random Forests: Desbalanceado	72,90
Random Forests: Balanceado	65,87
Gradient Boosting	72,26
SVM Kernel Radial	68,82

- O melhor desempenho foi obtido, sem tratar o desbalanceamento dos dados, utilizando Random Forests.
- A competição conta com 3563 participantes, nós estamos na posição 456.
- O líder atual obteve uma acurácia de 75,58%.

Obrigado!