

PROJETO - CATÁSTROFES NATURAIS

Integração e Processamento Analítico de Informação -
2022/2023

Ana Araújo nº 59457
Francisco Vicente nº 59363
João Faia nº 47051
Tomás Oom nº 59447

Índice

Introdução	3
1 Descrição de fonte de dados	4
1.1 EM-DAT (Emergency Events Database)	4
1.2 World Bank	8
1.2.1 World Data Bank - GDP	8
1.2.2 World Data Bank - Total Population	8
1.2.3 World Data Bank - Unemployment	8
1.2.4 World Data Bank - Turismo	8
1.2.5 FAOSTAT - Agricultura	9
2 Pré-processamento de Dados.....	10
2.1 Dataset - DESASTRES NATURAIS.....	10
2.1.1 Valores em falta.....	10
2.1.2 Remoção e renomeação de colunas irrelevantes.....	12
2.2 Dataset - AGRICULTURA	14
2.2.1 Valores em falta.....	14
2.2.2 Remoção e renomeação de colunas.....	14
2.3 Dataset - GDP.....	15
2.3.1 Valores em falta.....	15
2.3.2 Remoção e renomeação de colunas.....	15
2.4 Dataset - TOTAL POPULATION.....	16
2.4.1 Valores em falta.....	16
2.4.2 Remoção e renomeação de colunas irrelevantes.....	16
2.5 Dataset - TURISMO	17
2.5.1 Valores em falta.....	17
2.5.2 Remoção e renomeação de colunas irrelevantes.....	17
2.6 Dataset - DESEMPREGO.....	18
2.6.1 Valores em falta.....	18
2.6.2 Remoção e renomeação de colunas irrelevantes.....	18
2.7 Uniformização dos dados	19
3 Diagrama dos Datasets	20
4 Processo de Negócio.....	21
5 Questões analíticas.....	22
6 Modelação Dimensional	23

6.1 Tipo de tabela de factos e grão	24
6.2 Tabelas de Dimensões	29
6.2.1 Dimensão Tipo - dimType	29
6.2.2 Dimensão Data - dimDate	30
6.2.3 Dimensão Localização - dimLocation	33
6.2.4 Dimensão Evento - dimEvent	36
6.3 Diagrama em estrela.....	39
7 Sistema ETL e Relatórios.....	40
7.1 Extração	40
7.2 Transformação.....	40
7.2.1 Tratamento da Tabela DESASTRES NATURAIS.....	41
7.2.2 Tratamento da Tabela AGRICULTURA	41
7.2.3 Tratamento da Tabela DESEMPREGO.....	42
7.2.4 Tratamento da Tabela GDP	42
7.2.5 Tratamento da Tabela POPULAÇÃO	43
7.2.6 Tratamento da Tabela TURISMO	43
7.3 Construção do modelo dimensional.....	43
7.3.1 Dimensões:	43
7.3.2 Tabela de Factos:.....	45
7.4 Carregamento	45
7.4.1 Dimensões:	45
7.4.2 Factos:.....	47
7.5 Execução dos scripts	47
8 Diagrama de fluxo de dados	48
9 Cubo de dados	49
10 Respostas às questões analíticas	49
10.1 Primeira Pergunta analítica	50
10.2 Segunda Pergunta Analítica.....	73
10.3 Terceira Pergunta Analítica	80
Conclusão:	85
Referências	86

Introdução

A análise de dados é uma atividade essencial em diversas áreas, nomeadamente em finanças, ciência, saúde e política, entre outras. Com a crescente geração e armazenamento de dados, a construção de um data warehouse torna-se uma estratégia fundamental para gerenciar e explorar as informações de forma eficiente e eficaz. A integração de diferentes fontes de dados num único repositório, permitindo o acesso e análise de dados de várias perspectivas e dimensões, é um dos principais benefícios de um data warehouse.

No contexto deste projeto, o objetivo é explorar os dados relacionados a catástrofes naturais e entender o seu impacto em áreas cruciais, como agricultura, desenvolvimento dos países e turismo. A escolha dessas áreas é relevante, uma vez que esses setores são altamente influenciados pelas condições climáticas e pelos eventos naturais.

Na primeira etapa do projeto, a procura e a identificação de fontes de dados relevantes para o processo de negócio escolhido foram realizadas. A análise dos dados, considerando aspectos como quantidade, a sua dimensionalidade e a sua adequabilidade, foi conduzida para selecionar um conjunto de dados rico e multidimensional. A seleção de dados globais para um longo período de tempo foi realizada para proporcionar uma visão abrangente e histórica dos impactos das catástrofes naturais.

Durante a análise dos dados, foi utilizada a linguagem de programação python com a ferramenta deepnote. A representação das informações em um diagrama permitiu a visualização da conexão entre as diferentes fontes de dados e a identificação das hierarquias que fornecerão riqueza informativa para o projeto.

Um dos maiores desafios deste projeto é a complexidade dos dados relacionados às catástrofes naturais. São muitas as variáveis envolvidas, como a intensidade do evento, a localização, o tipo de desastre, entre outros. Além disso, é importante considerar a variabilidade e a heterogeneidade dos dados de diferentes fontes. Por isso, a seleção cuidadosa de fontes de dados relevantes e a análise minuciosa dos dados são fundamentais para o sucesso do projeto.

Espera-se que a construção e modelagem do data warehouse forneça uma visão mais completa e profunda sobre os impactos das catástrofes naturais na agricultura, desenvolvimento dos países e turismo. A estrutura de hierarquias identificadas oferecerá a riqueza necessária para a navegação e construção do projeto, permitindo uma análise multidimensional e uma visão integrada dos dados.

1 Descrição de fonte de dados

1.1 EM-DAT (Emergency Events Database)

EM-DAT (Emergency Events Database) é uma base de dados global de eventos de emergência, mantida pelo Centro de Investigação sobre Epidemiologia de Desastres (CRED) da Universidade Católica de Louvain, na Bélgica. A base de dados foi criada em 1988 e contém informações sobre desastres naturais desde 1900 até ao presente.

A base de dados EM-DAT é uma das mais completas e confiáveis fontes de informação sobre desastres naturais em todo o mundo. Esta fonte de dados inclui informações sobre desastres como terremotos, furacões, cheias, secas, incêndios florestais e deslizamentos de terra, entre outros. Para cada evento, a base de dados registra informações como a localização, data, tipo de desastre, número de vítimas, danos materiais e medidas tomadas para prevenir ou mitigar os efeitos do desastre.

A base de dados EM-DAT é utilizada por organizações internacionais, governos e agências de ajuda humanitária para monitorizar a evolução e as tendências dos desastres naturais em todo o mundo. Além disso, a base de dados é amplamente utilizada por investigadores e cientistas para estudar as causas e os efeitos dos desastres naturais e para desenvolver estratégias de prevenção.

O conjunto de dados EM-DAT de desastres naturais utilizado neste projeto é uma versão da base de dados EM-DAT disponibilizada no Kaggle¹, contendo informações sobre desastres naturais ocorridos entre 1900 e 2023. O conjunto de dados inclui informações sobre os eventos, tais como localização, data, tipo de desastre, número de mortos e danos materiais. O objectivo de utilizar este conjunto de dados serve para fazer análises e estudos de desastres naturais no mundo e estudar os seus impactos nas áreas da agricultura, desenvolvimento de países e turismo. A descrição de todas as variáveis do conjunto de dados, está especificada na tabela 1.

Tabela 1. Informação sobre os dados da dataset EM-DAT

Coluna	Descrição	Tipo de dados	
Dis No	Identificador único (Year+Seq+ISO)	Categórico	
Informação sobre a identificação da catástrofe	Disaster Group	Nível mais abrangente de classificação de desastres e divide-os em três grupos principais: desastres naturais, desastres tecnológicos e desastres complexos. Neste dataset apenas estão incluídos os desastres naturais.	Categóricos
	Disaster Subgroup	Divide os desastres naturais em sete subgrupos: climáticos, hidrológicos, meteorológicos, geofísicos, extraterrestres, biológicos e mistos.	Categóricos
	Disaster Type	Nível mais específico de classificação que identifica o tipo de desastre natural, como furacão, terremoto, inundação, seca, etc.	Categóricos
	Disaster Subtype	Nível mais detalhado de classificação e fornece informações adicionais sobre o tipo de desastre, como sua intensidade, duração, características específicas, entre outros.	Categóricos
	Event Name	Nome do evento específico que causou o desastre natural registrado.	Categóricos
	Entry Criteria	Indica os critérios que devem ser atendidos para que um evento seja registrado na base de dados	Categóricos
	Origin	Origem ou causa do desastre natural registrado	Categóricos
	Dis Mag Value	Refere-se à magnitude ou intensidade de um desastre	Numérico
	Dis Mag Scale	Fornece informações sobre a escala utilizada para medir a magnitude ou intensidade de um desastre na coluna "Dis Mag Value"	Categóricos
	Associated Dis	Evento relacionado que ocorre no mesmo período de tempo ou nas mesmas áreas geográficas afetadas pelo desastre registrado na linha correspondente	Categóricos
	Associated Dis2	Outro evento relacionado que ocorre no mesmo período de tempo ou nas mesmas áreas geográficas afetadas pelo desastre registrado na linha correspondente	Categóricos

Tabela 1. Informação sobre os dados da dataset EM-DAT (continua)

Coluna	Descrição	Tipo de dados
Informação sobre a acção política e humanitária	OFDA Response	Booleanos
	Appeal	Booleanos
	Declaration	Booleanos
	Aid Contribution	Numérico
Informação sobre a localização	Country	Categóricos
	ISO	Categóricos
	Region	Categóricos
	Continent	Categóricos
	Location	Categóricos
	Latitude	Numéricos
	Longitude	Numéricos
	Local Time	Data e hora
	River Basin	Categóricos

Tabela 1. Informação sobre os dados da dataset EM-DAT (continua)

Coluna	Descrição	Tipo de dados
Informação temporal	Year	Numérico
	Seq	Numérico
	Start Year	Numéricos
	Start Month	Numéricos
	Start Day	Numéricos
	End Year	Numéricos
	End Month	Numéricos
	End day	Numéricos
Informação sobre os afetados	Total Deaths	Numéricos
	No injured	Numéricos
	No Affected	Numéricos
	No Homeless	Numéricos
	Total Affected	Numéricos
Informação sobre os danos monetários	Reconstruction Costs ('000 US\$)	Numéricos
	Insured Damages ('000 US\$)	Numéricos
	Total Damages ('000 US\$)	Numéricos
	CPI	Numéricos

1.2 World Bank

O World Bank Open Data é uma plataforma online que disponibiliza gratuitamente dados e indicadores económicos, sociais e ambientais de todos os países membros. A plataforma permite acesso a vários indicadores de desenvolvimento, incluindo dados sobre a pobreza, educação, saúde, meio ambiente, finanças e comércio, entre outros.

1.2.1 World Data Bank - GDP

Este dataset refere-se ao indicador de Produto Interno Bruto (PIB) em dólares dos países, que é uma medida amplamente utilizada para medir o tamanho e a saúde económica de um país. O PIB é calculado somando o valor de todos os bens e serviços finais produzidos dentro das fronteiras de um país em um determinado período. A base de dados é mantida pelo World Bank² e inclui informações anuais do PIB para diversos países, bem como regiões e grupos.

1.2.2 World Data Bank - Total Population

Este dataset contém informações sobre a população total de países no mundo. A fonte de dados é o World data bank³ e a unidade de medida é o número total de habitantes. Os dados podem ser filtrados por país, região e ano. A base de dados é atualizada regularmente com novos dados à medida que se tornam disponíveis.

1.2.3 World Data Bank - Unemployment

Esta é uma base de dados do World Data Bank⁴ que fornece informações sobre a taxa de desemprego total em percentual da força de trabalho total para países no mundo. A base de dados permite a visualização dos dados e fazer uma comparação entre os países. Os dados estão disponíveis desde 1991 e são atualizados anualmente.

1.2.4 World Data Bank - Turismo

O conjunto de dados referente ao turismo foi extraído do World Data Bank⁵ e detalha o número de cidadãos internacionais que entram no país (número de chegadas). O dataset apresenta uma lacuna de dados desde 1960 até 1994, sendo estes anuais desde 1995 até 2020, para todos os países/regiões.

Tabela 2. Informação sobre as variáveis dos dataset retirados do World Data Bank

Variável	Descrição
Series Name	O nome da série de dados relativamente a GDP, População, Desemprego e Turismo
Series Code	O código da série de dados relativamente a GDP, População, Desemprego e Turismo
Country Name	O nome do país ou região.
Country Code	O código do país ou região.
Year	O ano em que os dados foram coletados ou relatados

1.2.5 FAOSTAT - Agricultura

A FAOSTAT é a base de dados estatísticos da Organização das Nações Unidas para Alimentação e Agricultura (FAO)⁶. Contém informações detalhadas sobre a produção, comércio e consumo de alimentos e produtos agrícolas no mundo. Os dados da FAOSTAT são de fontes nacionais e internacionais, incluindo governos, organizações não-governamentais e organizações internacionais, e abrangem uma ampla variedade de tópicos, como agricultura, pesca, segurança alimentar, nutrição, uso da terra e mudança climática.

Tabela 3. Informação sobre as variáveis do dataset AGRICULTURA

Variável	Descrição
Domain Code	Código do domínio a que pertence o conjunto de dados (por exemplo, "QI" para dados de qualidade e consumo de alimentos)
Domain	O nome do domínio a que pertence o conjunto de dados.
Area Code (M49)	Código numérico único para a área geográfica em que os dados foram coletados, conforme definido pelo sistema M49 das Nações Unidas.
Area	O nome da área geográfica em que os dados foram coletados.
Element Code	Código numérico único para o tipo de informação que está a ser relatada (por exemplo, "5510" para dados de uso de água em agricultura).
Element	O nome do tipo de informação que está a ser relatada.
Item Code (CPC)	Código numérico único para o tipo de produto ou recurso agrícola em que os dados estão sendo relatados, conforme definido pelo Sistema de Classificação de Produtos das Nações Unidas (CPC).
Item	O nome do produto ou recurso agrícola em que os dados estão a ser relatados.
Year Code	Código numérico único para o ano em que os dados foram adquiridos.
Year	O ano em que os dados foram adquiridos.
Unit	A unidade de medida usada para os dados relatados.
Value	O valor numérico dos dados relatados em forma de índice (2014-2016 = 100).
Flag	Um código de um ou dois caracteres que indica a qualidade ou a fonte dos dados relatados.
Flag Description	Uma descrição do código da Flag que indica a qualidade ou a fonte dos dados relatados.

2 Pré-processamento de Dados

Uma primeira análise aos datasets revela uma necessidade de processamento dos dados, de modo a que estes fiquem uniformizados. O pré-processamento de dados é uma etapa fundamental na construção de um data warehouse, especialmente quando se trata de integrar múltiplas fontes de dados. Este envolve várias etapas, incluindo a limpeza dos dados, a transformação de dados para um formato adequado para análise e a integração de dados de diferentes fontes. A limpeza dos dados envolve a remoção de valores duplicados, valores em falta (ou imputação), valores inconsistentes e outros erros que podem afetar negativamente a qualidade dos mesmos. A integração dos dados é importante para combinar informações de diferentes fontes num formato consistente e uniformizado, permitindo garantir a qualidade dos dados da warehouse e a validade das análises efetuadas.

Valores nulos

Para garantir a qualidade dos dados, é importante substituir os valores em branco e os valores considerados desconhecidos por um valor apropriado, como o NA. É essencial garantir que todas as tabelas utilizadas no processo estejam preenchidas corretamente. A abordagem da substituição dos valores nulos por “NA”. Essa abordagem garante que os dados sejam mais precisos e confiáveis, proporcionando uma base sólida para análise posterior.

2.1 Dataset - DESASTRES NATURAIS

2.1.1 Valores em falta

Neste dataset começamos por ver os valores nulos e foram obtidos os seguintes:

Dis No	0
Year	0
Seq	0
Glide	14845
Disaster Group	0
Disaster Subgroup	0
Disaster Type	0
Disaster Subtype	3295
Disaster Subsubtype	15464
Event Name	12605
Country	0
ISO	0
Region	0
Continent	0
Location	1810
Origin	12515
Associated Dis	12994
Associated Dis2	15809
OFDA Response	14852
Appeal	14008
Declaration	13239

AID Contribution ('000 US\$)	15792
Dis Mag Value	11514
Dis Mag Scale	1217
Latitude	13801
Longitude	13801
Local Time	15420
River Basin	15239
Start Year	0
Start Month	395
Start Day	3610
End Year	0
End Month	700
End Day	3529
Total Deaths	4778
No Injured	12453
No Affected	6946
No Homeless	14104
Total Affected	4488
Reconstruction Costs ('000 US\$)	16534
Reconstruction Costs Adjusted ('000 US\$)	16534
Insured Damages ('000 US\$)	15459
Insured Damages Adjusted ('000 US\$)	15459
Total Damages ('000 US\$)	11199
Total Damages Adjusted ('000 US\$)	11203
CPI	39

Existem várias razões pelas quais um registo de um desastre natural pode ter valores em branco no dataset de catástrofes naturais de EM-DAT:

A informação pode não estar disponível: o EM-DAT é uma fonte de dados global que depende de relatórios de várias fontes. Em alguns casos, a informação pode não estar disponível ou pode não ter sido relatada.

A existência de erros de comunicação, onde podem ser incluídos os atrasos na comunicação de informações sobre desastres naturais.

A disponibilidade limitada de dados históricos pode ser uma das razões, uma vez que muitos desastres naturais ocorreram antes do surgimento do EM-DAT ou antes que a cobertura de dados fosse tão ampla como nos dias de hoje. Portanto, é possível haver lacunas na informação histórica.

A variabilidade na qualidade dos dados no EM-DAT, que depende de muitas fontes diferentes para coletar informações sobre desastres naturais. Algumas fontes podem fornecer dados mais precisos do que outras, o que pode afectar a qualidade dos dados.

As diferentes definições de desastres naturais usadas pelas diferentes fontes, o que pode resultar em diferenças na contagem e na categorização dos eventos.

Estas podem ser algumas das causas de valores em branco no dataset EM-DAT. No entanto, é importante salientar que o EM-DAT é uma das fontes mais completas de dados sobre desastres naturais e, apesar das limitações, é amplamente utilizado para análises e estudos sobre desastres naturais em todo o mundo.

Os valores em branco da sample set foram substituídos por “NA”, à excepção do campo “OFDA”. De acordo com as diretrizes descritas na EM-DAT, para a coluna OFDA, os valores em branco foram substituídos por "No" e não por NA como nas demais colunas. Deve-se ao facto de que, inicialmente, o registo deste campo era feito apenas com a opção "Yes" e não havia registos "No". Com base nessas informações, os valores em branco na coluna OFDA foram substituídos por "No".

2.1.2 Remoção e renomeação de colunas irrelevantes

As colunas Glide, Adm Level, Admin1 Code, Admin2 Code e Geo Locations foram removidas da EM-DAT para simplificar e tornar mais uniforme o conjunto de dados. A coluna Glide não era relevante para a EM-DAT. As colunas Adm Level, Admin1 Code e Admin2 Code eram informações relacionadas à divisão administrativa dos países, que não eram úteis para a análise de dados de desastres naturais. A informação de localização geográfica também não era relevante para a análise, uma vez que a EM-DAT já fornece informações sobre a localização dos desastres naturais através das colunas Country e Location. A remoção dessas colunas permitiu que o conjunto de dados ficasse mais organizado e fácil de ser trabalhado para análises posteriores.

Para facilitar o trabalho com os dados, o nome das variáveis foi renomeado:

Tabela 4. Renomeação das variáveis do dataset EM-DAT

Nome da variável	Novo nome da variável
Dis No	DisasterID
Disaster Group	Disaster_group
Disaster Subgroup	Disaster_subgroup
Disaster Type	Disaster_type
Disaster Subtype	Disaster_subsubtype
Event Name	Event_name
Associated Dis	Associated_dis
Associated Dis2	Associated_dis2
OFDA Response	OFDA_response
AID Contribution ('000 US\$)	AID_contribution
Dis Mag Value	Mag_value
Dis Mag Scale	Mag_scale
Local Time	local_time

Tabela 4. Renomeação das variáveis do dataset EM-DAT (continua)

Nome da variável	Novo nome da variável
River Basin	River_basin
Start Year	Start_year
Start Month	Start_month
Start Day	Start_day
End Year	End_year
End Month	End_month
End Day	End_day
Total Deaths	Total_deaths
No Injured	N_injured
No Affected'	N_affected
No Homeless	N_homeless
Total Affected	Total_affected
Reconstruction Costs ('000 US\$)	Reconstruction_costs
Reconstruction Costs Adjusted ('000 US\$)	Reconstruction_costs_adjusted
Insured Damages ('000 US\$)	Insured_damages
Insured Damages Adjusted ('000 US\$)	Insured_damages_adjusted
Total Damages ('000 US\$)	Total_damages
Total Damages Adjusted ('000 US\$)	Total_damages_adjusted

Nomes como “Antilles” não se referem a um país específico. Por essa razão, foram removidos do dataset EM-DAT para evitar ambiguidades e garantir a precisão dos dados. Além disso, a remoção de Antilles não afetou significativamente os resultados da análise, uma vez que os eventos registados em Antilles eram relativamente poucos em comparação com outros países individuais.

2.2 Dataset - AGRICULTURA

2.2.1 Valores em falta

O conjunto de dados relativos à agricultura não apresentou valores nulos, no entanto foi necessário realizar a remoção e a renomeação de algumas colunas.

2.2.2 Remoção e renomeação de colunas

Foi necessário realizar a remoção das colunas Domain code, Domain, Element code, Year code, Unit, Flag e Flag description, visto que estas não se consideraram relevantes para análise futura nem acrescentam informação extra ao dataset.

A renomeação das colunas Area, Item Code (CPC) e Area code (M49) foi efetuada de modo a uniformizar esta tabela com as restantes.

Tabela 5. Renomeação das variáveis do dataset AGRICULTURA

Nome da variável	Novo nome da variável
Area	Country
Item Code (CPC)	Item_code
Area Code (M49)	Country_code

Além da renomeação das colunas, foi realizada a uniformização dos valores da variável Area code (M49) tendo em conta que esta possuía os códigos dos países em formato de número e as restantes tabelas em formato de três caracteres. Para este efeito e tendo em conta que a coluna dos países/regiões já tinha sido uniformizada anteriormente, apenas se substituíram os valores pelas siglas presentes na tabela dos desastres naturais. Neste caso, o país “China Mainland” não possuía correspondência com nenhuma sigla por ser exclusivo do dataset da agricultura, assim foi atribuída a sigla “CHM”.

2.3 Dataset - GDP

2.3.1 Valores em falta

O conjunto de dados relativo ao produto interno bruto inclui diversos valores em falta:

Series Name	0
Series Code	0
Country	0
Country Code	0
1960	132
...	
2017	9
2018	9
2019	11
2020	14
2021	21

No total, são 3336 valores nulos distribuídos ao longo dos anos, sendo que nos anos mais antigos este valor é superior aos anos mais recentes. Todos os valores em falta foram substituídos por NA.

2.3.2 Remoção e renomeação de colunas

As colunas series name e series code não apresentavam qualquer tipo de informação relevante para as análises futuras, tendo sido eliminadas.

Para uniformizar com os restantes datasets, a coluna “Country Name” foi renomeada para “Country” e “Country Code” para “Country_code”:

Tabela 6. Renomeação das variáveis do dataset GDP

Nome da variável	Novo nome da variável
Country Name	Country
Country Code	Country_code

2.4 Dataset - TOTAL POPULATION

2.4.1 Valores em falta

Neste dataset verificamos que não existem valores em branco.

```
Country Name      0
Country Code      0
1960 [YR1960]    0
1961 [YR1961]    0
1962 [YR1962]    0
...
2017 [YR2017]    0
2018 [YR2018]    0
2019 [YR2019]    0
2020 [YR2020]    0
2021 [YR2021]    0
Length: 64, dtype: int64
```

2.4.2 Remoção e renomeação de colunas irrelevantes

As colunas "Series Name" e "Series Code" do dataset do World Bank não contêm informações relevantes para análise ou modelagem de dados. Apenas fornecem a descrição e o código da variável correspondente, que são úteis apenas para fins de identificação e referência. Uma vez que essas informações podem ser obtidas facilmente a partir do site do World Data Bank, elas não são necessárias para análise de dados e, portanto, podem ser removidas do conjunto de dados para simplificar e reduzir o tamanho do arquivo.

A coluna “Country Name” foi renomeada para “Country” e a variável “Country Code” para “Country_code”, de forma a ser mais fácil o manuseamento dos dados:

Tabela 7. Renomeação das variáveis do dataset TOTAL POPULATION

Nome da variável	Novo nome da variável
Country Name	Country
Country Code	Country_code

2.5 Dataset - TURISMO

2.5.1 Valores em falta

Foi verificado a existência de valores nulos (“nan”) no dataset. Estes valores estão principalmente entre os anos 1960 e 1994, inclusive, e 2021.

Country Name	0
Country Code	0
Indicator Name	0
Indicator Code	0
1960	266
...	
2017	32
2018	36
2019	43
2020	134
2021	266

Existem um total de 10692 de valores nulos, que dizem respeito aos primeiros anos. As colunas em branco nos primeiros anos no dataset de turismo podem ser devidas a vários factores, incluindo falta de dados disponíveis, mudanças na metodologia de recolha de dados ou falta de interesse ou capacidade de alguns países em reportar os seus dados. Também é possível que algumas regiões geográficas não tenham sido registadas como países independentes até anos mais recentes, o que pode resultar em valores em falta para essas regiões nos anos anteriores. Além disso, alguns países podem não ter tido uma indústria de turismo significativa ou não tiveram um sistema de recolha de dados robusto nos seus primeiros anos de existência. Esses fatores podem ter contribuído para os valores nulos nas primeiras colunas do dataset. A substituição foi efetuada de forma similar aos restantes datasets, introduzindo “NA”.

2.5.2 Remoção e renomeação de colunas irrelevantes

No seguimento da análise de valores nulos no dataset, prosseguiu-se à remoção das colunas entre 1960 e 1994 e 2021, por não incluírem qualquer tipo de informação. Além do referido, a coluna Indicator name e Indicator code, foram também removidas por serem apenas uma identificação.

A coluna “Country Name” foi renomeada para “Country” e a coluna “Country Code” foi renomeada para “Country_code” de forma a ser mais fácil o manuseamento dos dados:

Tabela 8. Renomeação das variáveis do dataset TURISMO

Nome da variável	Novo nome da variável
Country Name	Country
Country Code	Country_code

2.6 Dataset - DESEMPREGO

2.6.1 Valores em falta

```
Country Name      0
1991 [YR1991]    31
1992 [YR1992]    31
...
2018 [YR2018]    31
2019 [YR2019]    31
2020 [YR2020]    31
2021 [YR2021]    33
dtype: int64
```

Existem várias razões pelas quais um valor pode ser nulo no dataset de taxa de desemprego total do World Bank:

- Falta de dados: Em alguns países, pode haver falta de dados para um determinado ano ou período. Isso pode ocorrer porque o país não recolheu as informações necessárias ou porque os dados não foram disponibilizados pelo governo ou outras fontes.
- Metodologia de cálculo: A taxa de desemprego é calculada com uma metodologia específica, que pode variar de país para país ou ao longo do tempo. Se houver mudanças na metodologia de cálculo ou se diferentes fontes forem usadas, os dados podem não estar disponíveis para um determinado período.
- Erros de entrada de dados: É possível que os dados tenham sido inseridos incorretamente ou que tenham ocorrido erros durante a transmissão ou processamento dos dados.
- Diferenças culturais: As definições de desemprego podem variar entre países, o que pode levar a diferenças nos dados reportados. Alguns países podem ter uma definição mais ampla de desemprego, enquanto outros podem ter uma definição mais restrita.

Os dados em branco foram alterados para “NA”.

2.6.2 Remoção e renomeação de colunas irrelevantes

As colunas removidas deste conjunto de dados foram a Series Name e Series Code e os anos de 1960 a 1990, devido à alta quantidade de valores “NA”. De forma semelhante aos outros datasets, Country Name foi renomeado para Country, Country Code para Country_code e foram retirados os códigos dos anos.

Tabela 9. Renomeação das variáveis do dataset DESEMPREGO

Nome da variável	Novo nome da variável
Country Name	Country
Country Code	Country_code
1960 [YR1960]	1960

2.7 Uniformização dos dados

O conjunto de dados dos desastres naturais possui duas variáveis que conectam com as restantes tabelas, country e country_code, sendo necessário verificar a concordância entre estes em todas as tabelas onde ocorrem de modo a uniformizar o data warehouse. Após uma profunda análise dos países presentes em todos os datasets, concluiu-se que existem países que não estão presentes em alguns datasets mas estão presentes noutras e vice-versa. Estes países são, na sua maioria, colónias ou arquipélagos que não são reconhecidos como países independentes em todas as fontes de dados. Para realizar a uniformização dos dados, foi necessário avaliar quais as diferenças nos nomes dos países presentes em todos os datasets e trocar em todas as tabelas de modo a torná-los comparáveis.

As tabelas de dados foram extraídas da mesma fonte (World Data) e, por essa razão, foi-se verificar se as colunas que contêm nomes dos países têm valores idênticos, para que possam ser comparadas posteriormente. Assim, foi também realizada a uniformização dos valores errados num dos datasets e renomeou-se alguns países noutras datasets e comparou-se as colunas com os nomes dos países em diferentes datasets de forma a verificar se têm valores correspondentes.

Assim, criou-se uma lista de países exclusivos (sem repetições) para cada dataframe de dados (EM-DAT, AGRICULTURA, GDP, POPULAÇÃO, TURISMO e DESEMPREGO). Em seguida, verificou-se se os países em do dataset GDP são iguais aos países noutras datasets. Foi criada duas novas listas chamadas "diff1" e "diff2", que contêm países que estão presentes no dataset GDP e não no dataset TURISMO e vice-versa.

Depois, foi uniformizado alguns valores errados na coluna "Country" do dataset TURISMO. Em seguida, verificou-se os nomes dos países do dataset GDP correspondem aos nomes dos países nos outros datasets. Foi criado novamente duas listas "diff1" e "diff2", que agora contêm países que estão presentes no dataset GDP e não no dataset EM-DAT e vice-versa.

Por fim, foi renomeado alguns países na coluna "Country" do dataset EM-DAT. Algumas dessas alterações incluíram a remoção de erros de digitação, a mudança nomes de países para que correspondam aos nomes do dataset GDP e a alteração do nome da antiga Tchecoslováquia para o nome atual dos seus estados sucessores.

3 Diagrama dos Datasets

De modo a tornar mais evidente as relações entre cada dataset e facilitar o processo de criação da data warehouse, foi criado um diagrama com todos os conjuntos de dados onde se conectam as variáveis partilhadas pelos diversos datasets.

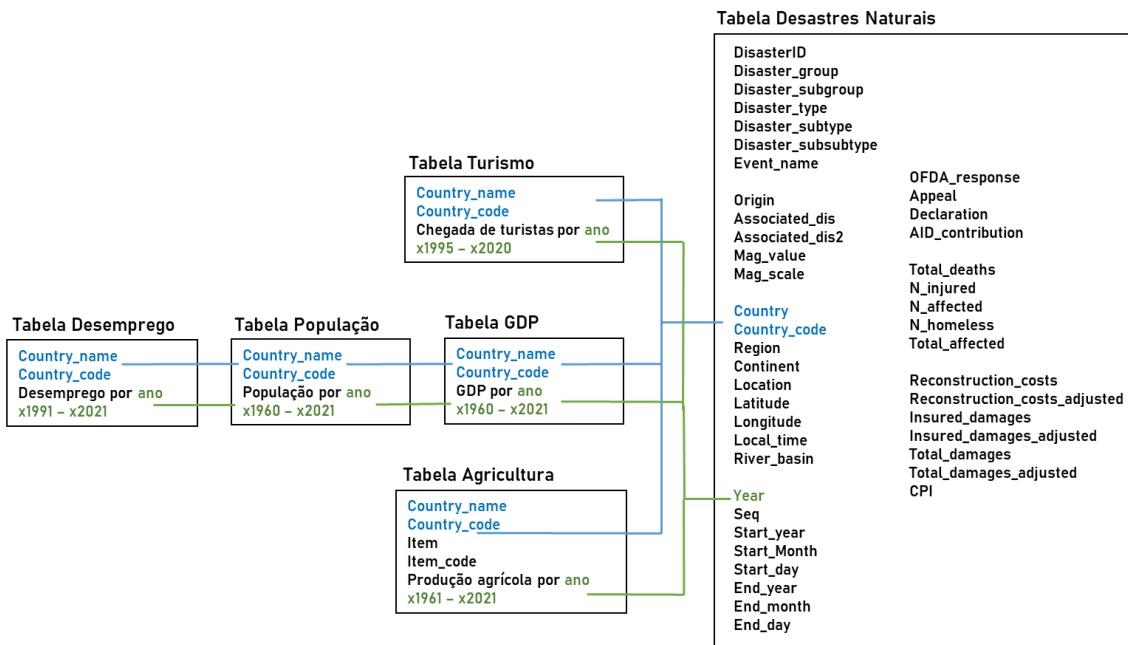


Figura 1. Diagrama dos datasets

4 Processo de Negócio

As catástrofes naturais têm um grande impacto na sociedade e no ambiente. Estas ocorrências podem ter consequências devastadoras, tais como a perda de vidas, destruição de propriedades e infraestruturas, impacto na economia e no desenvolvimento dos países. A fim de compreender e gerir melhor essas situações, é importante registar e monitorizar esses eventos.

Um dos principais objetivos de manter um registo de catástrofes naturais é estabelecer padrões que possam ajudar a prevenir ou mitigar os impactos desses eventos no futuro. Com a ajuda de dados recolhidos de diferentes fontes, podemos analisar e identificar tendências e padrões relacionados com o tipo de desastres, a sua localização geográfica e a sua intensidade. Estas informações são valiosas para tomar decisões estratégicas e planeamento a longo prazo para prevenir ou minimizar o impacto de futuras catástrofes naturais.

Por exemplo, os dados sobre catástrofes naturais podem ser utilizados para identificar áreas que são mais propensas a serem afetadas por eventos climáticos extremos, ajudando assim os governos a adotar medidas preventivas. Também podem ser utilizados para avaliar a viabilidade de um determinado setor, como o turismo ou a agricultura, em uma determinada região ou país, tendo em conta a probabilidade de ocorrerem eventos naturais adversos.

No entanto, é importante ter em conta a dimensionalidade dos dados dos *datasets* utilizados. Os dados sobre catástrofes naturais podem ser muito complexos e podem incluir informações de várias fontes diferentes. A interpretação desses dados pode ser difícil e requer um bom conhecimento do domínio. Além disso, é importante garantir a qualidade dos dados, pois informações imprecisas ou incompletas podem levar a tomadas de decisão erradas. Adicionalmente, neste *dataset* existe muita falta de informação para dados antes de 1960. A inexistência de informação para este período para levar a um estudo diferente e, consequentemente, a uma tomada de decisão inadequada.

Em resumo, a compilação de dados sobre catástrofes naturais é fundamental para entender melhor os eventos adversos e o seu impacto na sociedade e no meio ambiente. A análise desses dados pode levar a decisões mais informadas e estratégicas, que podem ter um impacto significativo na prevenção e mitigação de eventos naturais adversos no futuro. É crucial garantir a qualidade dos dados e estar ciente da complexidade dos mesmos, para que possam ser interpretados corretamente e utilizados para tomadas de decisão eficazes.

5 Questões analíticas

No seguimento dos objetivos do projeto proposto, foram elaboradas três perguntas analíticas que nos ajudem a estudar o impacto que as catástrofes naturais podem ter na agricultura, turismo e no desenvolvimento económico de um país:

- 1) Como as catástrofes naturais afetam a produção agrícola num determinado país e qual é o impacto na sua produção? Que tipo de desastre naturais são mais propensos a afetar a agricultura e quais os tipos de matérias-primas mais afetadas?
- 2) Existe alguma relação entre a afluência de turistas de um país e a frequência de ocorrência de catástrofes naturais? Se sim, como essas catástrofes afetam o desenvolvimento do setor turístico e a economia do país? Quais os tipos de desastres naturais que influenciam mais o turismo?
- 3) Qual é o impacto das catástrofes naturais no índice de desenvolvimento de crescimento de um país e como fica afetada a taxa de desemprego e a produtividade da população? A inflação tem um valor de expressão maior nos países não desenvolvidos?

6 Modelação Dimensional

A modelação dimensional é uma técnica utilizada no armazenamento de dados que organiza os dados em dimensões e medidas, criando uma representação multidimensional dos dados. O objetivo da modelação dimensional é facilitar aos utilizadores finais a análise e compreensão de dados complexos, fornecendo uma estrutura que reflete os processos empresariais e a forma como os utilizadores pensam sobre os dados. A modelação dimensional baseia-se no conceito de um cubo de dados, que representa a natureza multidimensional dos dados.

Na modelação dimensional, os dados são organizados em dois tipos de tabelas: tabelas de factos e tabelas de dimensões. As tabelas de factos contêm dados quantitativos aditivos, tais como o número de mortes provocado por cada desastre ou o custo de reconstrução gasto pelo país, enquanto as tabelas de dimensões fornecem o contexto para os dados na tabela de factos, tais como a data ou a localização desses desastres. Ao combinar tabelas de factos e tabelas de dimensões, é possível criar um modelo de dados flexível e poderoso que pode suportar uma vasta gama de consultas analíticas. As tabelas de dimensões conectam-se à tabela de factos através de chaves estrangeiras, permitindo a cada facto da tabela de factos (com a sua própria chave primária) estar conectado às diversas dimensões. Deste modo, a tabela de factos juntamente com as tabelas de dimensões, permitem perceber quais as características relacionadas com cada desastre natural e analisar as diversas medidas de forma integrada.

As linhas da tabela de factos incluem o seu próprio ID, ou uma chave primária, que identifica unicamente um desastre natural, conectada às diversas dimensões pelas suas chaves substitutas. Primeiramente, a tabela de factos está associada à dimensão Tipo onde explicita o tipo de desastre natural, seguida da dimensão Data onde está incluída a informação no tempo do acontecimento (esta não está conectada diretamente à tabela de factos, pois existe uma vista SQL com os atributos da data de início do desastre e a data de fim do mesmo). Além das referidas, a dimensão Localização também é de extrema relevância com detalhes importantes sobre o país ou região onde se desenrolou o desastre e, adicionalmente, com informação relativa ao rendimento do país, o nível de turismo ou o seu desempenho na agricultura. Por fim, a última ligação da tabela de factos é com a dimensão Evento, onde inclui diversas características que descrevem o evento como a sua origem (se deriva de outro desastre), se tem desastres associados (terremoto gerar um *tsunami*) ou se o país recebeu ajuda após o desastre.

6.1 Tipo de tabela de factos e grão

A tabela de factos é do tipo transacional, uma vez que capta, em cada linha, uma transação do processo de negócio, sendo possível calcular a sua frequência e por vezes a duração de transações. Deste modo, é possível identificar o grão da tabela de factos como cada linha, ou cada transação, da tabela de factos. O grão permite identificar o nível de detalhe apresentado na linha da tabela de factos, sendo este diretamente associado ao número de dimensões conectadas e, por conseguinte, associadas ao número de chaves estrangeiras presentes nos atributos da tabela de factos.

Com base na informação anterior, a tabela de factos (tabela 10), factDisaster, irá incluir 4 dimensões com atributos filtrados de acordo com as necessidades do processo de negócio e a capacidade da *data warehouse*. Estas são a dimType, a dimDate, a dimLocation e a dimEvent. As medidas presentes na tabela de factos são essenciais para responder às questões propostas na etapa anterior e futuras questões que hipotéticos decisores possam necessitar de responder. Estas incluem as mortes provocadas pelo desastre (Total_deaths), o número de feridos (N_injured), o número de afetados (N_affected), o número de desalojados (N_homeless), a contribuição monetária de outros países (AID_contribution), o custo de reconstrução, tanto com e sem a inflação considerada (Reconstruction_costs e Reconstruction_costs_adjusted), o custo em danos cobertos pelos seguros, tanto com e sem a inflação considerada (Insured_damages e Insured_damages_adjusted) e por último o custo total de danos (Total_damages). A tabela de factos possui uma dimensão de 16568 linhas, sendo este o número de desastres registados. Tendo em conta que, teoricamente, a tabela de factos deveria conter 90% do tamanho armazenado da data warehouse, é possível prever, desde já, que não será o caso desta estratégia de modelação, existindo dimensões com maior número de linhas. Outra questão relevante para esta etapa do projeto e que terá de ser resolvida na próxima etapa é a presença de valores nulos na tabela de factos (ver Figuras 2 e 3). Tendo em conta que são de uma magnitude bastante elevada, a imputação de valores não fazia sentido. Além disso, não seria benéfico retirarmos as medidas por completo só por apresentarem bastantes valores nulos, visto que são variáveis que fazem sentido para caracterizar os desastres ocorridos. Por exemplo, um desastre natural do tipo “Drought” não vai ter qualquer tipo de impacto nas infraestruturas do país, ao contrário de um desastre do tipo “Tornado”.

Tabela 10. Tabela de factos factDisaster - Descrição dos campos de dados: Identificadores, atributos e medidas aditivas.

Campo	Descrição dos dados	Origem dos dados	Valores
DisasterID	Chave Primária (Identificador do desastre natural)	Tabela DESASTRES NATURAIS (ID gerado sequencialmente)	N/AP
TypeKEY	Chave estrangeira (dimensão Tipo)	N/AP	N/AP
LocationKEY	Chave estrangeira (dimensão Localização)	N/AP	N/AP
StartDateKEY	Chave estrangeira (dimensão Data)	N/AP	N/AP
EndDateKEY	Chave estrangeira (dimensão Data)	N/AP	N/AP
EventKEY	Chave estrangeira (dimensão Evento)	N/AP	N/AP
Total_deaths	Número total de pessoas que morreram em consequência do desastre natural	Tabela DESASTRES NATURAIS	1-3700000 Média: 2758.87

Tabela 10. Tabela de factos factDisaster - Descrição dos campos de dados: Identificadores, atributos e medidas aditivas. (continua)

Campo	Descrição dos dados	Origem dos dados	Valores
N_injured	Número de pessoas que ficaram feridas como resultado do desastre natural	Tabela DESASTRES NATURAIS	1-1800000 Média: 2559.52
N_affected	Número de pessoas afetadas pelos desastres naturais	Tabela DESASTRES NATURAIS	1-330000000 Média: 873742.22
N_homeless	Número de pessoas que ficaram desalojadas como resultado do desastre natural.	Tabela DESASTRES NATURAIS	3-15850000 Média: 72361.08
Total_affected	Número total de pessoas afetadas pelo desastre natural	Tabela DESASTRES NATURAIS	1-330000000 Média: 711587.56
AID_contribution	Ajuda financeira por países ou organizações para ajudar no alívio ou recuperação de uma situação de desastre.	Tabela DESASTRES NATURAIS	1.0-3518530.0 Média: 20039.59
Reconstruction_costs	Custo estimado para reconstruir as áreas afetadas pelo desastre	Tabela DESASTRES NATURAIS	84.0-25000000.0 Média: 2543789.88
Reconstruction_costs_adjusted	custo estimado de reconstrução dos danos causados pelo desastre, ajustado à inflação e à variação cambial	Tabela DESASTRES NATURAIS	126.0-43922383.0 Média: 3700856.44
Insured_damage	Valor estimado dos danos segurados causados pelo desastre	Tabela DESASTRES NATURAIS	34.0-60000000.0 Média: 909690.38
Insured_damages_adjusted	Valor dos danos causados pelo desastre que foram segurados pelas vítimas ou por empresas, ajustado à inflação e à variação cambial para refletir os valores atuais.	Tabela DESASTRES NATURAIS	46.0-89913156.0 Média: 1248945.49
Total_damages	Custo total estimado dos danos causados pelo desastre	Tabela DESASTRES NATURAIS	2.0-210000000.0 Média: 780384.98

Tamanho: 16568 linhas

Tabela 11. Amostra exemplificativa da tabela de factos factDisaster.

DisasterID	object	TypeKey	int64	LocationKey	int64	StartDateKey	int...	Reconstruction...	Insured_damag...	Insured_damag...	Total_damages
1900-9002-...	_ 0%	0 - 53		0 - 6145		0 - 9814					
1900-9001-...	_ 0%										
16566 others	_ 100%										
1900-9002-CPV		0		0		0		nan	nan	nan	nan
1900-9001-IND		0		1		0		nan	nan	nan	nan
1902-0012-GTM		1		2		1		nan	nan	nan	25000.0
1902-0003-GTM		2		2		2		nan	nan	nan	nan
1902-0010-GTM		2		2		3		nan	nan	nan	nan
...											
2023-0110-ZMB		27		6137		9807		nan	nan	nan	nan
2023-0068-ZMB		6		6137		9813		nan	nan	nan	nan
2023-0095-ZWE		4		6144		9806		nan	nan	nan	nan
2023-0022-SRB		10		6145		9814		nan	nan	nan	nan

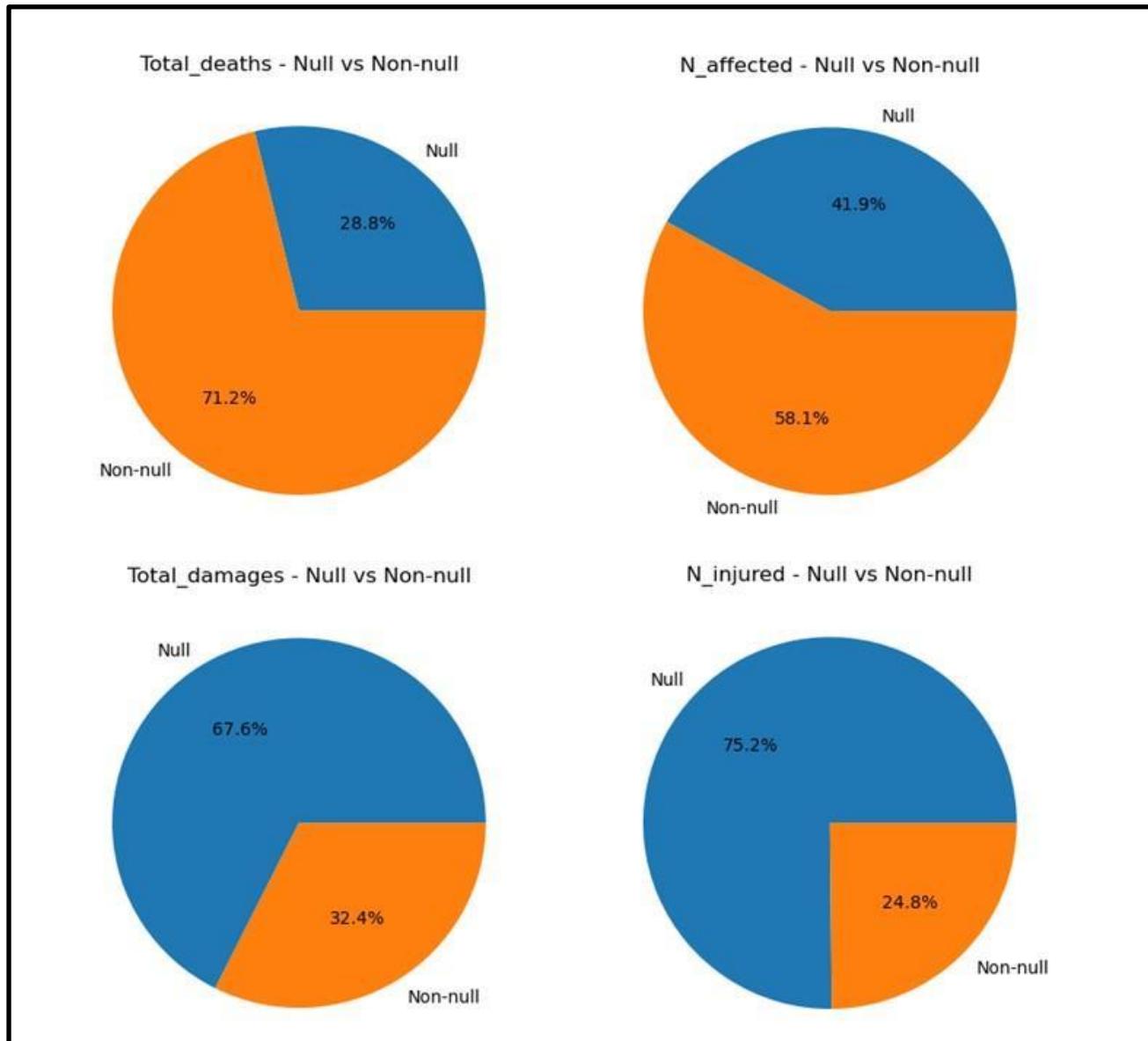


Figura 2. Gráficos de valores nulos para “Total_deaths”, “N_affected”, “Total_damage” e “N_injured” da tabela factos factDisaster.

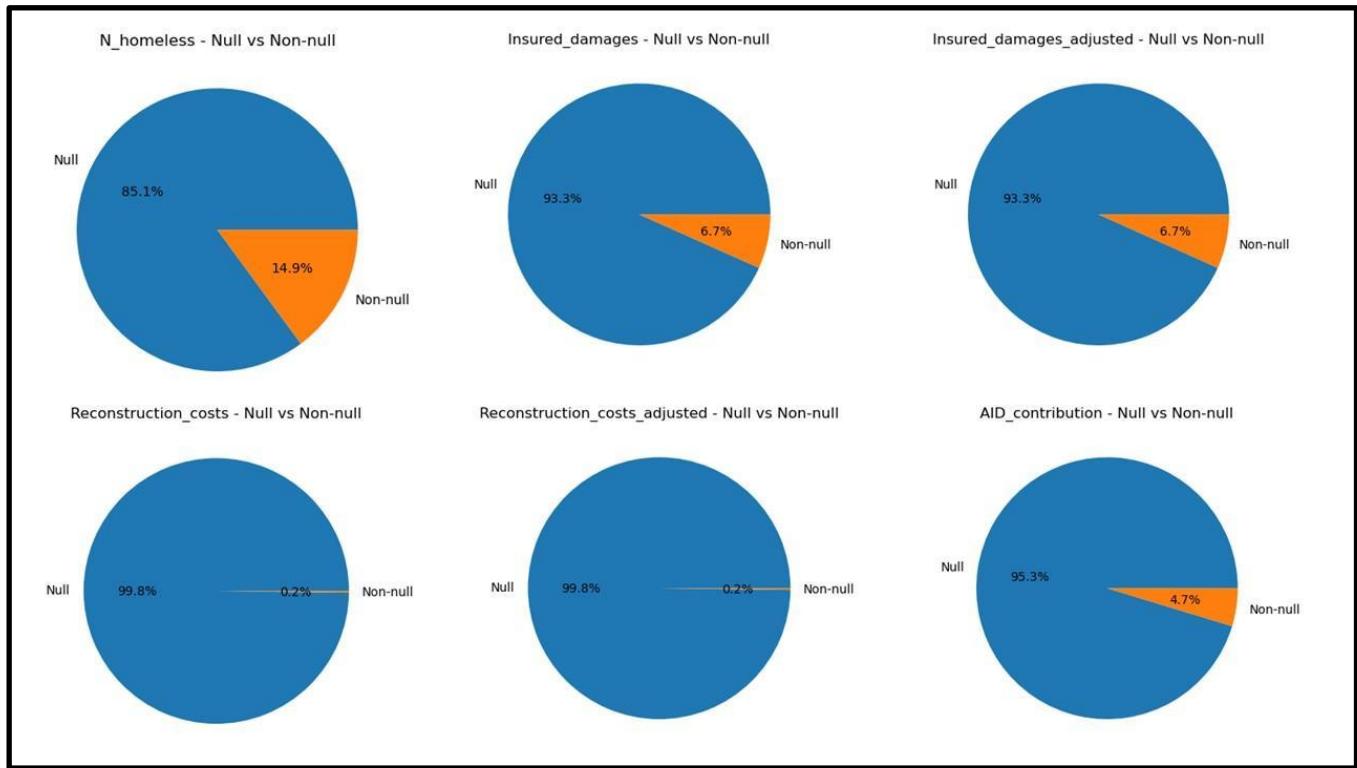


Figura 3. Gráficos circulares de valores nulos para “N_homeless”, “N_damaged”, “Insured_damage_adjusted”, “Reconstruction_costs”, “Reconstruction_costs_adjusted” e “AID_contribution” da tabela factos factDisaster.

6.2 Tabelas de Dimensões

A criação das dimensões são um passo essencial na modelação, tendo em conta que é isso que define o nível de detalhe refletido na tabela de factos. Como referido anteriormente, foram criadas quatro dimensões com características diferentes que permitem caracterizar cada linha da tabela de factos. O modelo de dados é completamente dependente desta fase do projeto, sendo de seguida caracterizado cada dimensão de modo a elucidar o seu papel e as suas características no data warehouse.

6.2.1 Dimensão Tipo - dimType

A dimensão tipo tem como objetivo a caracterização do desastre da tabela de factos identificando o grupo à qual pertence (Meteorológico, geofísico, etc.), o tipo dentro desse grupo (Tempestade, atividade vulcânica, etc.), o subtípico e o “subsubtipo”. Foi gerado um ID, TypeKey, como chave substituta para identificar cada linha da tabela e irá ser usada pela tabela de factos para identificar cada tipo de desastre. Originalmente no dataset dos desastres naturais, existia uma variável que identificava o grupo dos desastres (Disaster_group), no entanto, esta não continha nenhuma informação relevante, uma vez que apenas continha a categoria “Natural disaster”. Deste modo, esta variável não foi incluída na dimensão tipo, mas, tendo em conta que a nomenclatura seria bastante útil para tornar a dimensão mais compreensível, decidiu-se renomear a variável Disaster_subgroup para Disaster_group e manter as Disaster_type intactas. Assim, a dimensão tipo originou uma tabela com quatro atributos: Disaster_group, Disaster_type, Disaster_subtype e Disaster_subsubtype.

É possível compreender por observação direta dos atributos que existe uma hierarquia associada onde vai do mais abrangente, Disaster_group até ao mais detalhado, Disaster_subsubtype. Este tipo de hierarquias permitem realizar as operações de drill-down e roll-up permitindo ao utilizador escolher o nível de detalhe pretendido.

A dimensão Tipo originou uma tabela com 54 linhas e com diversos valores nulos, sendo estes mais abrangentes no último nível da hierarquia, Disaster_subsubtype. Isto deve-se ao facto de este atributo ser muito específico e a maior parte dos desastres apenas serem caracterizados por um grupo, tipo e subtípico.

Tabela 12. Tabela de dimensão Tipo - dimType. Descrição dos campos de dados: Identificadores, atributos e exemplos de valores.

Campo	Descrição dos dados	Origem dos dados	Valores
TypeKey	Chave Substituta	ID gerado sequencialmente	1-54
Disaster_group	Classificação mais abrangente de desastres naturais	Tabela DESASTRES NATURAIS	Ex: Meteorological
Disaster_type	Identifica o tipo de desastre natural	Tabela DESASTRES NATURAIS	Ex: Storm
Disaster_subtype	Nível mais detalhado de classificação de desastre	Tabela DESASTRES NATURAIS	Ex: Convective storm NA: 12
Disaster_subsubtype	Nível mais detalhado de classificação de desastre	Tabela DESASTRES NATURAIS	Ex: Tornado NA: 42

Hierarquias: Disaster_group > Disaster_type > Disaster_subtype > Disaster_subsubtype

Tamanho: 54 linhas

Tabela 13. Amostra exemplificativa da tabela de dimensão Tipo - dimType

TypeKey int64	Disaster_group	Disaster_type	Disaster_subtype	Disaster_subsub...
0 - 53	Meteorologi...	Storm	22.2%	Not Applica...
	Hydrological ..	Landslide	13%	Convective ...
	4 others	13 others	64.8%	27 others
0	Climatological	Drought	Drought	Not Applicable
1	Geophysical	Earthquake	Ground movement	Not Applicable
2	Geophysical	Volcanic activity	Ash fall	Not Applicable
3	Geophysical	Mass movement (dry)	Rockfall	Not Applicable
4	Meteorological	Storm	Tropical cyclone	Not Applicable
...				
50	Meteorological	Storm	Convective storm	Derecho
51	Geophysical	Volcanic activity	Pyroclastic flow	Not Applicable
52	Climatological	Drought	Not Applicable	Not Applicable
53	Climatological	Glacial lake outburst	Not Applicable	Not Applicable

6.2.2 Dimensão Data - dimDate

A dimensão data tem um papel crucial na maioria dos data warehouses, sendo esta que determina todas as características temporais associadas a uma transação. Na modelação apresentada neste projeto, a dimensão data define temporalmente um desastre natural através dos seus atributos, tendo sido gerada uma chave substituta para ser associada a cada linha da tabela de factos (DateKey). Os atributos presentes na dimensão foram gerados tendo por base a informação temporal presente no dataset dos desastres naturais, no entanto, diversos novos atributos foram gerados a partir dos existentes para criar alguma redundância e oferecer uma maior flexibilidade ao utilizador, enriquecendo o presente modelo. Os atributos estabelecidos na dimensão e retirados diretamente do dataset foram o ano, o mês, o dia e o horário local (continham a hora e os minutos do desastre). A partir destes, foram gerados o nome do mês, a década (tendo em conta que se trata de desastres e alguns podem acontecer raramente), o semestre, ou trimestre, o dia da semana (número), o dia da semana (nome), as horas, os minutos e a estação do ano (importante para desastres como a seca ou tempestades). A análise destes atributos permite identificar duas hierarquias que também oferecem alguma flexibilidade ao utilizador. A primeira e de maior extensão começa no atributo mais geral, neste caso, a década e passa pelo ano, semestre, trimestre, mês (com o nome do mês associado), dia (com o número e nome do dia da semana associado), horas e minutos. A segunda hierarquia começa na década, passa pelo ano e acaba na estação do ano.

Tendo em conta a natureza dos desastres naturais, onde alguns podem ser momentâneos como um terramoto e outros podem durar dias, meses ou mesmo anos como uma pandemia, existe, não só, a necessidade de indicar a data de início do desastre, mas também, a data do fim do desastre. O facto de não ser possível utilizar duas chaves para uma mesma dimensão, obriga à criação de uma vista SQL, onde, numa fase futura, a dimensão data irá estar associada à dimStartDate e dimEndDate. Estas duas novas vistas vão estar conectadas à tabela de factos pelas chaves StartDateKey e EndDateKey, permitindo a existência de duas datas para o mesmo desastre. O facto de se

utilizar a data em dois contextos diferentes é um caso de role-playing, onde a data tem o papel de início e fim de um desastre.

A dimensão data originou uma tabela com 19221 linhas que acabou por se estender para além da tabela de factos. Esta situação deriva do caso de role-playing descrito anteriormente, onde o facto dos desastres terem duas datas associadas fez com que existissem mais combinações a ter em conta na dimensão data. A dimensão só tem mais 2653 linhas do que a tabela de factos, não sendo um valor exageradamente grande e, por isso, não é considerada uma dimensão monstro. Para além do número de linhas, a presença de valores nulos na dimensão data é uma situação indesejável, mas inevitável no contexto deste problema. Neste caso, tendo novamente em conta a natureza dos desastres naturais, é plausível que existam valores nulos pelo simples facto que alguns desastres não possuem minutos, horas, dias ou meses pois são desastres graduais como a seca. Nestes casos, também não faria sentido imputar os valores, tendo os valores nulos sido substituídos por “Unknown”.

Tabela 14. Tabela de dimensão Data - dimDate. Descrição dos campos de dados: Identificadores, atributos e medidas aditivas.

Campo	Descrição dos dados	Origem dos dados	Valores
DateKey	Chave Substituta	ID gerado sequencialmente	1-19220
Year	Número do ano	Tabela DESASTRES NATURAIS	1900-2023
Month	Número do mês	Tabela DESASTRES NATURAIS	1-12 NA: 193
MonthName	Nome do mês	Dado gerado com recurso a linguagem de programação	Ex: January NA: 193
Day	Dia do mês	Tabela DESASTRES NATURAIS	Ex: 1 NA: 1850
Decade	Década onde ocorreu o desastre	Dado gerado com recurso a linguagem de programação	Ex: 1980s
Semester	Número do semestre no respectivo ano	Dado gerado com recurso a linguagem de programação	Ex: 1st Semester NA: 193
Season	Estação do ano	Dado gerado com recurso a linguagem de programação	Ex: Winter NA: 1851
Quarter	Número do trimestre no respectivo ano	Dado gerado com recurso a linguagem de programação	Ex: 2nd Quarter NA: 193
Weekday	Época sazonal	Dado gerado com recurso a linguagem de programação	1-7 NA: 1851
WeekDayName	Nome dia da semana	Dado gerado com recurso a linguagem de programação	Ex: Monday NA: 1851
LocalTime	Hora do incidente	Dado gerado com recurso a linguagem de programação	Ex: 20:20 NA: 8683
Hour	Hora em que ocorreu	Dado gerado com recurso a linguagem de programação	0-24 NA: 8683
Minutes	Hora em que ocorreu	Dado gerado com recurso a linguagem de programação	0-59 NA: 8683

Hierarquias: Decade > Year > Semester > Quarter > Month (MonthName) > Day (Weekday) > Hour > Minutes
 Decade > Year > Season

Tamanho: 19221 linhas

Tabela 15. Amostra exemplificativa da tabela de dimensão Data - dimDate

DateKey int64	Year object	Month object	MonthName obj...	WeekdayName o...	Local_time object	Hour object	Minutes object
0 - 19220	2002	2.7%	8	9.6%	August	9.6%	
	2021	2.6%	9	9.5%	September	9.5%	
	122 others	94.8%	11 others	80.9%	11 others	80.9%	
0	1900	Unknown	Unknown	Unknown	Unknown	Unknown	Unknown
1	1902	4	April	Friday	20:20	20	20
2	1902	4	April	Tuesday	Unknown	Unknown	Unknown
3	1902	10	October	Friday	Unknown	Unknown	Unknown
4	1903	4	April	Wednesday	Unknown	Unknown	Unknown
19217	2023	2	February	Monday	Unknown	Unknown	Unknown
19218	2023	1	January	Thursday	Unknown	Unknown	Unknown
19219	2023	3	March	Thursday	Unknown	Unknown	Unknown
19220	2023	1	January	Friday	Unknown	Unknown	Unknown

6.2.3 Dimensão Localização - dimLocation

A dimensão localização é outra de extrema importância para o presente modelo de negócio, pois cada desastre tem de estar associado a um local (ou a vários locais). Além do referido, a importância desta dimensão torna-se muito mais elevada quando também estão associados aos locais diversas variáveis que avaliam as características desses países/regiões, como o produto interno bruto, a população, o desemprego, o turismo e a agricultura, atributos esses que são do interesse deste projeto, pois as questões produzidas futuramente na análise são em parte dependentes dos mesmos. Como referido nas restantes dimensões, foi gerada uma chave substituta (LocationKey) para identificar cada linha e se associar a um desastre. Os atributos criados nesta dimensão, provêm do dataset dos desastres naturais e dos restantes que contém as variáveis mencionadas acima. Decidiu-se não incluir a coluna “Location” como um atributo, visto que não seria útil para análises futuras considerando que os restantes atributos (PIB, população, agricultura, etc) só têm valores para o país. As colunas “RiverBasin”, “Longitude” e “Latitude” também foram retiradas pois só se aplicavam a desastres com origem na água ou a terramoto onde se eram registados os epicentros. A renomeação das colunas foi também um passo importante na criação desta dimensão, para tornar a sua interpretação mais compreensível. A coluna “GDP” foi alterada para o atributo “CountryGDP”, a coluna “Tourism” foi alterada para “CountryTourism”, a coluna “Unemployment” foi alterada para “CountryUnemployment”, a coluna “population” foi alterada para “CountryPopulation”, a coluna “Food” foi alterada para “CountryAgricultureFood” e a “Non Food” foi alterada para “CountryAgricultureNonFood”. Os restantes atributos criados foram o país (“Country”), o código do país (“CountryCode”), a região (“Region”) e o continente (“Continent”).

Após uma análise à dimensão, é possível perceber que se trata de uma dimensão de mudança lenta, visto que, devido à adição dos atributos referentes ao PIB, turismo, desemprego, etc. que são atualizados uma vez por ano, esta dimensão também só será atualizada uma vez por ano. Além disso, todas estas variáveis possuem valores anuais históricos desde 1960 o que indica que existe uma atualização por ano. Para resolver este problema, uma estratégia do tipo 2 foi aplicada, sendo criada uma linha para cada novo ano e três novos atributos:

“RowEffectiveDate”, onde se detalha a data do início do ano, “RowExpirationDate”, onde se detalha a data do fim do ano e “CurrentRowIndicator”, onde se torna explícito se o ano é o atual ou não (através da categoria “Current” ou “Expired”). Tendo em conta que a chave substituta seria nova para cada linha gerada com este novo sistema, foi necessário atribuir uma chave supernatural para agrupar as linhas que apenas foram alterando anualmente, sendo utilizado o atributo do país, que era único nesta dimensão. Outro possível problema que esta solução criou foi a possibilidade de esta dimensão se tornar monstro, pois cada atributo foi multiplicado pelo número de anos presentes no histórico, mas a sua dimensão (13418 linhas) acabou por não passar o número de linhas da tabela de factos.

Outra filtragem que foi efetuada nesta fase foi a decisão de apenas incluir os atributos “CountryAgricultureFood” e “CountryAgricultureNonFood” do dataset da agricultura. Considerou-se que a inclusão de todas as categorias disponíveis iria tornar a dimensão bastante complexa e prejudicar a simplicidade do modelo. Foi considerado realizar um *outtrigger*, no entanto, a realização de *snowflaking* não seria o indicado nesta situação, tendo em conta que esse nível de detalhe não compensaria a complexidade acrescentada ao modelo.

Os valores nulos nesta dimensão foram inevitáveis, tendo em conta que os atributos do PIB, população, desemprego, etc. só tinham dados a partir de 1960, mas os desastres estão registados desde 1900. Assim, para se manter uma coerência, mantivemos essas linhas na dimensão para que os desastres estivessem conectados sempre à linha referente ao ano a que se deu o desastre. Outros valores nulos não existem por possível falta de medição para certos países, não fazendo sentido, novamente, imputar os valores em falta. Neste caso, os valores nulos foram substituídos por “Unknown”.

Por último, foi também identificada uma hierarquia nesta dimensão, começando no continente, passando pela região (parte do continente) e terminando no país que inclui todos os restantes atributos associados ao mesmo.

Tabela 16. Tabela de dimensão Localização - dimLocation. Descrição dos campos de dados: Identificadores, atributos e medidas aditivas.

Campo	Descrição dos dados	Origem dos dados	Valores
LocationKey	Chave Substituta	ID gerado sequencialmente	1-6146
Country	País em que se insere a localização	Tabela DESASTRES NATURAIS	Ex: Índia
CountryCode	Código do país onde ocorreu o desastre	Tabela DESASTRES NATURAIS	Ex: IND
Region	Região onde ocorreu o desastre	Tabela DESASTRES NATURAIS	Ex: Southern Asia
Continent	Continente onde ocorreu o desastre	Tabela DESASTRES NATURAIS	Ex: Asia
CountryGDP	Produto interno bruto, registado em dólares americanos	Tabela GDP	Ex: 159594492.3 NA: 3551
CountryTourism	Número de cidadãos internacionais que entram no país	Tabela TURISMO	Ex: 2403074088.0 NA: 9045
CountryUnemployment	Taxa de desemprego total em percentual da força de trabalho total para países no mundo.	Tabela DESEMPREGO	Ex: 5.543163701 NA: 7747
CountryPopulation	População total no país	Tabela POPULATION	Ex: 81200.0 NA: 1018
CountryAgricultureFood	Produção per capita de produtos alimentares	Tabela AGRICULTURA	Ex: 76.33 NA: 2875
CountryAgricultureNonFood	Produção per capita de produtos não-alimentares	Tabela AGRICULTURA	Ex: 163.88 NA: 4131
RowEffectiveDate	Ano em que iniciou o desastre	Dado gerado com recurso a linguagem de programação	Ex: 2012-01-01
RowExpirationDate	Mês em que acabou o desastre	Dado gerado com recurso a linguagem de programação	Ex: 2012-12-31
CurrentRowIndicator	Indica o estado do desastre (se o desastre ainda está acontecer ou não)	Dado gerado com recurso a linguagem de programação	Ex: expired

Hierarquias: Continent > Region > Country (CountryCode, CountryGDP, CountryTourism, CountryPopulation, CountryAgricultureType, CountryAgricultureValue, CountryUnemployment)

Tamanho: 13418 linhas

Tabela 17. Amostra exemplificativa da tabela de dimensão Localização - dimLocation

LocationKey int64	Country object	Country_code o...	Region object	CountryAgricult...	RowEffectiveDate	RowExpirationD...	CurrentRowIndi...
0 ~ 13417	China 0.8% United States ... 0.8% 226 others 98.5%	CHN 0.8% USA 0.8% 226 others 98.5%	Caribbean 9.1% Western Asia ... 8.3% 21 others 82.5%				
0	Cabo Verde	CPV	Western Africa	Unknown	1900-01-01	1900-12-31	expired
1	India	IND	Southern Asia	Unknown	1900-01-01	1900-12-31	expired
2	Guatemala	GTM	Central America	Unknown	1902-01-01	1902-12-31	expired
3	Canada	CAN	Northern America	Unknown	1903-01-01	1903-12-31	expired
4	Comoros	COM	Eastern Africa	Unknown	1903-01-01	1903-12-31	expired
...							
13414	Zimbabwe	ZWE	Eastern Africa	72.64	2004-01-01	2004-12-31	expired
13415	Zimbabwe	ZWE	Eastern Africa	46.68	2006-01-01	2006-12-31	expired
13416	Zimbabwe	ZWE	Eastern Africa	85.24	2012-01-01	2012-12-31	expired
13417	Zimbabwe	ZWE	Eastern Africa	113.73	2020-01-01	2020-12-31	expired

6.2.4 Dimensão Evento - dimEvent

A tabela de dimensão Evento contém informações mais detalhadas sobre cada evento, incluindo o nome atribuído, a origem, outros desastres associados, a escala e magnitude do evento, se houve ajuda humanitária e se o estado de emergência foi declarado.

Na tabela de dimensão Evento, fez-se uma renomeação das colunas para uma maior compreensão. Aqui considerámos as seguintes renomeações: 'Event_name' por 'EventName', 'Associated_dis' por 'AssociatedDisaster', 'Associated_dis2' por 'AssociatedDisaster2', 'Mag_value' por 'MagnitudeValue', 'Mag_scale' por 'MagnitudeScale' e 'OFDA_response' por 'OFDA_Response'.

A dimensão "Evento" contém 8809 linhas, com vários valores nulos. No campo "EventName", os valores nulos foram substituídos por "Name not specified", já que há desastres registados que não receberam um nome específico atribuído. Para os campos "Origin", "MagnitudeValue", "Appeal" e "Declaration", os valores nulos foram designados como desconhecidos, pois não existem registos desses dados. Para os campos "AssociatedDisaster" e "AssociatedDisaster2", os valores nulos foram substituídos por "No Disaster Associated", pois não foram registados outros desastres naturais que ocorressem simultaneamente e que fossem causados pelo desastre natural em questão. Para os campos "Magnitude", "Scale" e "OFDA_response", os valores nulos foram substituídos por "Not Applicable", pois não se aplica atribuir esses atributos em determinadas situações. Para o campo "CPI", os valores nulos foram substituídos por "Not Determined", pois em certos anos não foi calculada a variação dos preços dos bens e serviços consumidos pelas famílias ao longo do tempo do desastre em questão.

Existe uma hierarquia associada onde vai do mais abrangente, Continente, até ao mais detalhado, Country. Este tipo de hierarquia é uma hierarquia de profundidade, permitindo realizar as operações de drill-down e roll-up e escolher o nível de detalhe pretendido.

Tabela 18. Tabela de dimensão Localização - dimEvent. Descrição dos campos de dados: Identificadores, atributos e medidas aditivas.

Campo	Descrição dos dados	Origem dos dados	Valores
EventKey	Chave Substituta	ID gerado sequencialmente	1 - 8809
EventName	Nome atribuído ao evento	Tabela DESASTRES NATURAIS	Ex: Thomas NA: 5866
Origin	Origem ou causa do desastre	Tabela DESASTRES NATURAIS	Ex: Earthquake NA: 5987
AssociatedDisaster	Evento relacionado que ocorre no mesmo período de tempo ou nas mesmas áreas geográficas afetadas pelo desastre	Tabela DESASTRES NATURAIS	Ex: Slide (land, mud, snow, rock) NA: 5944
AssociatedDisaster2	Outro evento relacionado que ocorre no mesmo período de tempo ou nas mesmas áreas geográficas afetadas pelo desastre	Tabela DESASTRES NATURAIS	Ex: Tsunami/Tidal wave NA: 8100
MagnitudeValue	magnitude ou intensidade de um desastre	Tabela DESASTRES NATURAIS	Ex: 145.0 NA: 4630
MagnitudeScale	informações sobre a escala utilizada para medir a magnitude	Tabela DESASTRES NATURAIS	Ex: Ritcher NA: 620
OFDA_Response	Intervenções da Office of U.S. Foreign Disaster Assistance (OFDA),	Tabela DESASTRES NATURAIS	Ex: No
Appeal	Indica se houve um apelo de ajuda humanitária	Tabela DESASTRES NATURAIS	Ex: No NA: 7158
Declaration	Indica se houve uma declaração oficial de estado de emergência	Tabela DESASTRES NATURAIS	Ex: No NA: 6524
CPI	Mede a variação dos preços de um conjunto de bens e serviços consumidos pelas famílias ao longo do tempo	Tabela DESASTRES NATURAIS	Ex: 3,077011162 NA: 31

Hierarquias: Origin > AssociatedDisaster > AssociatedDisaster2

Tamanho: 8809 linhas

Tabela 19. Amostra exemplificativa da tabela de dimensão Evento - dimEvent

EventKey int64	EventName object	Origin object	AssociatedDisaster... object	OFDA_Response object	Appeal object	Declaration object	CPI object
0 - 8808	<p>Name not specified 66.6%</p> <p>Unknown Origin 68%</p> <p>Cholera 1.6%</p> <p>Heavy rains 8.2%</p> <p>Slide (land, ... 10.9%</p> <p>1612 others 31.8%</p> <p>705 others 23.9%</p>	<p>No Disaster Associated 67.5%</p> <p>Slide (land, ... 10.9%</p> <p>30 others 21.7%</p>	<p>No 85.8%</p> <p>Yes 14.2%</p>	<p>Unknown 81.3%</p> <p>No 16.3%</p> <p>Yes 2.4%</p>	<p>Unknown 74.1%</p> <p>No 16%</p> <p>Yes 9.9%</p>	<p>70,848,7927 3.8%</p> <p>66,731,05799 3.5%</p> <p>114 others 92.7%</p>	
0	Name not specified	Unknown Origin	Famine	No	No	No	2,849,084,409
1	Name not specified	Unknown Origin	No Disaster Associated	No	No	No	2,849,084,409
2	Name not specified	Unknown Origin	Tsunami/Tidal wave	No	Unknown	Unknown	2,963,047,785
3	Santa Maria	Unknown Origin	No Disaster Associated	No	Unknown	Unknown	2,963,047,785
4	Name not specified	Unknown Origin	No Disaster Associated	No	Unknown	Unknown	3,077,011,162
8805				No	Unknown	Unknown	Not Determined
8806				No	Unknown	Yes	Not Determined
8807				No	Unknown	Unknown	Not Determined
8808				No	Unknown	Unknown	Not Determined

6.3 Diagrama em estrela

O diagrama em estrela ajuda a visualizar as relações entre as tabelas de factos e dimensões criadas e, assim, facilita a compreensão das informações armazenadas. Aqui, conseguimos visualizar a relação entre as tabelas dimLocation, dimDate, dimEvent e dimType com a tabela factDisaster.

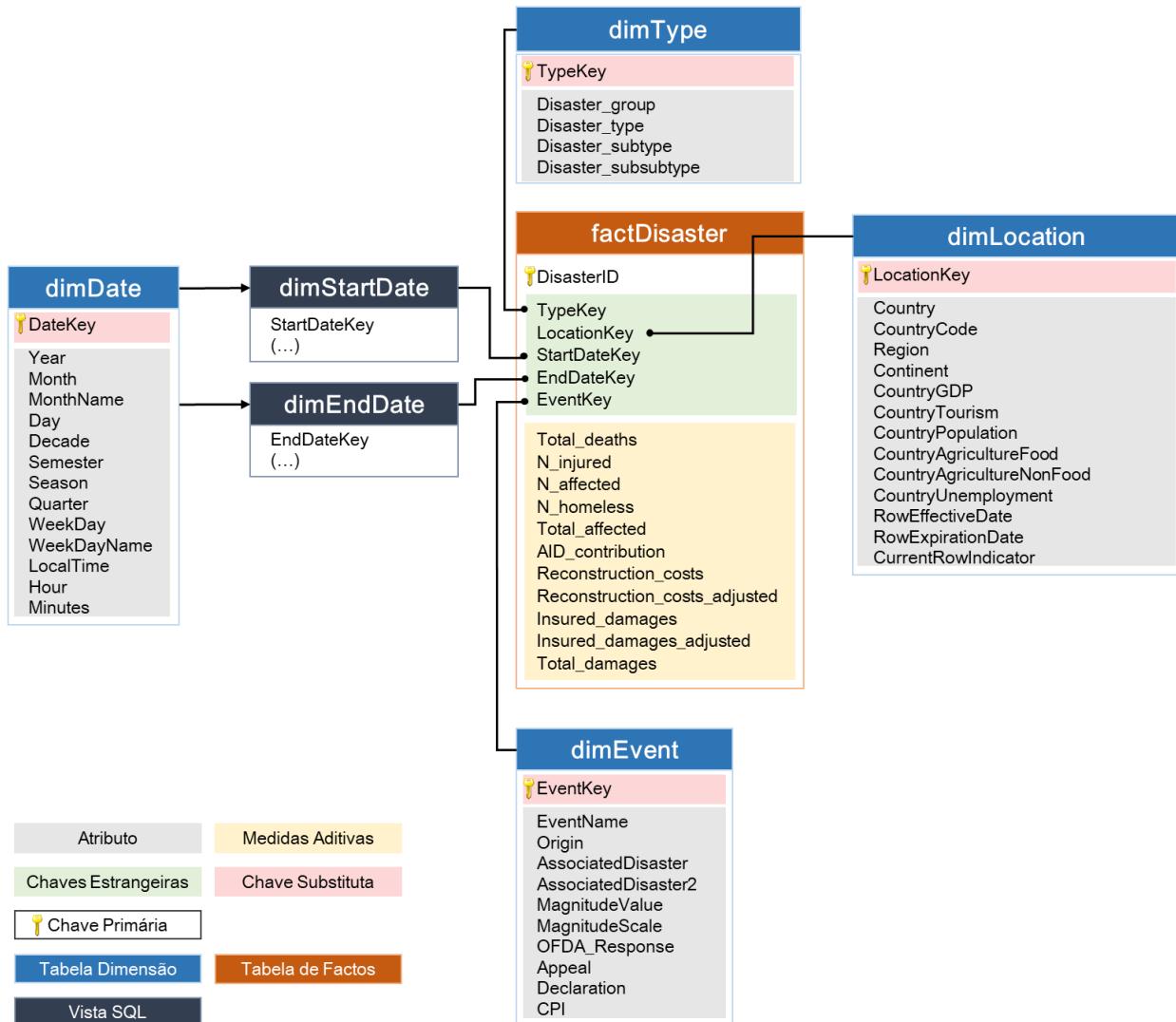


Figura 4. Diagrama em estrela para o modelo de negócio, incluindo a tabela de factos factDisaster e as respetivas dimensões.

7 Sistema ETL e Relatórios

Um sistema ETL (Extração, Transformação e Carregamento) atua como a espinha dorsal da integração de dados, permitindo a extração dos mesmos de várias fontes diferentes, a sua transformação num formato padronizado e o seu carregamento para o data warehouse. Ao automatizar estes processos essenciais, um sistema ETL facilita a integração dos dados, assegura a sua qualidade e consistência e permite a sua utilização da forma mais acertada para fins de análise, relatórios e *business intelligence*. O foco desta etapa é a implementação do sistema ETL que agrupa e otimiza as duas etapas anteriores relativas à extração e transformação, acrescentando a fase de carregamento na base de dados e análise dos dados no PowerBI.

7.1 Extração

O processo de extração é o primeiro passo do processo de ETL, onde é usado para integrar e consolidar dados de várias fontes num único local para análise. Após se procurarem e estabelecerem as fontes de dados, reúnem-se os dados para serem preparados para a fase de transformação. É um passo crítico, pois a qualidade e a integridade dos dados extraídos afetam diretamente os resultados e a eficácia de análises posteriores.

Neste projeto, a extração de dados foi realizada a partir de várias fontes relevantes, como o EM-DAT, World Bank e FAOSTAT, com o objetivo de explorar o tema dos desastres naturais. Os desastres naturais são um assunto fascinante e versátil, a fonte de dados do EM-DAT revelou-se valiosa para o nosso propósito.

Ao analisar o conteúdo dessa fonte de dados, identificamos a necessidade de entender como os desastres naturais podem afetar diferentes áreas, como agricultura, crescimento da população, turismo e desenvolvimento económico de um país. Para obter uma perspetiva abrangente, procuramos obter mais dados relacionados a essas áreas.

Neste projeto, optamos por utilizar fontes de dados que forneciam dados históricos sobre desastres naturais em todo o mundo. Essas fontes de dados permitiram que obtivéssemos uma visão abrangente e global sobre o tema. A ligação entre as várias fontes de dados foi estabelecida com base na localização geográfica, como região, continente ou país, e no período de tempo a que os dados se referiam, como a década, ano ou mês. Esta abordagem permitiu-nos agrupar e analisar os dados de acordo com critérios específicos, facilitando a comparação e a compreensão das tendências e padrões.

7.2 Transformação

Nesta segunda etapa do sistema ETL, procedemos à transformação dos dados, que envolve a sua uniformização e conformação. Aqui, fizemos a limpeza dos erros e harmonizamos os dados, para facilitar a ligação entre as várias tabelas. Para este propósito, os dados dos ficheiros foram transformados em scripts em python. A apresentação destes scripts são uma ordem sequencial do processo nesta etapa. Primeiro fez-se um script com as funções para o tratamento dos dados (script “funcoes.py”), como a remoção de colunas não relevantes, a uniformização dos nomes das tabelas, adição de novas colunas e remoção de valores nulos ou incorretos. A partir deste script, os dados das tabelas foram processados e, posteriormente, criadas as tabelas de dimensões e de factos.

7.2.1 Tratamento da Tabela DESASTRES NATURAIS

- script: transformDisasters.py

Responsabilidade:

Nesta script, o objetivo é remover colunas que não são relevantes, renomear os nomes das colunas, renomear o nome dos países e corrigir o tipo de colunas considerando a conversão para o tipo de dados corretos. Para as colunas com valores inválidos e infinitos, será feita a substituição por NaN e, posteriormente, para o tipo de inteiro nulo. Além disso, irá criar novas colunas para um maior detalhe em relação à data do evento, como o semestre, década, trimestre, estação do ano, dia da semana hora e minutos.

Input	Output
<p>Ficheiro 'DatasetDesastres.csv' com:</p> <ul style="list-style-type: none"> • Colunas não relevantes ('Glide', 'Adm Level','Admin1 Code','Admin2 Code','Geo Locations' e 'Disaster Group') • Nome das colunas não conformes com outras tabelas • Nome dos países não conformes com os nomes dos países das outras tabelas • Colunas com erros • Tabela sem colunas para a hora, minutos e estação, semestre, trimestre, década e dia da semana. • Colunas com valores nulos. 	<p>Ficheiro 'DatasetDesastres.csv' com:</p> <ul style="list-style-type: none"> • Tabela sem colunas relevantes ('Glide', 'Adm Level','Admin1 Code','Admin2 Code','Geo Locations' e 'Disaster Group') • Colunas com alteração do nome e em conformidade com as restantes tabelas. • Nomes dos países corrigidos e em uniformidade com os restantes países. • Correção do tipo de colunas para o tipo de inteiros nulos para as colunas referentes à data (ano, mês, dia). • Tabela com colunas para a hora, minutos e estação, semestre, trimestre, década e dia da semana. • Colunas sem valores nulos.

7.2.2 Tratamento da Tabela AGRICULTURA

- script: transformAgricultura.py

Responsabilidade:

Seguindo a mesma abordagem, nesta tabela, removemos colunas que não são relevantes e renomeamos os nomes das colunas para uma maior facilidade de manipulação dos dados.

Input	Output
<p>Ficheiro 'DatasetAgricultura.csv' com:</p> <ul style="list-style-type: none"> • Colunas não relevantes ('Domain Code', 'Domain', 'Element Code', 'Element', 'Year Code', 'Unit', 'Flag' e 'Flag Description') • Nome das colunas não conformes com outras tabelas • Nome dos países não conformes com os nomes dos países das outras tabelas 	<p>Ficheiro 'DatasetAgricultura.csv' com:</p> <ul style="list-style-type: none"> • Tabela sem colunas relevantes ('Domain Code', 'Domain', 'Element Code', 'Element', 'Year Code', 'Unit', 'Flag' e 'Flag Description') • Colunas com alteração do nome e em conformidade com as restantes tabelas. • Nomes dos países corrigidos e em uniformidade com os restantes países.

7.2.3 Tratamento da Tabela DESEMPREGO

- script: transformDesempregados.py

Responsabilidade:

Nesta tabela, removemos colunas que não são relevantes e renomeamos os nomes das colunas para uma maior facilidade de manipulação dos dados.

Input	Output
Ficheiro 'DatasetDesemprego.csv' com: <ul style="list-style-type: none"> ● Colunas não relevantes 'Series Name' e 'Series Code') ● Nome das colunas não conformes com outras tabelas ● Colunas com valores nulos ou incorrectos 	Ficheiro 'DatasetDesemprego.csv' com: <ul style="list-style-type: none"> ● Tabela sem colunas relevantes ('Series Name' e 'Series Code') ● Colunas com alteração do nome e em conformidade com as restantes tabelas. ● Valores nulos ou incorrectos substituídos por 'NA'

7.2.4 Tratamento da Tabela GDP

- script: transformGDP.py

Responsabilidade:

Nesta tabela, removemos colunas que não são relevantes e renomeamos os nomes das colunas para uma maior facilidade de manipulação dos dados.

Input	Output
Ficheiro 'DatasetGDP.csv' com: <ul style="list-style-type: none"> ● Colunas não relevantes 'Series Name' e 'Series Code') ● Nome das colunas não conformes com outras tabelas ● Colunas com valores nulos ou incorrectos 	Ficheiro 'DatasetGDP.csv' com: <ul style="list-style-type: none"> ● Tabela sem colunas relevantes ('Series Name' e 'Series Code') ● Colunas com alteração do nome e em conformidade com as restantes tabelas. ● Valores nulos ou incorrectos substituídos por 'NA'

7.2.5 Tratamento da Tabela POPULAÇÃO

- [script: transformPopulacao.py](#)

Responsabilidade:

Nesta tabela, removemos colunas que não são relevantes e renomeamos os nomes das colunas para uma maior facilidade de manipulação dos dados.

Input	Output
Ficheiro 'DatasetPopulação.csv' com: <ul style="list-style-type: none"> • Colunas não relevantes 'Series Name' e 'Series Code' • Nome das colunas não conformes com outras tabelas • Colunas com valores nulos ou incorrectos 	Ficheiro 'DatasetPopulação.csv' com: <ul style="list-style-type: none"> • Tabela sem colunas relevantes ('Series Name' e 'Series Code') • Colunas com alteração do nome e em conformidade com as restantes tabelas. • Valores nulos ou incorrectos substituídos por 'NA'

7.2.6 Tratamento da Tabela TURISMO

- [script: transformTurismo.py](#)

Nesta tabela, removemos colunas que não são relevantes e renomeamos os nomes das colunas para uma maior facilidade de manipulação dos dados.

Input	Output
Ficheiro 'DatasetTurismo.csv' com: <ul style="list-style-type: none"> • Colunas não relevantes ('Indicator Name' e 'Indicator Code') • Nome das colunas não conformes com outras tabelas • Nome dos países não conformes com os nomes dos países das outras tabelas 	Ficheiro 'DatasetTurismo.csv' com: <ul style="list-style-type: none"> • Tabela sem colunas relevantes ('Indicator Name' e 'Indicator Code') • Colunas com alteração do nome e em conformidade com as restantes tabelas. • Nomes dos países corrigidos e em uniformidade com os restantes países.

7.3 Construção do modelo dimensional

7.3.1 Dimensões:

- [script: dimData.py](#)

Responsabilidade:

O objectivo é uma tabela de dimensão data com uma chave substituta e colunas com as informações sobre a altura da ocorrência dos desastres que permita a correspondência entre as datas e os dados da tabela DESASTRES NATURAIS.

Input	Output
Ficheiro de dados 'df_desastres.csv'	Ficheiro de dados dimData.csv que corresponde à tabela de dimensão Data

- script: dimLocation.py

Responsabilidade:

O objetivo é criar uma tabela de dimensão Localização, gerando a chave substituta e as colunas com as informações sobre a localização. Neste script transformou-se as colunas “years” em linhas para todas as tabelas GDP, TURISMO,DESEMPREGO e POPULAÇÃO. Fez-se a junção destas tabelas (df_stats) e ordenou-se por ano e país. Em relação à tabela AGRICULTURA, estabelecemos as colunas Country_code, ‘Country’ e ‘Year’ como índices. De seguida, retiramos as colunas de produtos, ficando apenas os dados sobre localização e ano e juntámos à tabela criada anteriormente, criando a tabela “dimLocation” e adicionou-se a LocationKey. Foi criada uma linha para cada novo ano e três novos atributos: “RowEffectiveDate”, onde se detalha a data do início do ano, “RowExpirationDate”, onde se detalha a data do fim do ano e “CurrentRowIndicator”, onde se torna explícito se o ano é o atual ou não (através da categoria “Current” ou “Expired”).

Input	Output
Ficheiro de dados ‘df_desastres.csv’	Ficheiro de dados dimLocation.csv que corresponde à tabela de dimensão Localização
Ficheiro de dados ‘df_agricultura.csv’	
Ficheiro de dados ‘df_populacao.csv’	
Ficheiro de dados ‘df_turismo.csv’	
Ficheiro de dados ‘df_desemprego.csv’	
Ficheiro de dados ‘df_gdp.csv’	

- script: dimEvent.py

Responsabilidade:

Neste script, o objectivo é criar a tabela dimEvent com uma chave substituta e colunas com informação sobre as ocorrências dos desastres naturais, que permita a correspondência entre os eventos e os dados da tabela DESASTRES NATURAIS.

Input	Output
Ficheiro de dados ‘df_desastres.csv’	Ficheiro de dados dimEvent.csv que corresponde à tabela de dimensão Evento

- script: dimType.py

Responsabilidade:

Neste script, o objectivo é criar a tabela dimType com uma chave substituta e colunas com informação sobre o tipo de desastres naturais, que permita a correspondência entre os tipos de desastres e os dados da tabela DESASTRES NATURAIS.

Input	Output
Ficheiro de dados ‘df_desastres.csv’	Ficheiro de dados dimType.csv que corresponde à tabela de dimensão Evento

7.3.2 Tabela de Factos:

- script: factDesaster.py

Responsabilidade:

Neste script, o objetivo é criar a tabela de factos, com as colunas para armazenar as chaves referentes a cada tabela de dimensões e as medidas do processo de negócio.

Input	Output
Ficheiro de dados 'df_desastres.csv'	Ficheiro de dados factDisaster.csv que corresponde à tabela de factos

7.4 Carregamento

O Carregamento dos dados extraídos e transformados corresponde à terceira etapa, para a data presentation area. Nesta etapa os dados preparados são carregados numa base de dados para análise e geração de relatórios. Através de um script em python, estabeleceu-se a conexão com a base de dados PostgreSQL para o carregamento dos dados.

7.4.1 Dimensões:

7.4.1.1 Carregamento da dimensão data

- script: loadData.py

Responsabilidade:

O objetivo principal deste script é criar a tabela dimensão "dimData" na base de dados e preenchê-la com os valores que já foram previamente tratados e armazenados na tabela "dimData.csv".

A tabela "dimData" será criada se ainda não existir na base de dados. Em seguida, os valores que passaram por todas as transformações necessárias serão copiados da tabela "dimData.csv" e inseridos na tabela "dimData".

Devido ao facto de termos optado por criar 2 vistas SQL para fazerem a relação com a tabela de factos, neste script também foram criadas estas 2 vistas, uma relativa às "Start Dates" e outras às "EndDates".

Input	Output
Ficheiro de dados 'dimData.csv'	Criação da tabela dimensão 'dimData', e views startDate e endDate, que são adicionadas à base de dados. A nova tabela é preenchida com os com os seus respectivos valores.

7.4.1.2 Carregamento da dimensão evento

- script: loadEvent.py

Responsabilidade:

O objetivo principal deste script é criar a tabela dimensão "dimEvento" na base de dados e preenchê-la com os valores que já foram previamente tratados e armazenados na tabela "dimEvento.csv".

A tabela "dimEvento" será criada se ainda não existir na base de dados. Em seguida, os valores que passaram por todas as transformações necessárias serão copiados da tabela "dimEvento.csv" e inseridos na tabela "dimEvento".

Input	Output
Ficheiro de dados 'dimEvento.csv'	Criação da tabela dimensão 'dimEvento' que é adicionada à base de dados. A nova tabela é preenchida com os seus respetivos valores.

7.4.1.3 Carregamento da dimensão evento

- script: loadLocation.py

Responsabilidade:

O objetivo principal deste script é criar a tabela dimensão "dimLocation" na base de dados e preenchê-la com os valores que já foram previamente tratados e armazenados na tabela "dimLocation.csv".

A tabela "dimLocation" será criada se ainda não existir na base de dados. Em seguida, os valores que passaram por todas as transformações necessárias serão copiados da tabela "dimLocation.csv" e inseridos na tabela "dimEvento".

Input	Output
Ficheiro de dados 'dimLocation.csv'	Criação da tabela dimensão 'dimLocation' que é adicionada à base de dados. A nova tabela é preenchida com os seus respetivos valores.

7.4.1.4 Carregamento da dimensão evento

- script: loadType.py

Responsabilidade:

O objetivo principal deste script é criar a tabela dimensão "dimType" na base de dados e preenchê-la com os valores que já foram previamente tratados e armazenados na tabela "dimLocation.csv".

A tabela "dimType" será criada se ainda não existir na base de dados. Em seguida, os valores que passaram por todas as transformações necessárias serão copiados da tabela "dimType.csv" e inseridos na tabela "dimEvento".

Input	Output
Ficheiro de dados 'dimType.csv'	Criação da tabela dimensão 'dimType' que é adicionada à base de dados. A nova tabela é preenchida com os seus respetivos valores.

7.4.2 Factos:

7.4.2.1 Carregamento da facto desastres

- [script:loadDisaster.py](#)

Responsabilidade:

O objetivo principal deste script é criar a tabela de fatos "factDisaster" na base de dados e preenchê-la com os valores que foram previamente tratados e armazenados na tabela "factDisaster.csv". É fundamental destacar que, devido às dependências existentes entre a tabela de fatos e as tabelas dimensionais, a criação da tabela de fatos deve ser realizada por último, a fim de garantir que as ligações entre as tabelas sejam estabelecidas corretamente.

Com a execução deste script, a base de dados estará completa e pronta para ser analisada e explorada por meio de perguntas analíticas relacionadas ao tema em questão.

Input	Output
Ficheiro de dados 'factDisaster.csv'	Criação e preenchimento da tabela de factos 'factDisaster' com os seus respetivos valores.

7.5 Execução dos scripts

O script "run_scripts_windows.bat" foi criado com o objetivo de facilitar a execução de todo o sistema ETL (Extração, Transformação e Carregamento), englobando todos os scripts das várias etapas pela sua devida ordem. A sua principal função é executar a transformação e carregamento dos dados, permitindo que essa tarefa seja realizada com apenas a execução de um único script.

8 Diagrama de fluxo de dados

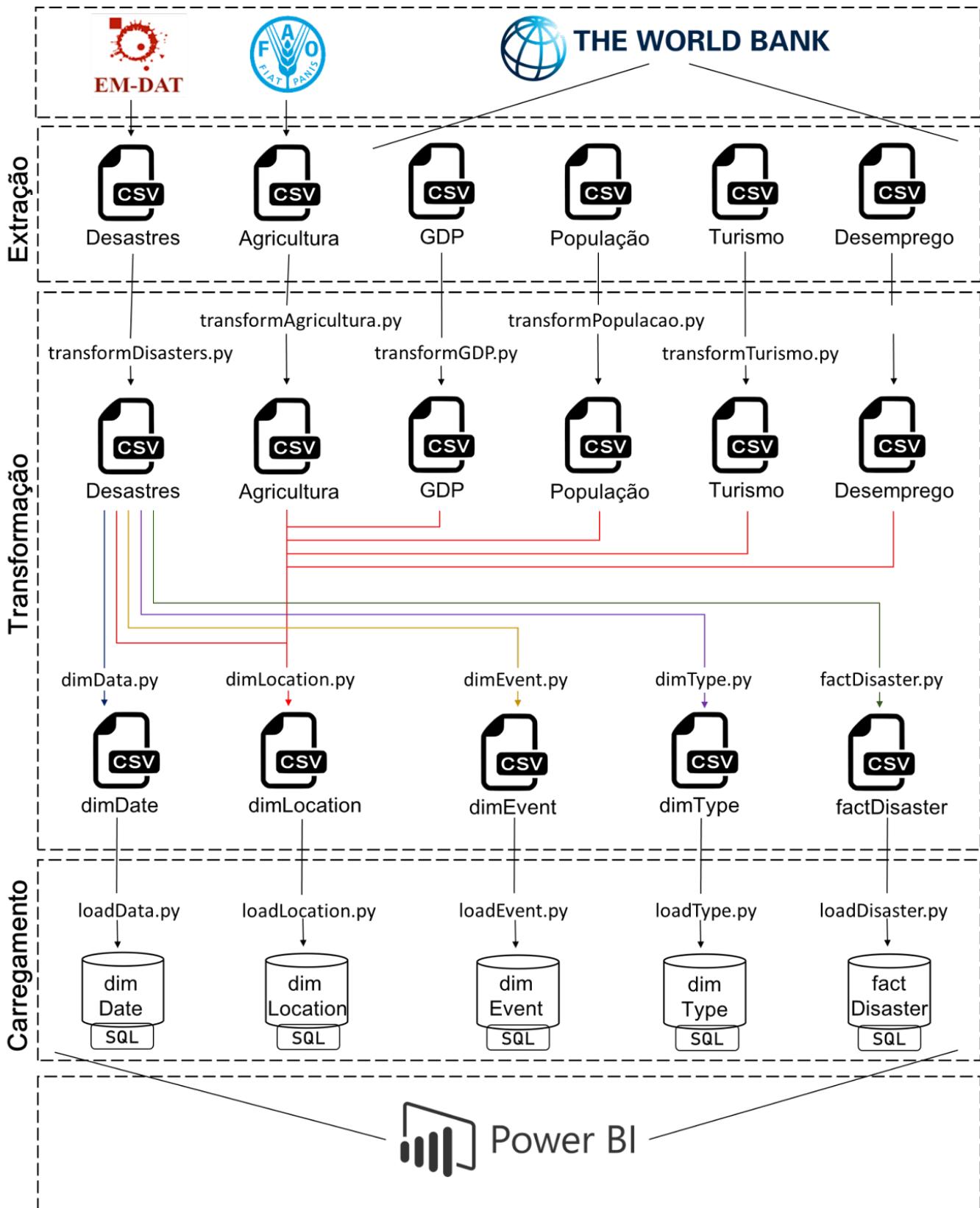


Figura 5 - Diagrama de Fluxo de Dados

9 Cubo de dados

De forma a gerar e visualizar o respetivo cubo de dados, foi utilizada a ferramenta do PowerBI. Para isso conectou-se a este a Base de Dados no PostgreSQL, carregando todas as dimensões e tabela de factos. Foi ainda necessário assegurar que o roleplaying era garantido, criando relações entre as vistas da dimensão Data com a tabela de factos, pelo que não tinha ainda sido feito anteriormente.

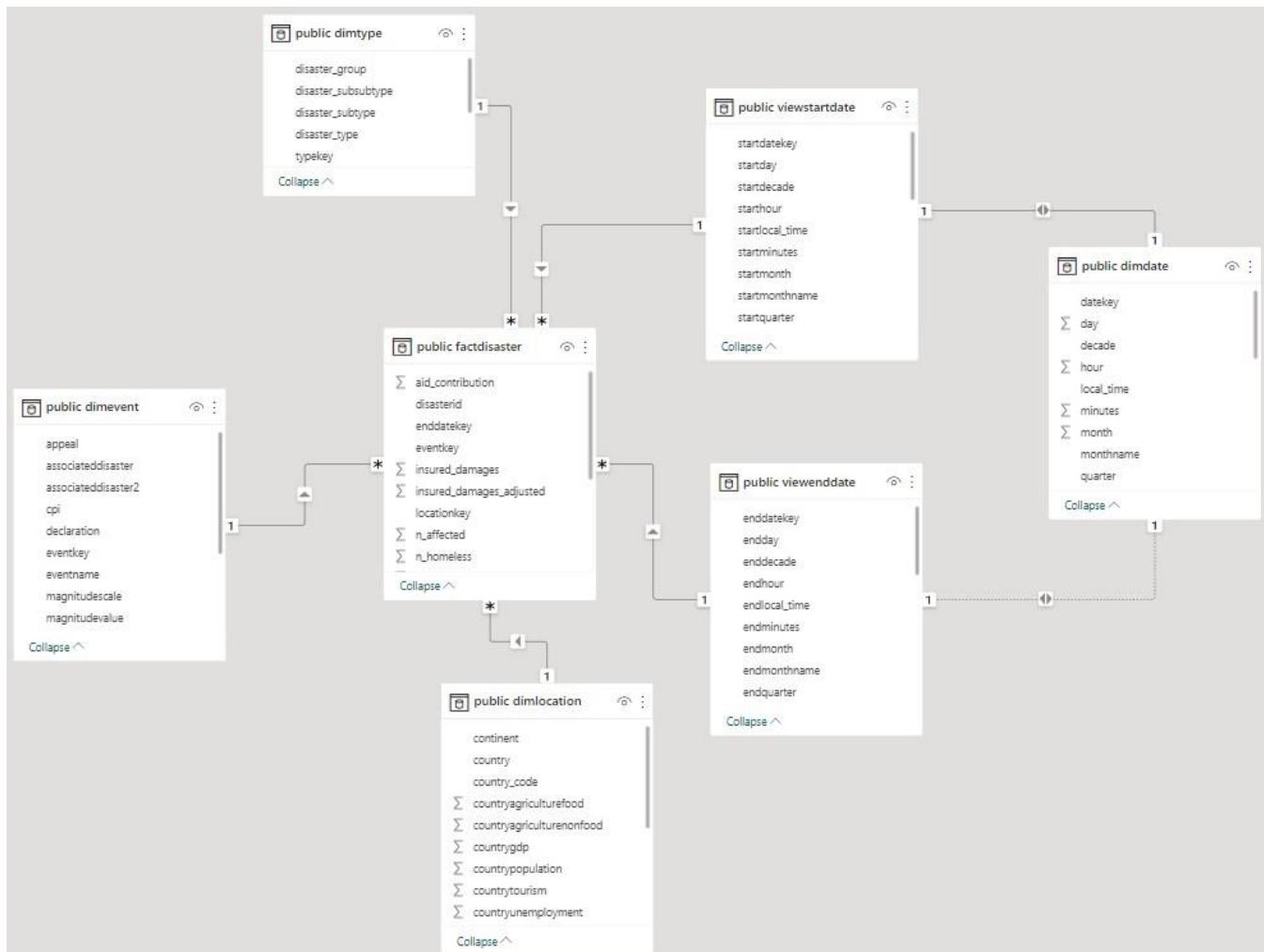


Figura 6 – Cubo de dados

10 Respostas às questões analíticas

O software escolhido para a produção e relatórios foi o PowerBI. Este é uma plataforma de business intelligence da Microsoft que oferece recursos avançados de visualização e análise de dados, permitindo a criação de relatórios interativos e painéis de controlo dinâmicos. É uma ferramenta popular para a produção de relatórios e apresentação de dados de forma intuitiva e visualmente atraente.

10.1 Primeira Pergunta analítica

“Como as catástrofes naturais afetam a produção agrícola num determinado país e qual é o impacto na sua produção? Que tipo de desastre naturais são mais propensos a afetar a agricultura e quais os tipos de matérias-primas mais afetadas?”

Primeiro, começamos por analisar quais as 10 ocorrências de desastres naturais que tiverem mais danos. Para os primeiros 4 desses países, fomos analisar um a um, se a produção de agricultura foi afetada.

Tabela 20. Tabela com as ocorrências de desastres naturais com mais danos

País	Tipo de Desastre	Ano	Mês	Danos totais
Japão	Sismo	2011	Março	210000000
Estados Unidos	Tempestade	2005	Agosto	125000000
Japão	Sismo	1995	Janeiro	100000000
Estados Unidos	Tempestade	2022	Setembro	100000000
Estados Unidos	Tempestade	2017	Agosto	95000000
China	Sismo	2008	Maio	85000000
Porto Rico	Tempestade	2017	Setembro	68000000
Estados Unidos	Tempestade	2021	Agosto	65000000
Estados Unidos	Tempestade	2017	Setembro	57000000
Estados Unidos	Tempestade	2012	Outubro	50000000
Tailândia	Cheias	2011	Agosto	40000000

Japão:

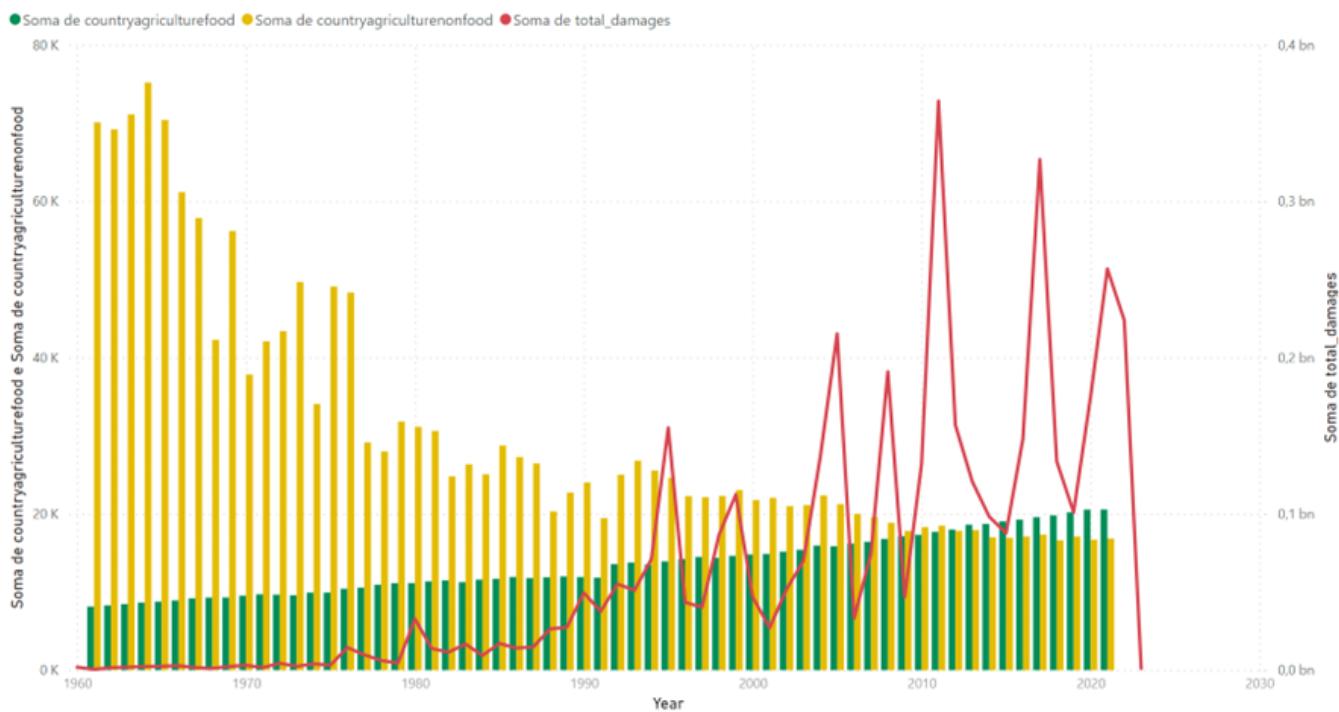


Figura 7. Produção agrícola vs danos totais no Japão.

Estados Unidos da América:

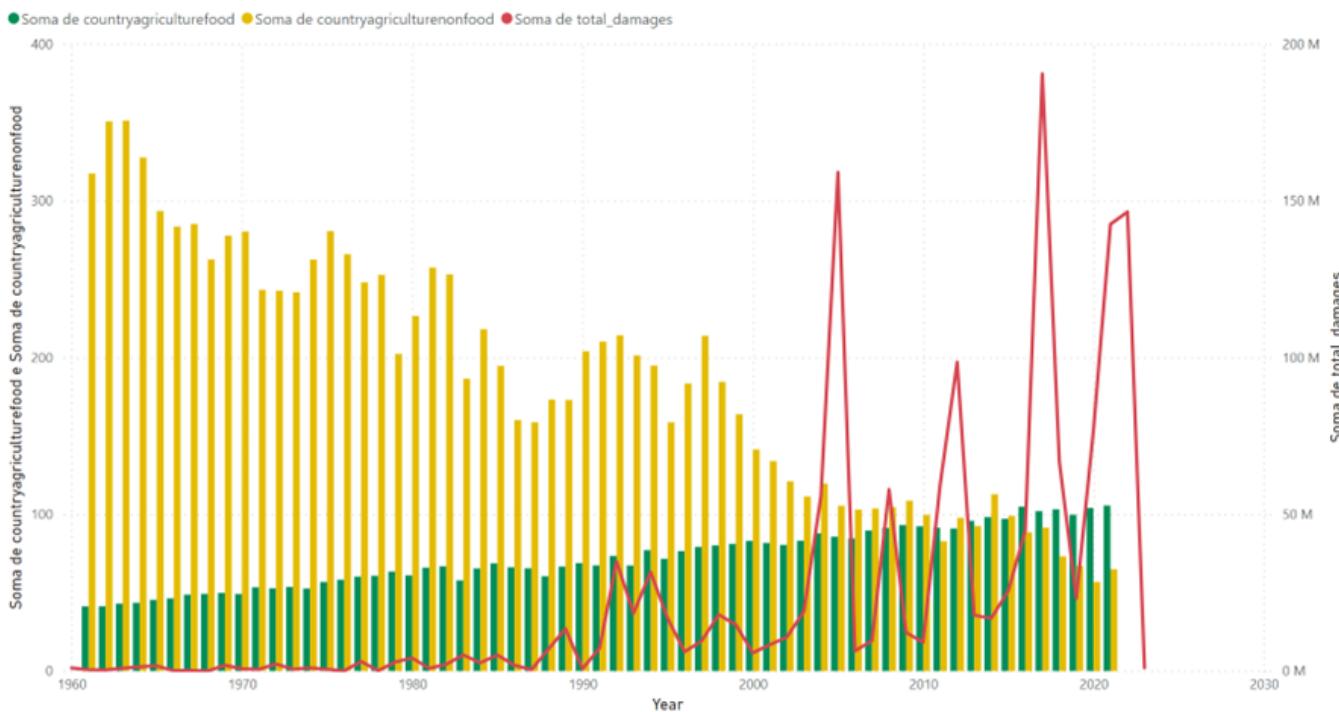


Figura 8. Produção agrícola vs danos totais nos Estados Unidos da América.

China:



Figura 9. Produção agrícola vs danos totais na China

Tailândia:

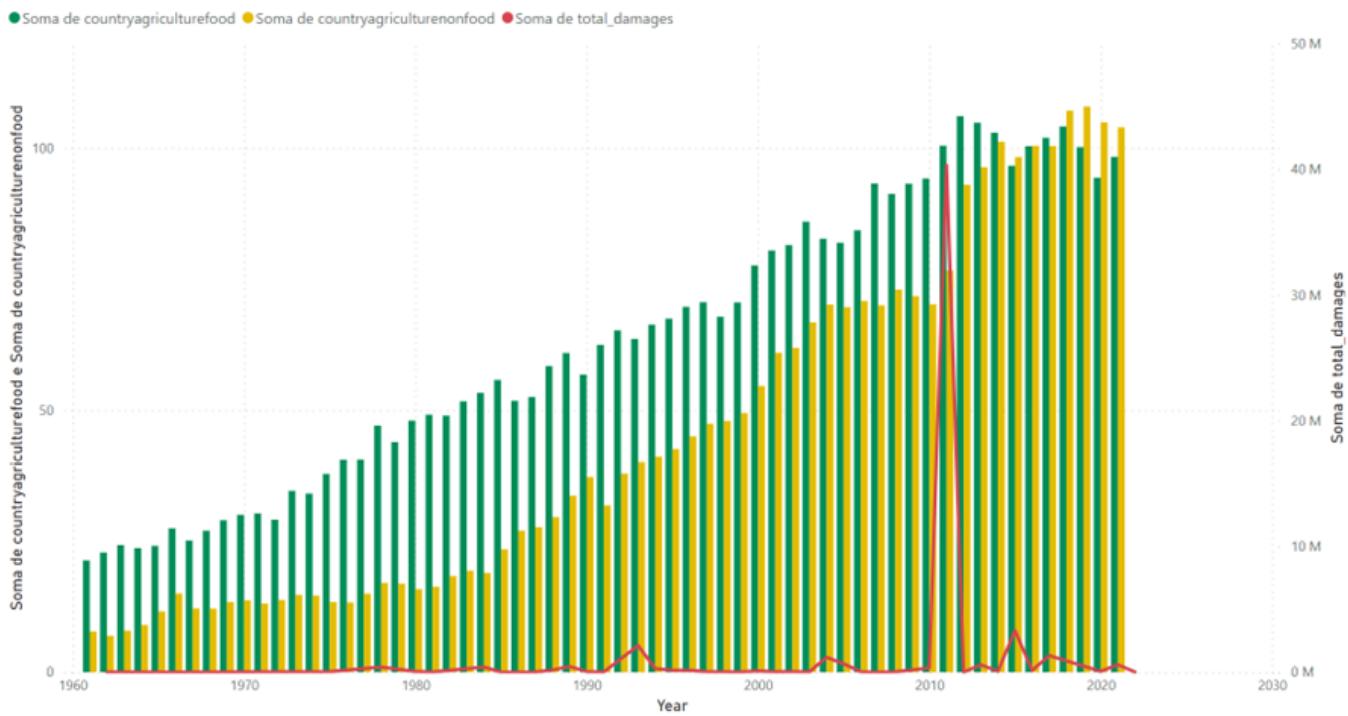


Figura 10. Produção agrícola vs danos totais na Tailândia

No Japão e nos Estados Unidos, podemos ver uma tendência decrescente da produção de produtos não alimentares e um aumento de danos originados pelos desastres. Enquanto a produção de produtos alimentares vai subindo ligeiramente de ano para ano, em ambos os países.

Na China e na Tailândia, podemos afirmar que a produção agrícola de produtos alimentares e não-alimentares cresce exponencialmente e os danos também são mais acentuados nas últimas duas décadas, tendo maior expressão na China.

No Japão e Estados Unidos, a ocorrência de desastres naturais pode ter impacto na produção de produtos não alimentares, pois estes podem levar à interrupção da produção, destruição de fábricas, perda de equipamentos e incapacidade de atender à demanda. No entanto, com esta análise não conseguimos encontrar uma relação inequívoca entre a produção Agrícola e a ocorrência de desastres naturais. Tanto no Japão como nos Estados Unidos, a diminuição da produção de produtos não alimentares pode ser devido a:

- Mudança de estratégia da economia, direcionando o investimento para outros sectores diferentes de produtos não alimentares que podem ter origem nas exigências do mercado, tecnologia ou políticas governamentais;
- Importação: à medida que a globalização avança, é possível que os países tenham optado por importar produtos não alimentares de outras regiões, em vez de produzi-los internamente. Isso pode ser influenciado por fatores como custo, eficiência ou acesso a recursos específicos.
- Mudanças no setor industrial: Alterações no setor industrial, como avanços tecnológicos, automação e aumento da eficiência produtiva, podem ter impacto na produção *per capita* de produtos não alimentares. Essas mudanças podem resultar em redução da necessidade de mão-de-obra ou em uma produção mais concentrada em empresas especializadas.
- Fatores ambientais: Algumas indústrias de produtos não alimentares podem ter impactos ambientais significativos. Como resultado, os governos podem ter implementado políticas ambientais mais rigorosas, levando a uma diminuição da produção de certos produtos ou à adoção de tecnologias mais limpas.

Após esta análise inicial, formos avaliar os países com mais mortes, onde os desastres tiveram um impacto mais negativo no estado de saúde das pessoas e avaliar a tendência da produção agrícola.

Tabela 21. Tabela com as ocorrências de desastres naturais com maior nº de mortes

País	Tipo de Desastre	Ano	Mês	Nº de Mortes
Índia	Seca	1965	Desconhecido	1500000
Bangladesh	Tempestade	1970	Novembro	300000
Etiópia	Seca	1983	Maio	300000
China	Sismo	1976	Julho	242000
Haiti	Sismo	2010	Janeiro	222570
Indonésia	Sismo	2004	Dezembro	165708
Sudão	Seca	1983	Abril	150000
Bangladesh	Tempestade	1991	Abril	138866
Myanmar	Tempestade	2008	Maio	138366
Etiópia	Seca	1973	Dezembro	100000
Moçambique	Seca	1981	Desconhecido	100000

Fomos avaliar a produção agrícola na Índia, Bangladesh, Haiti e China, que são os 4 países onde houve mais mortes devido a uma ocorrência de desastre natural. Em baixo, encontram-se os gráficos, onde avaliamos o nº de mortes (linha vermelha) vs o a produção agrícola não alimentar (cor amarela) e alimentar (cor verde) ao longo dos anos.

Índia:

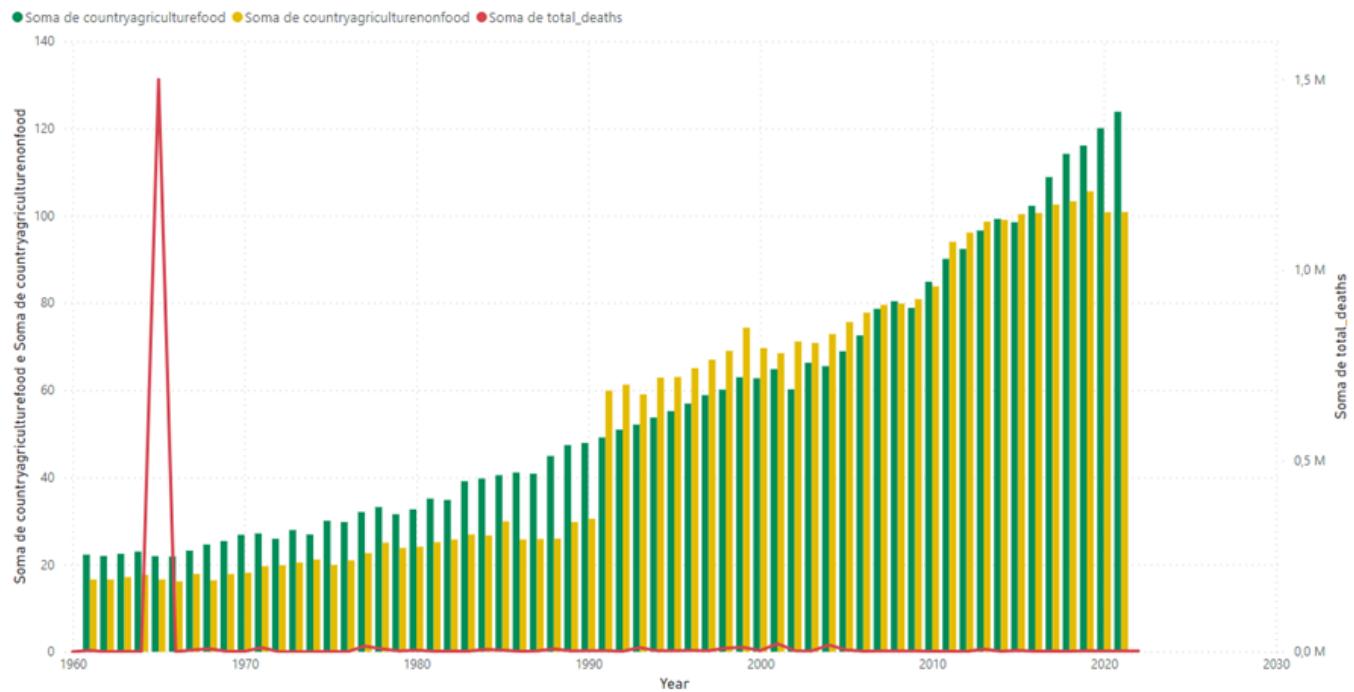


Figura 11. Produção agrícola vs nº de mortes na Índia

Bangladesh:

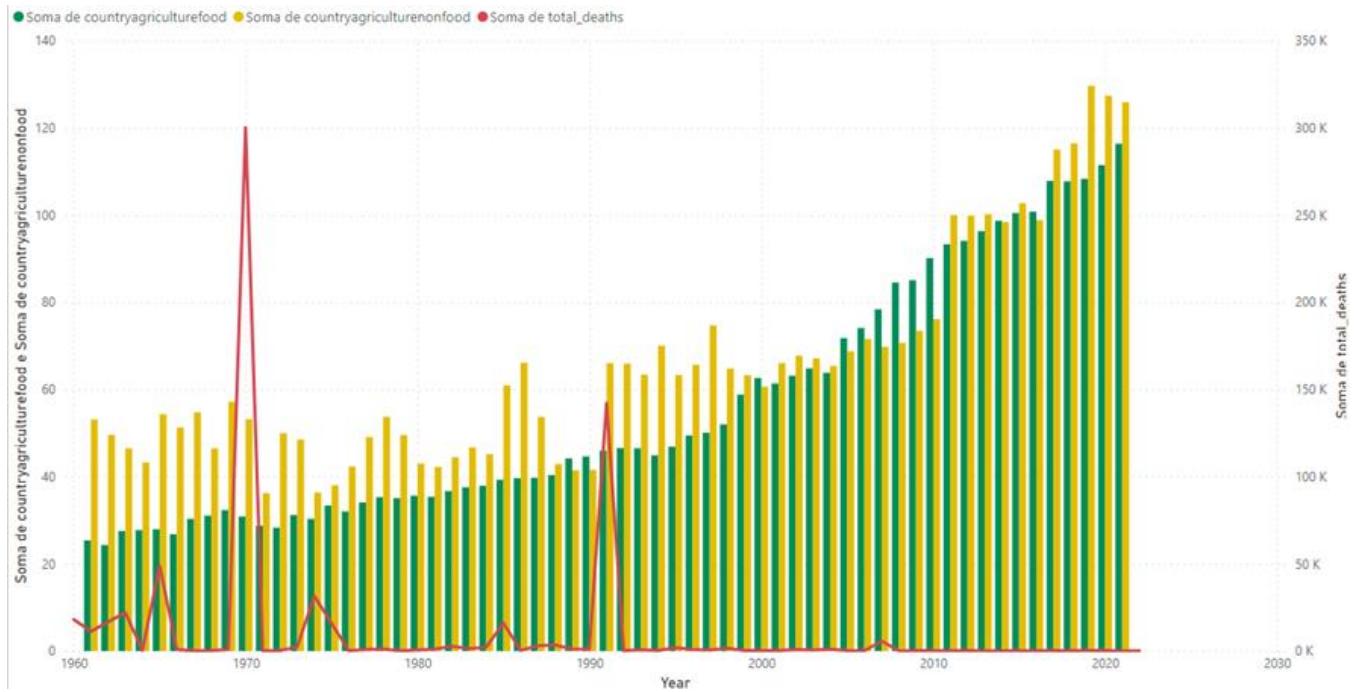


Figura 12. Produção agrícola vs nº de mortes no Bangladesh

Haiti:

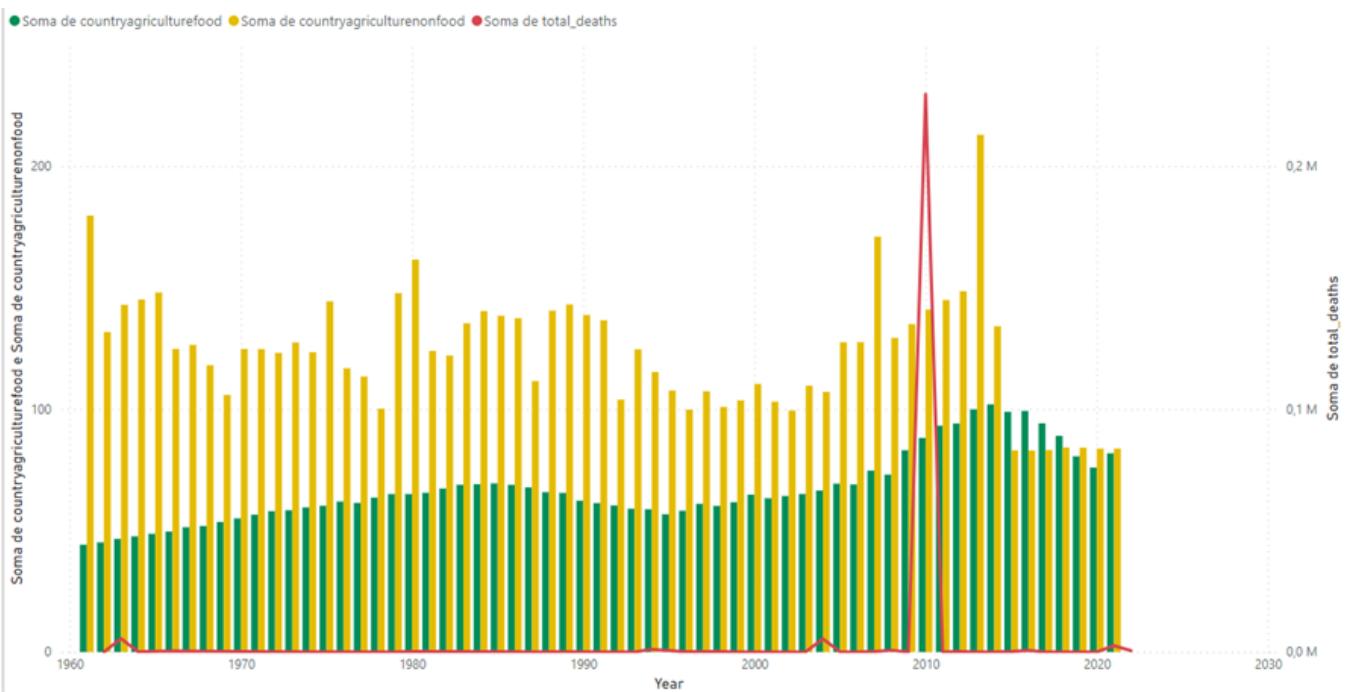


Figura 13. Produção agrícola vs nº de mortes no Haiti

China:

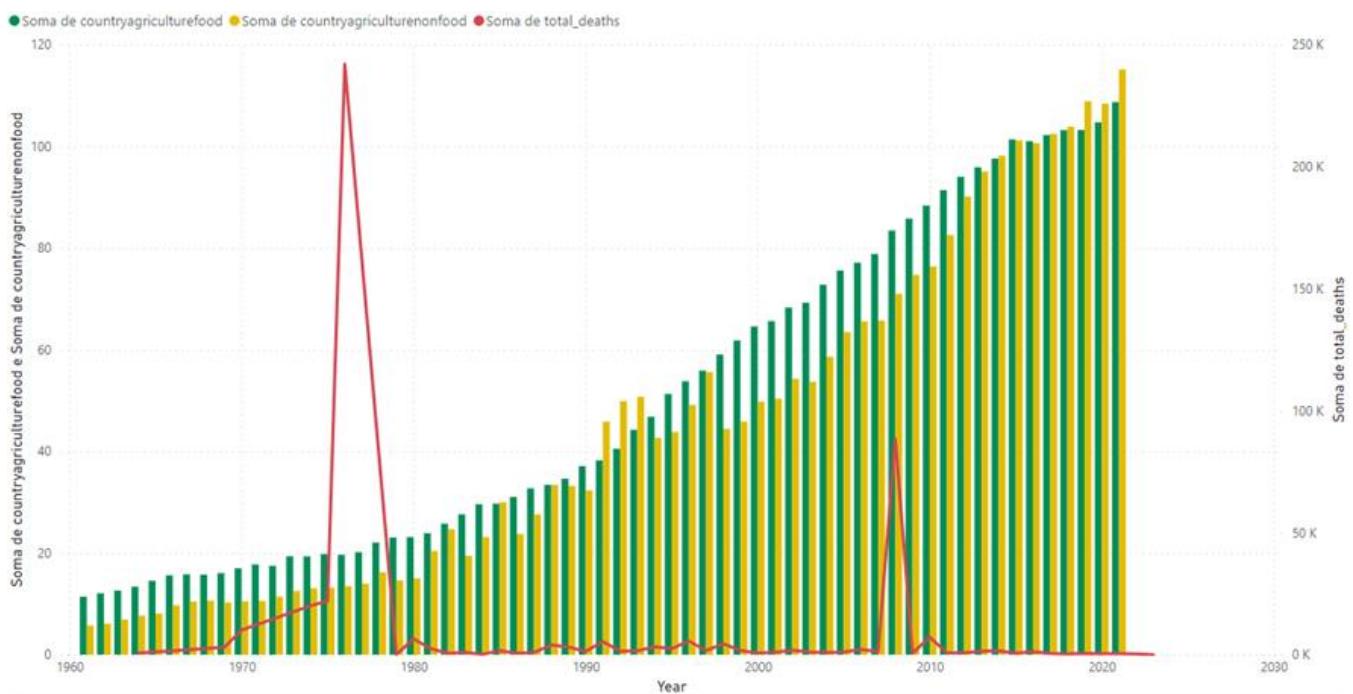


Figura 14. Produção agrícola vs nº de mortes na China

À semelhança da China, a Índia tem um crescimento de produção agrícola de ano para ano. Enquanto no Bangladesh e Haiti, a produção de produtos não alimentares varia ao longo do tempo, tendo decrescido e estabilizado nos últimos anos no Haiti. Em relação à produção dos produtos alimentares, este tem crescido no Bangladesh, mas decrescido no Haiti. Dos gráficos acima, não conseguimos encontrar nenhuma relação entre os desastres naturais com mais impacto na vida humana e na produção agrícola.

De seguida, fomos analisar os desastres naturais que podem influenciar a produção agrícola. Na figura 15, expomos os gráficos dos desastres mais frequentes para ver os que podem afetar a produção agrícola. A linha vermelha corresponde ao nº de desastres que ocorreram e as colunas verde e amarela correspondem à produção agrícola alimentar e não alimentar, respetivamente.

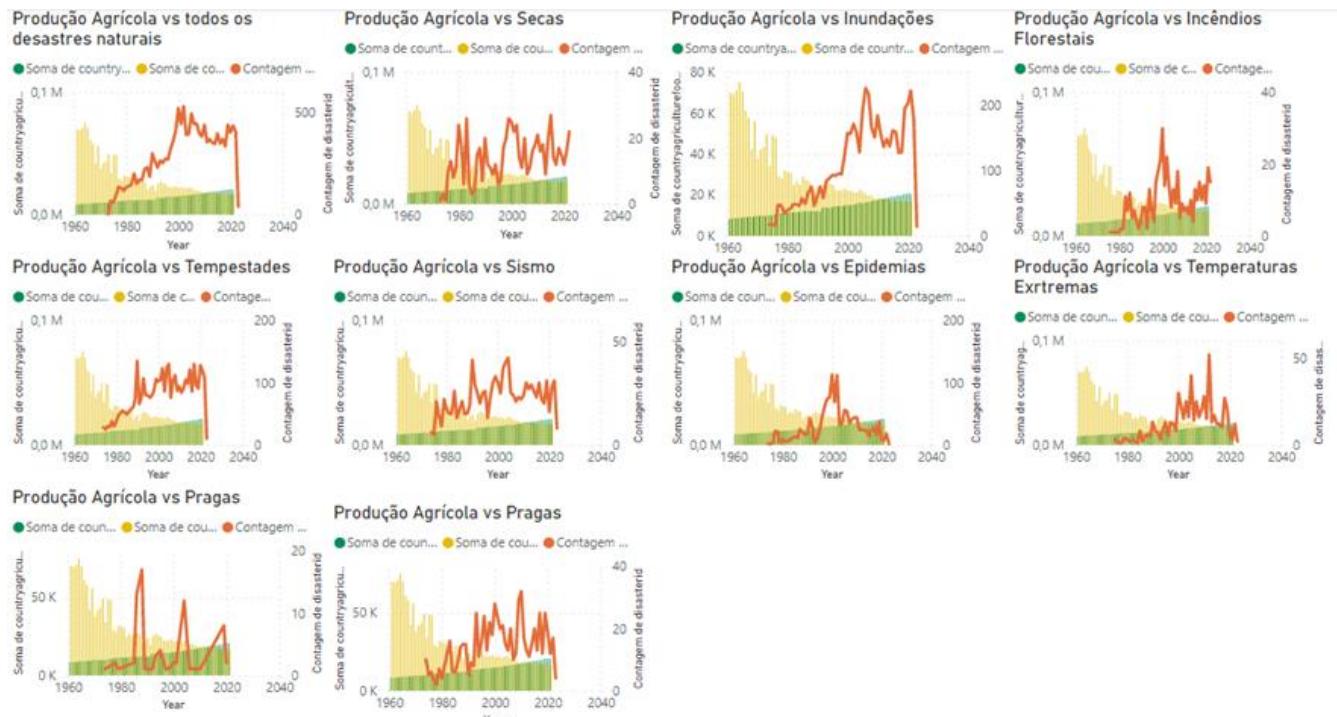


Figura 15. Produção agrícola vs os tipos de desastre naturais

De acordo com a Figura 15, podemos observar que a partir do ano 2000 existe um grande número de desastres que ocorreram e que nesse período de tempo afeta a produção agrícola não alimentar. Ao observar cuidadosamente os gráficos, podemos ver que os desastres naturais que mais podem influenciar a produção agrícola não alimentar são as secas, as inundações, as tempestades e os incêndios florestais, pois à medida que este acontece, a produção baixa no ano e/ou no ano a seguir.

Secas:

A falta de chuvas adequadas pode resultar em escassez de água para irrigação e crescimento das plantas, resultando em redução da produtividade agrícola.

Fomos consultar quais os 10 países onde ocorrem mais inundações e ver se a produção agrícola é afetada por este tipo de desastre para os 4 países mais afetados.

Tabela 22. Tabela com os países onde ocorreram mais cheias

País	Nº de Ocorrências de Secas
China	39
Brasil	21
Etiópia	18
Quénia	18
Somália	17
Moçambique	17
Índia	15
Bolívia	14
Burkina Faso	13
Mauritânia	13
Nigéria	13

Em baixo apresentamos os gráficos em que avaliamos a produção agrícola alimentar (cor verde) e não alimentar (cor amarela) vs o nº de secas (linha laranja), para os 4 países onde ocorreram mais secas.

China - Secas vs Produção agrícola:

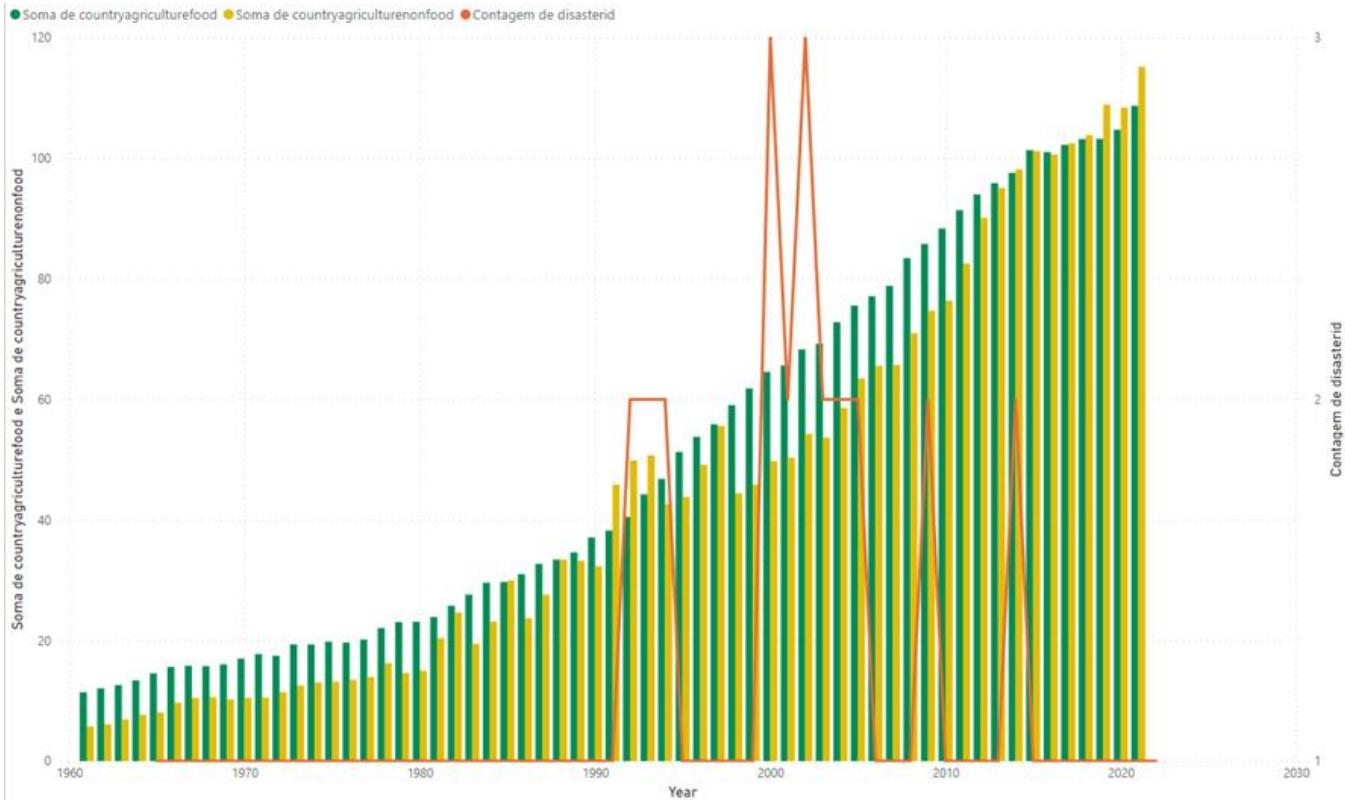


Figura 16. Produção agrícola vs Secas na China

Brasil - Secas vs Produção agrícola:

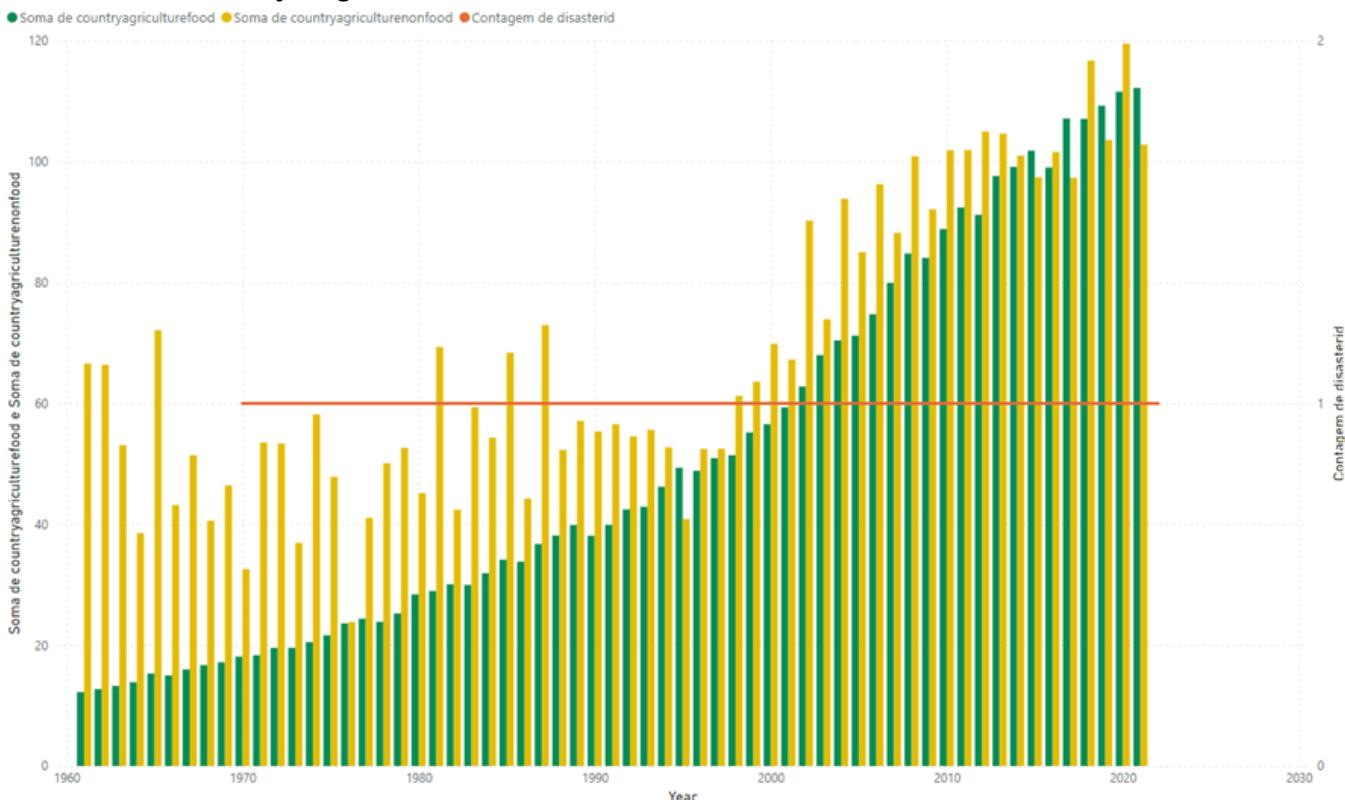


Figura 17. Produção agrícola vs Secas no Brasil

Etiópia - Secas vs Produção agrícola:

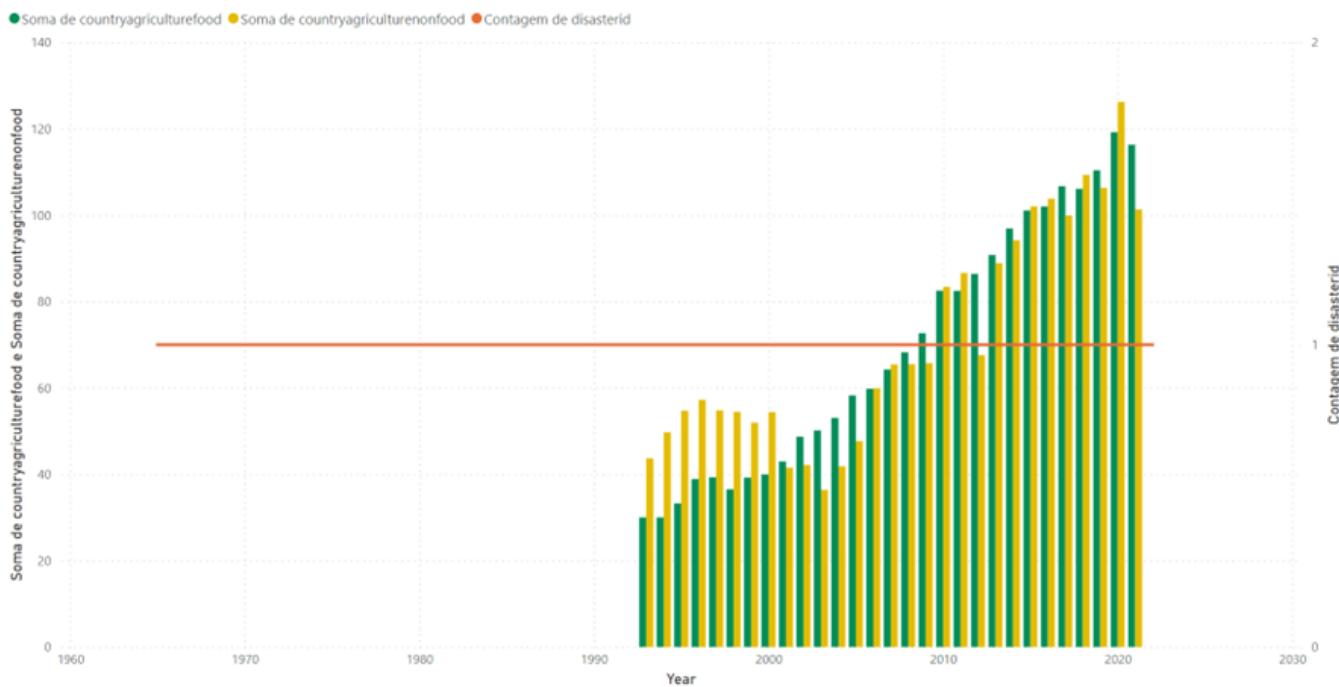


Figura 18. Produção agrícola vs Secas na Etiópia

Estados Unidos: Secas vs Produção agrícola:

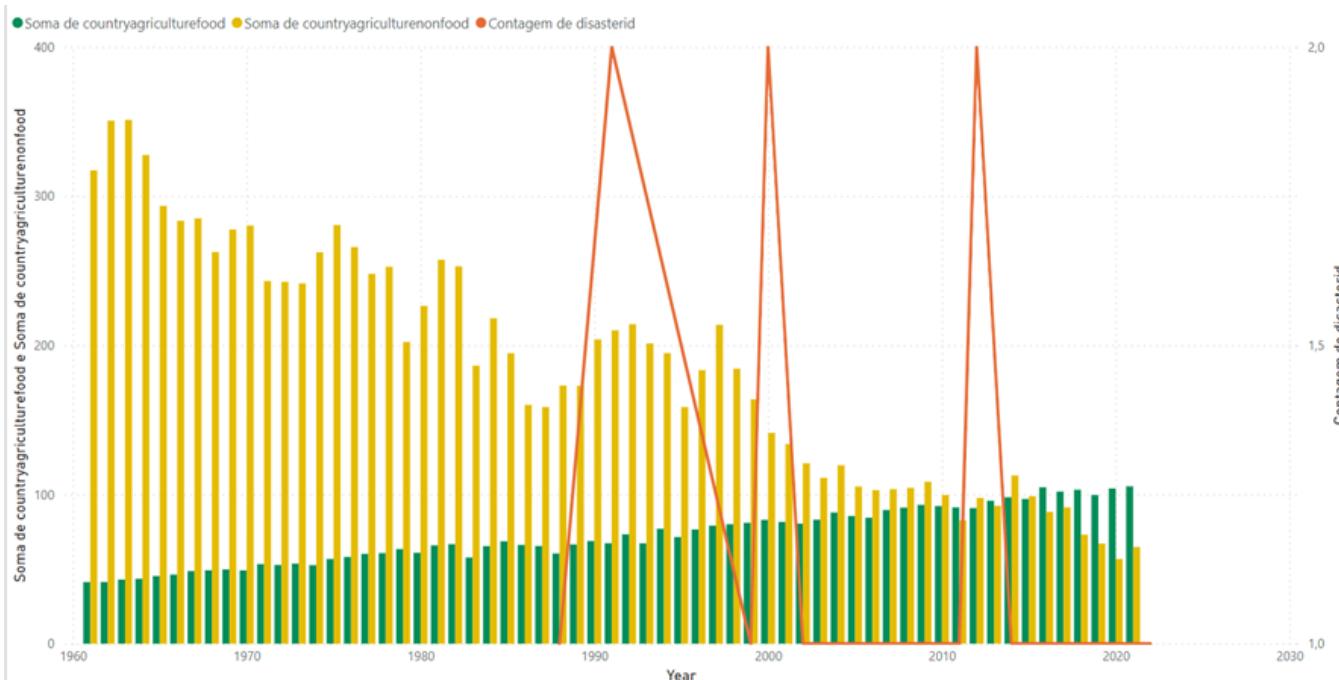


Figura 19. Produção agrícola vs Secas nos Estados Unidos

Produção agrícola vs Secas a nível global:

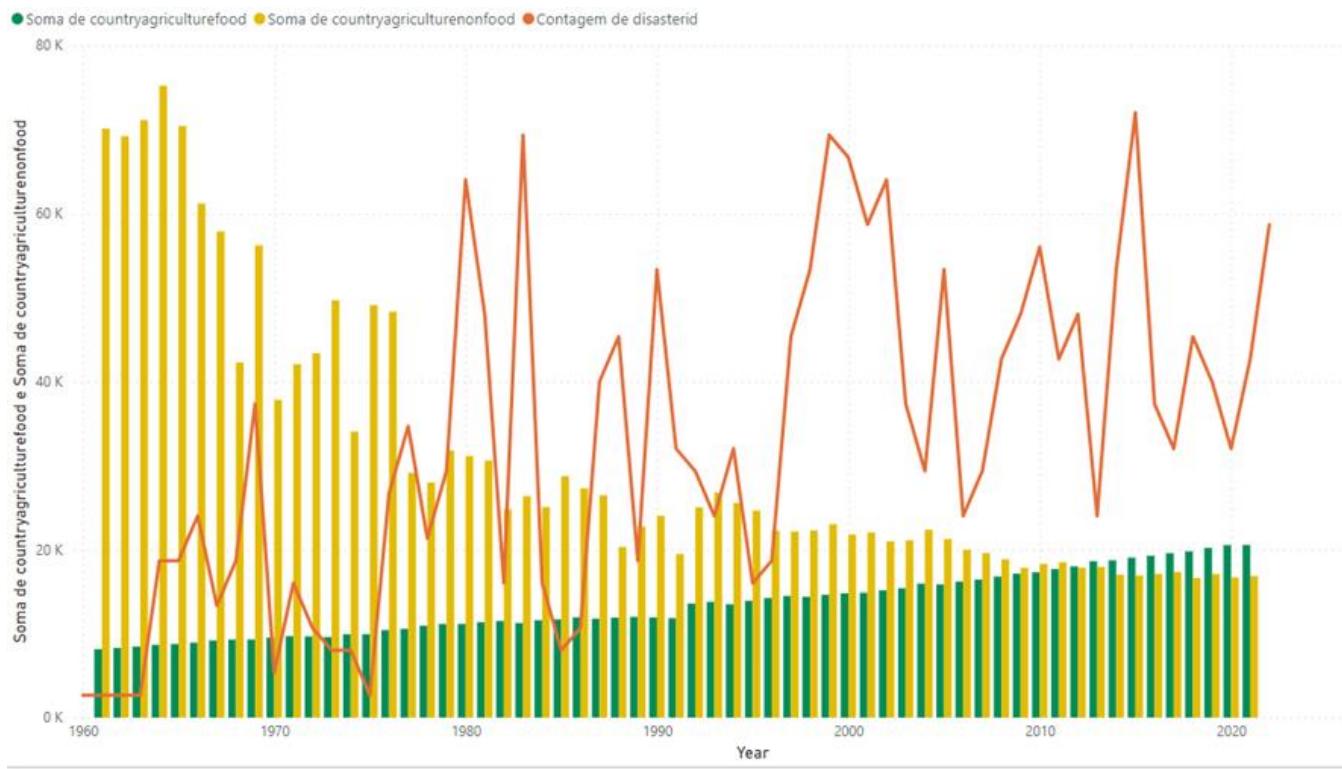


Figura 20. Produção agrícola vs Secas no Mundo.

Embora as secas possam ter impacto na produção agrícola nos países acima, não conseguimos tirar alguma relação entre a ocorrência destas e a produção agrícola. Para se poder estudar esta questão mais afundo em cada país, teríamos de ver com mais pormenor os vários produtos que podem ser afetados. Neste projeto, apenas temos os produtos divididos em produtos alimentares e não alimentares. Para conseguirmos ter alguma conclusão sobre a influência das secas, teríamos de ver dentro destas duas categorias, o histórico da produção de cada produto individualmente, pois estes podem ser afetados de forma diferente por este tipo de ocorrência.

No entanto, na Figura 20, podemos ver que a nível global, o nº de secas aumentou e houve um decréscimo da produção não alimentar. As secas podem ter influência neste tipo de produto, no entanto, com o avanço tecnológico, este tipo de produção agrícola pode ter decrescido.

Inundações:

Inundações repentinas podem destruir plantações e causar a perda de animais. O excesso de água também pode afetar a qualidade do solo e retardar o plantio ou a colheita.

Fomos consultar quais os 10 países onde ocorrem mais inundações e ver se a produção agrícola é afetada por este tipo de desastre para os 4 países mais afetados.

Tabela 23. Tabela com os países onde ocorreram mais inundações

País	Nº de Ocorrências de Inundações
Índia	301
China	297
Indonésia	264
Estados Unidos	194
Brasil	165
Filipinas	161
Colômbia	108
Paquistão	105
Vietname	102
Afeganistão	101
Bangladesh	98

Em baixo apresentamos os gráficos em que avaliamos a produção agrícola alimentar (cor verde) e não alimentar (cor amarela) vs o nº de inundações (linha laranja), para os 4 países onde ocorreram mais inundações.

Índia - Inundações vs Produção agrícola:

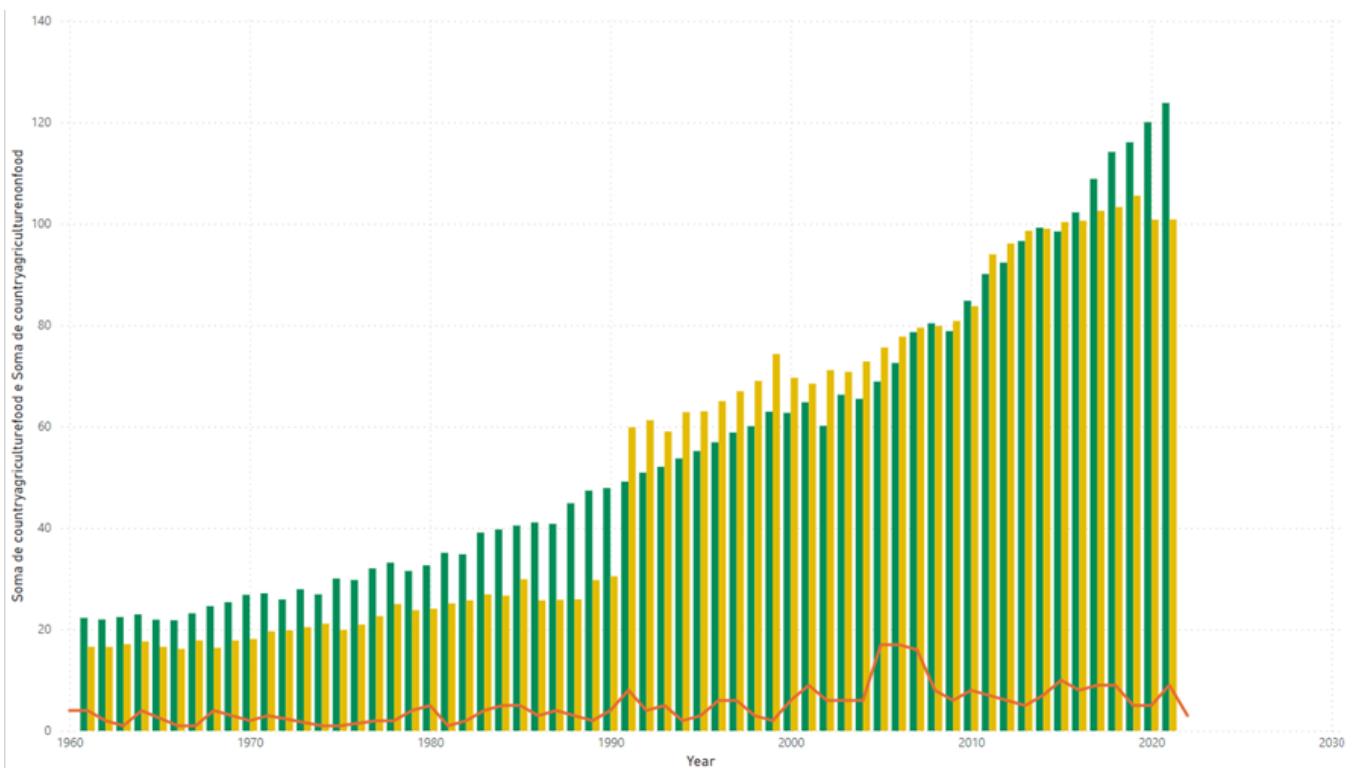


Figura 21. Produção agrícola vs Inundações na Índia.

China - Inundações vs Produção agrícola:

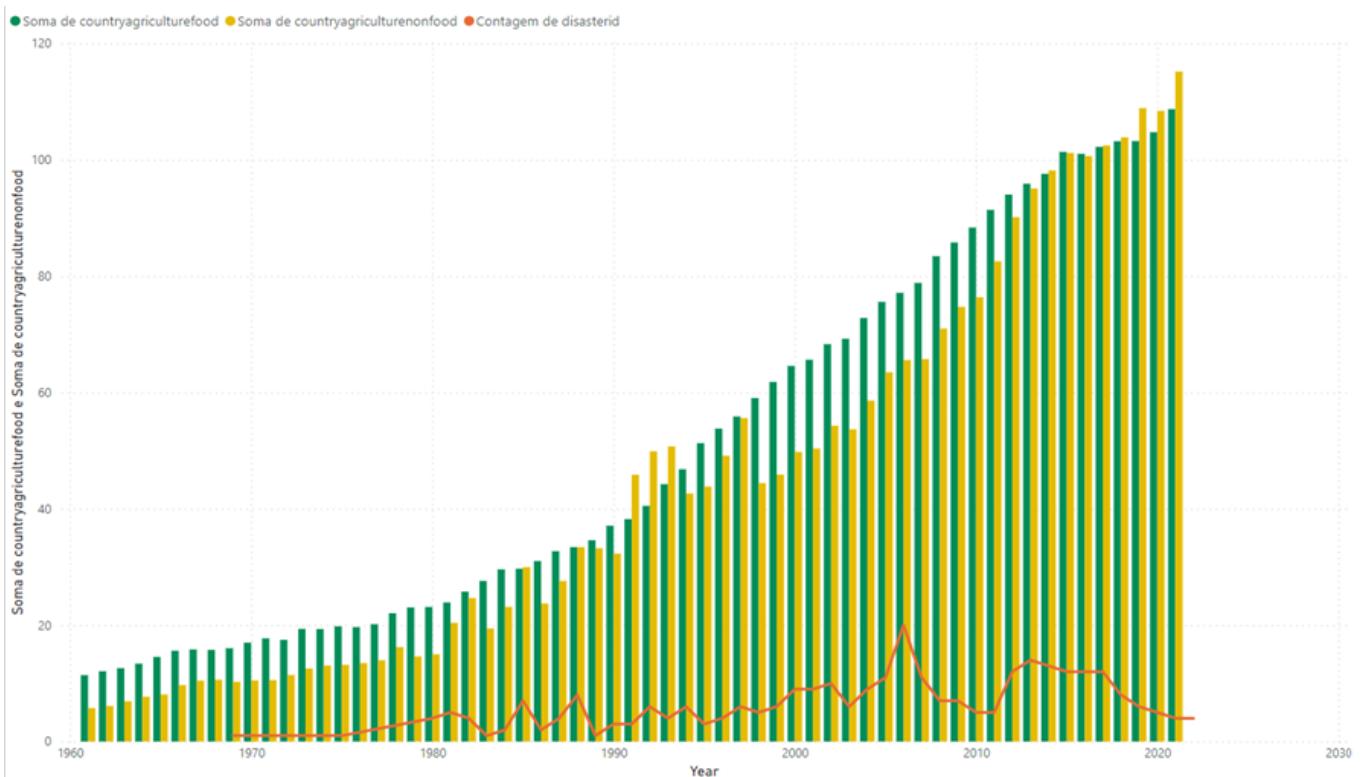


Figura 22. Produção agrícola vs Inundações na China

Indonésia - Inundações vs Produção agrícola:

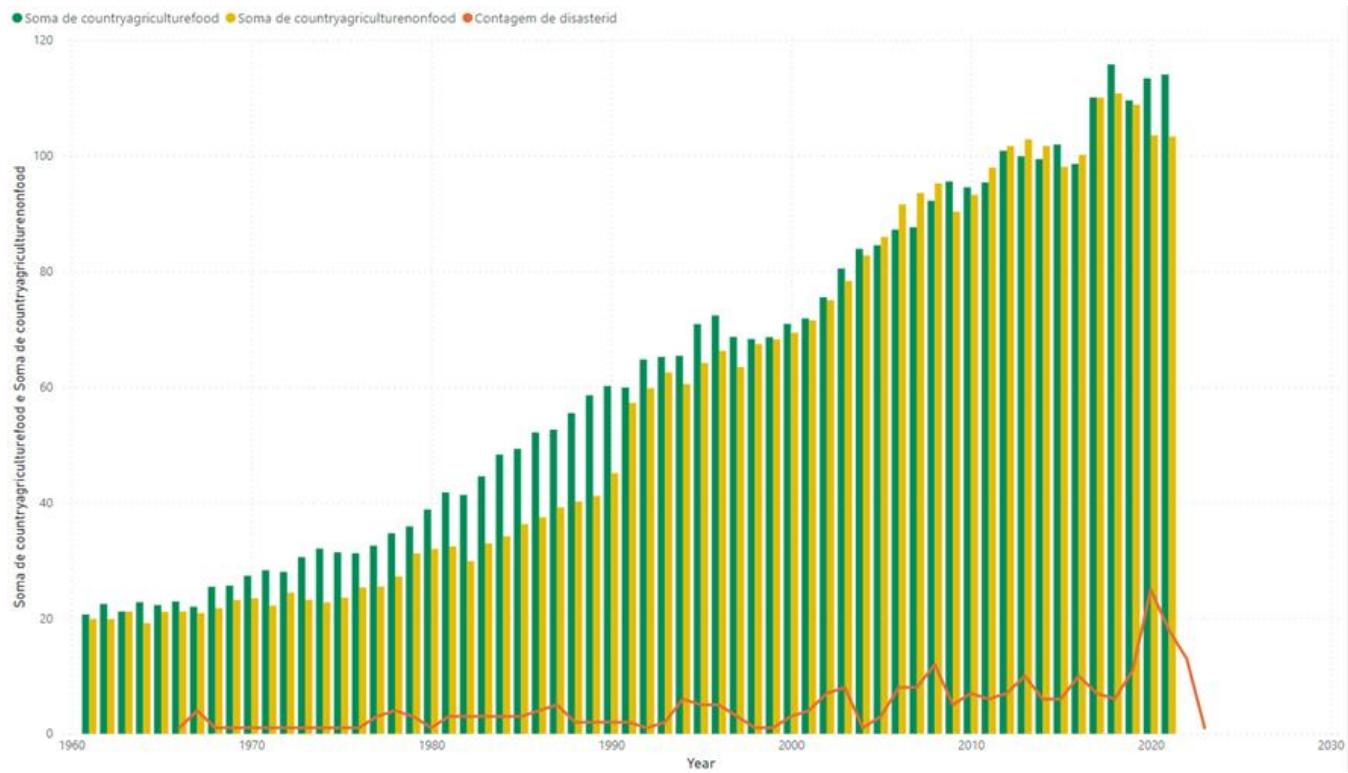


Figura 23. Produção agrícola vs Inundações na Indonésia

Estados Unidos - Inundações vs Produção agrícola:

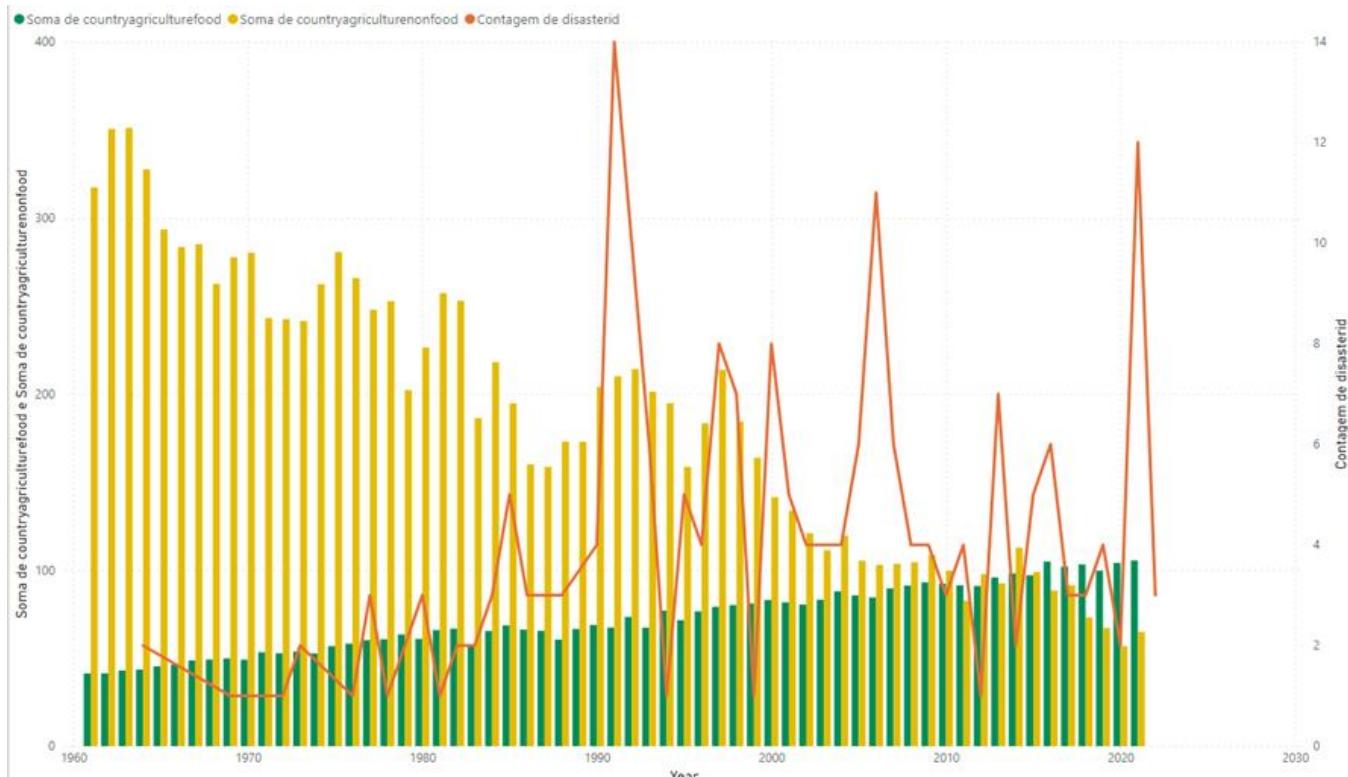


Figura 24. Produção agrícola vs Inundações nos Estados Unidos

Produção agrícola vs Inundações a nível global:

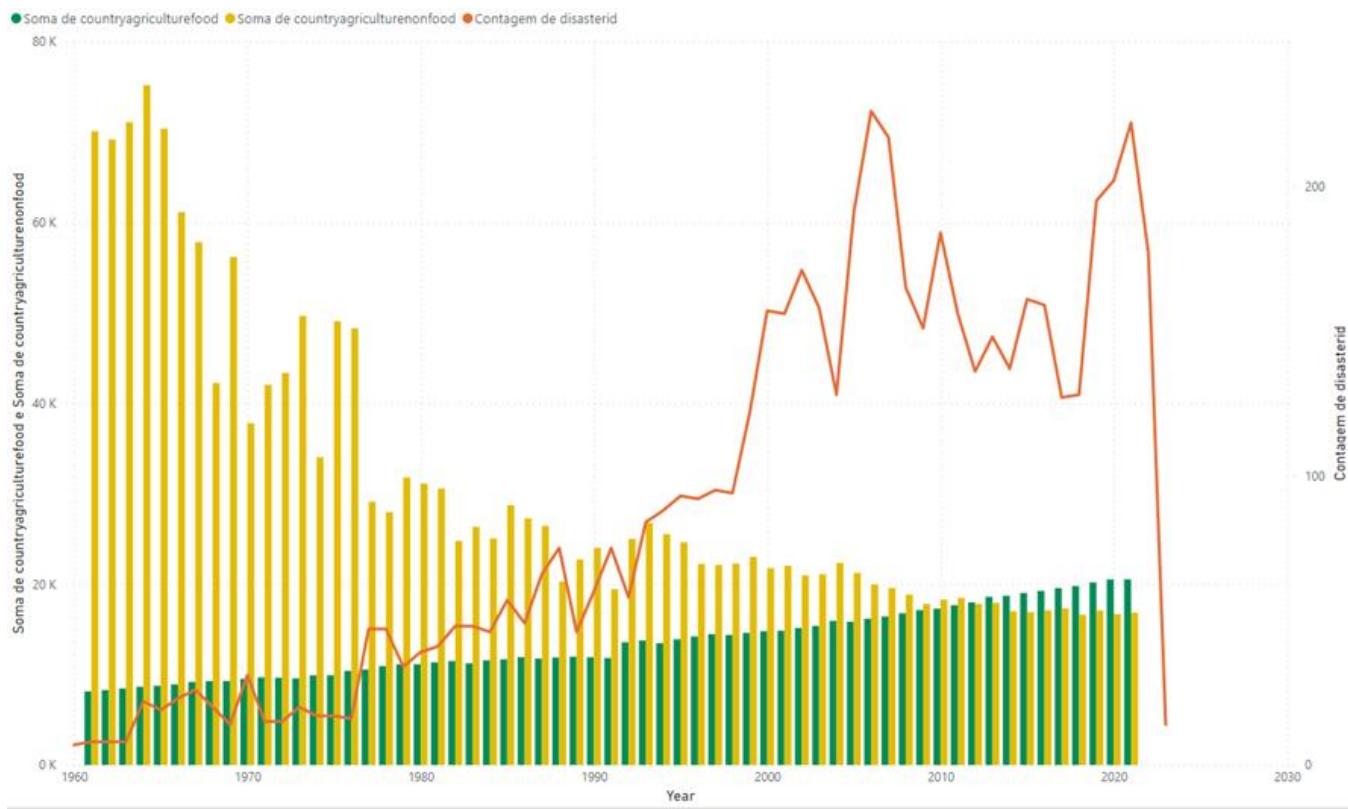


Figura 25. Produção agrícola vs Inundações no Mundo.

À semelhança das Secas, as inundações podem ter impacto na produção agrícola nos países acima, no entanto, não conseguimos tirar alguma relação entre a ocorrência destas e a produção agrícola, à exceção dos Estados Unidos. Para se poder estudar esta questão mais afundo em cada país, teríamos de ver com mais pormenor os vários produtos que podem ser afetados.

Nos Estados Unidos, podemos ver que as inundações afetaram a produção agrícola, pois a partir de 2000, à medida que começaram a ser mais frequentes, a produção agrícola diminuiu tanto para produtos alimentares como para não alimentares.

Na Figura 25, podemos ver que a nível global, o número de inundações aumentou e houve um decréscimo da produção não alimentar. As inundações podem ter influência neste tipo de produto, no entanto, com o avanço tecnológico, este tipo de produção agrícola pode ter decrescido.

Tempestades:

Tempestades, como furacões, tornados e ciclones, podem causar danos significativos às plantações, estruturas agrícolas e sistemas de armazenamento.

Fomos consultar quais os 10 países onde ocorrem mais tempestades e ver se a produção agrícola é afetada por este tipo de desastre para os 4 países mais afetados.

Tabela 24. Tabela com os países onde ocorreram mais tempestades

País	Nº de Ocorrências de tempestades
Estados Unidos	647
Filipinas	359
China	313
Índia	183
Bangladesh	168
Japão	159
Vietname	129
México	120
Austrália	85
Taiwan	81
França	76

Em baixo apresentamos os gráficos em que avaliamos a produção agrícola alimentar (cor verde) e não alimentar (cor amarela) vs o nº de tempestades (linha laranja), para os 4 países onde ocorreram mais tempestades.

Estados Unidos - Tempestades vs Produção agrícola:

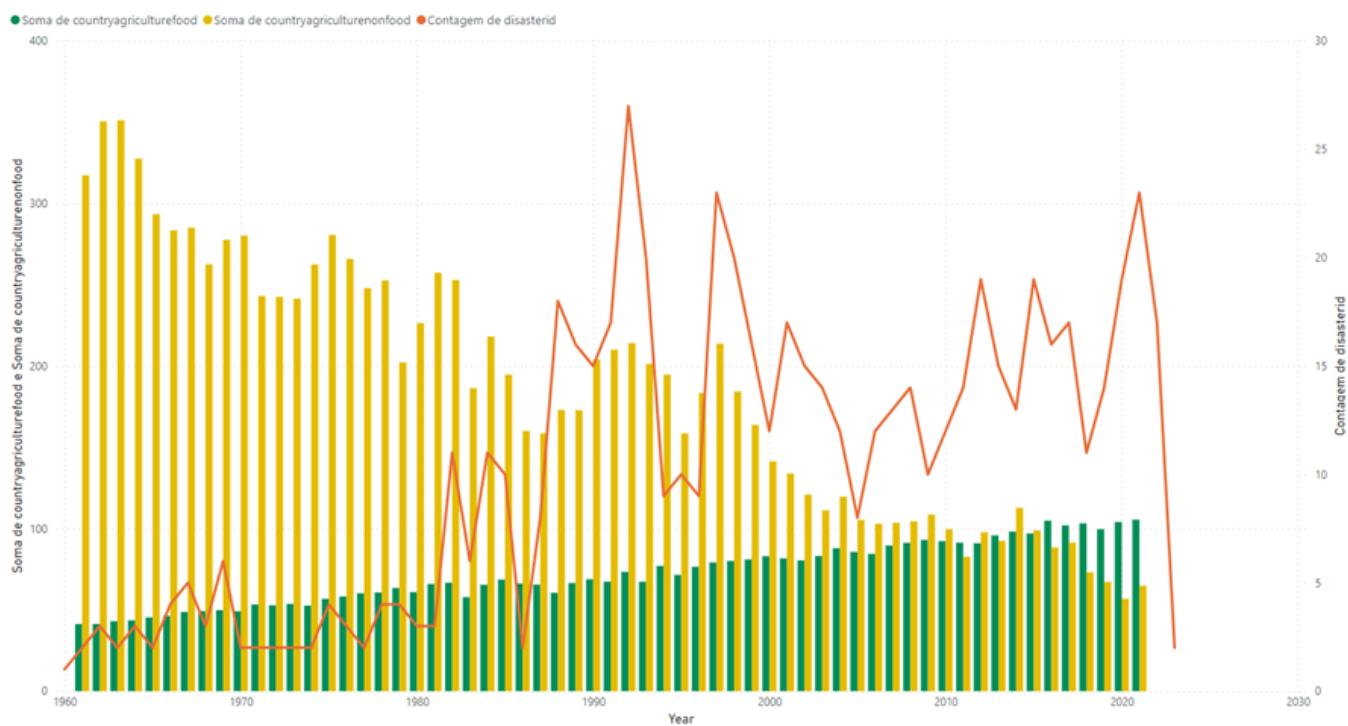


Figura 26. Produção agrícola vs Inundações nos Estados Unidos.

Filipinas - Tempestades vs Produção agrícola:

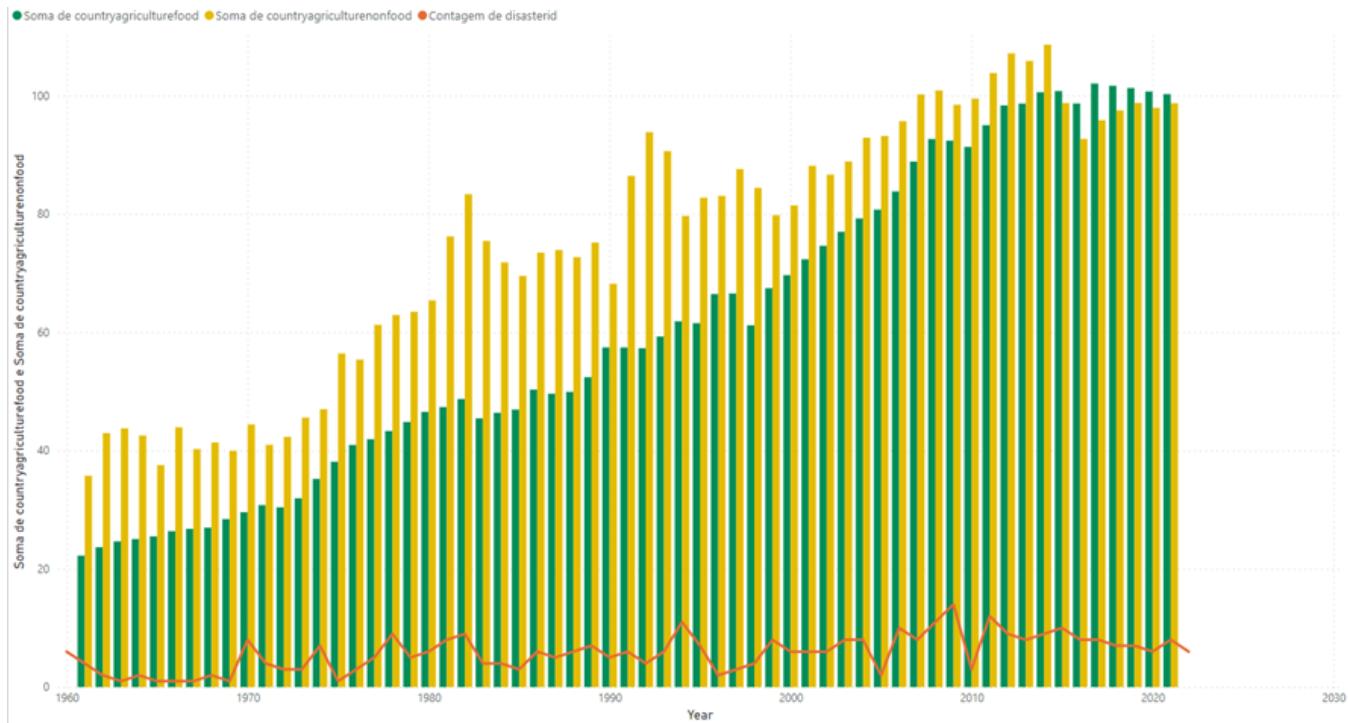


Figura 27. Produção agrícola vs Inundações nas Filipinas.

China - Tempestades vs Produção agrícola:

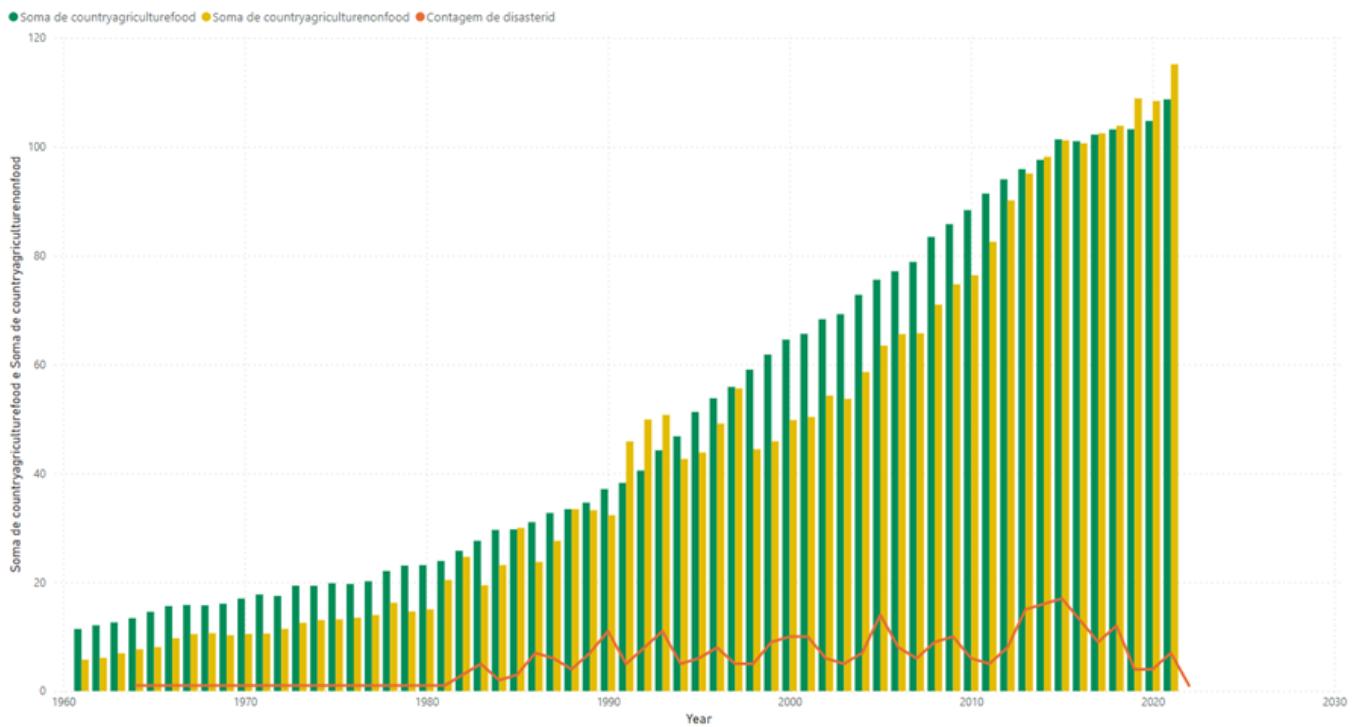


Figura 28. Produção agrícola vs Inundações na China.

Índia - Tempestades vs Produção agrícola:

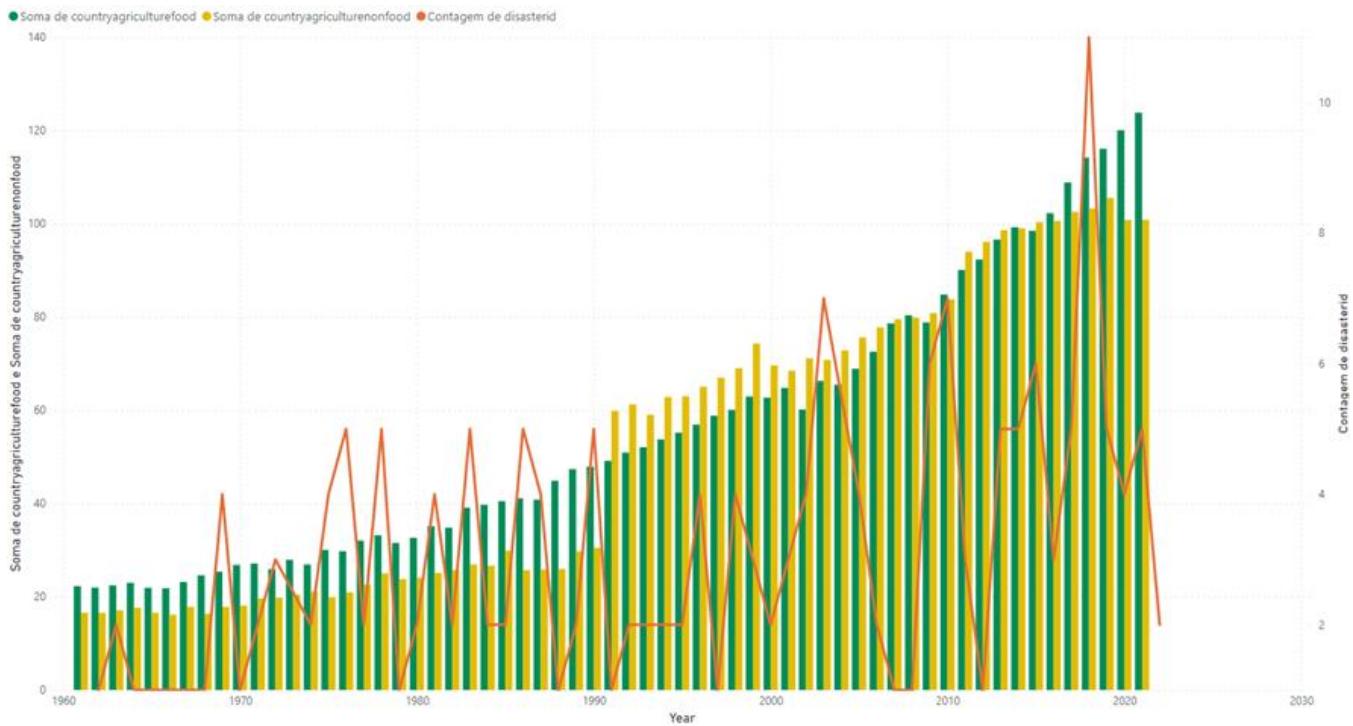


Figura 29. Produção agrícola vs Inundações na Índia.

Produção agrícola vs Tempestades a nível global:

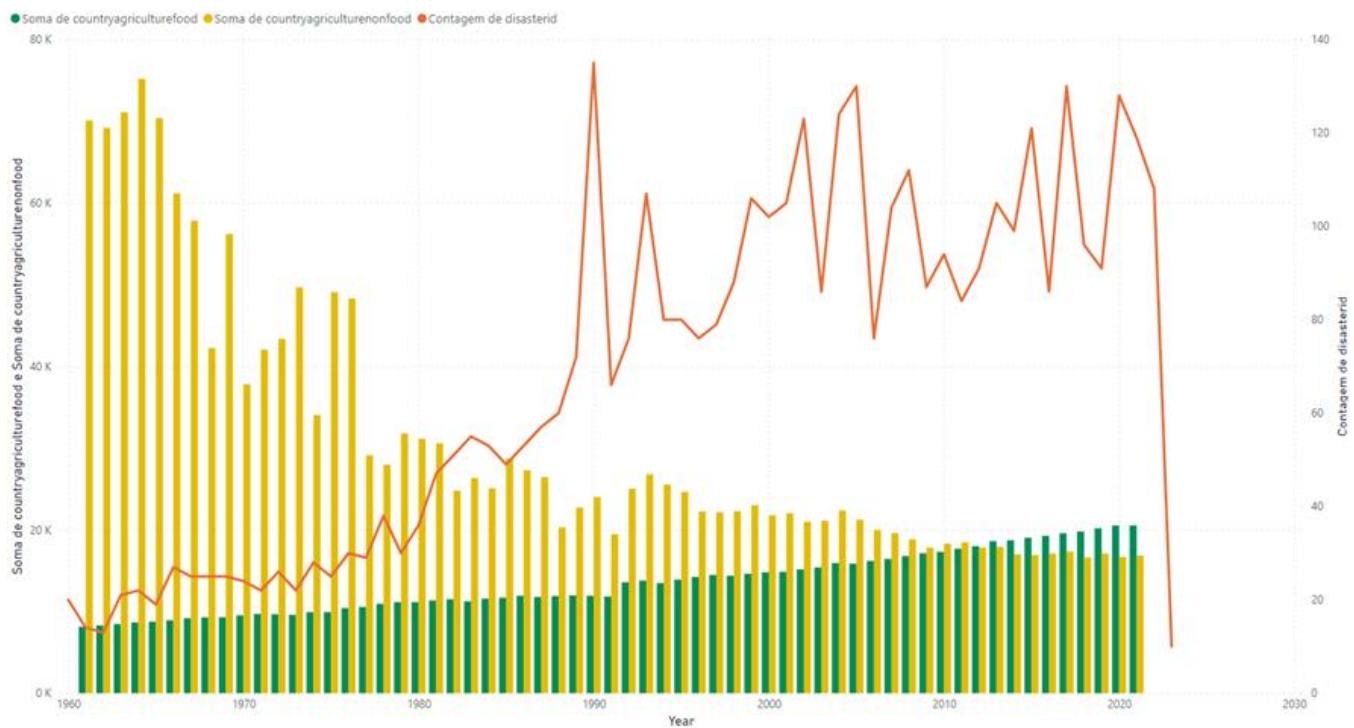


Figura 30. Produção agrícola vs tempestades no Mundo.

À semelhança das Secas e inundações, as tempestades podem ter impacto na produção agrícola nos países acima, no entanto, não conseguimos tirar alguma relação entre a ocorrência destas e a produção agrícola, à exceção dos Estados Unidos. Para se poder estudar esta questão mais a fundo em cada país, teríamos de ver com mais pormenor os vários produtos que podem ser afetados.

Nos Estados Unidos, podemos ver que, tal como as inundações, as tempestades afetaram a produção agrícola, pois a partir de 2000, à medida que começaram a ser mais frequentes, a produção agrícola diminuiu tanto para produtos alimentares como para não alimentares.

Na Figura 30, podemos ver que a nível global, o nº de tempestades aumentou e houve um decréscimo da produção não alimentar. As tempestades podem ter influência neste tipo de produto, no entanto, com o avanço tecnológico, este tipo de produção agrícola pode ter decrescido.

Incêndios florestais:

Incêndios descontrolados podem consumir grandes áreas de terras agrícolas, destruindo plantações, matando animais e contaminando o solo.

Fomos consultar quais os 10 países onde ocorrem mais incêndios florestais e ver se a produção agrícola é afetada por este tipo de desastre para os 4 países mais afetados.

Tabela 25. Tabela com os países onde ocorreram mais tempestades

País	Nº de Ocorrências de tempestades
Estados Unidos	92
Austrália	35
Rússia	24
Canadá	23
Espanha	18
Grécia	16
Chile	15
Portugal	15
França	12
Indonésia	11

Em baixo apresentamos os gráficos em que avaliamos a produção agrícola alimentar (cor verde) e não alimentar (cor amarela) vs o nº de incêndios florestais (linha laranja), para os 4 países onde ocorreram mais tempestades.

Estados Unidos – Incêndios florestais vs Produção agrícola:

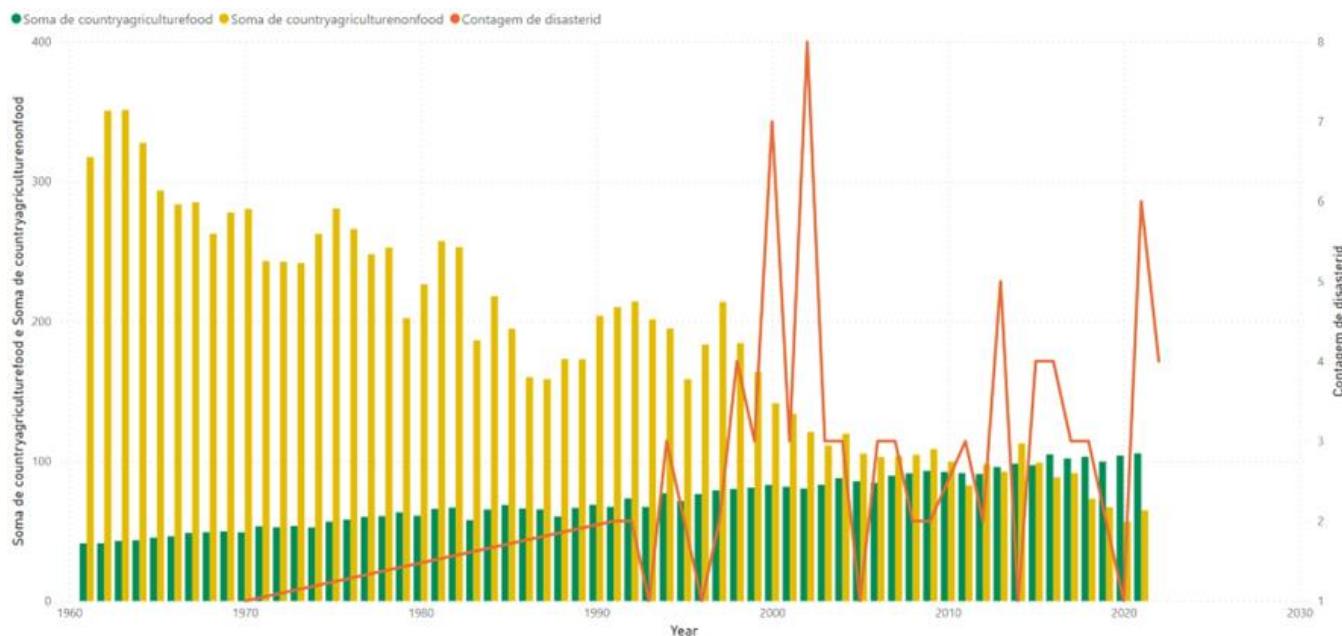


Figura 31. Produção agrícola vs incêndios florestais nos Estados Unidos.

Austrália – Incêndios florestais vs Produção agrícola:



Figura 32. Produção agrícola vs incêndios florestais na Austrália.

Rússia – Incêndios florestais vs Produção agrícola:

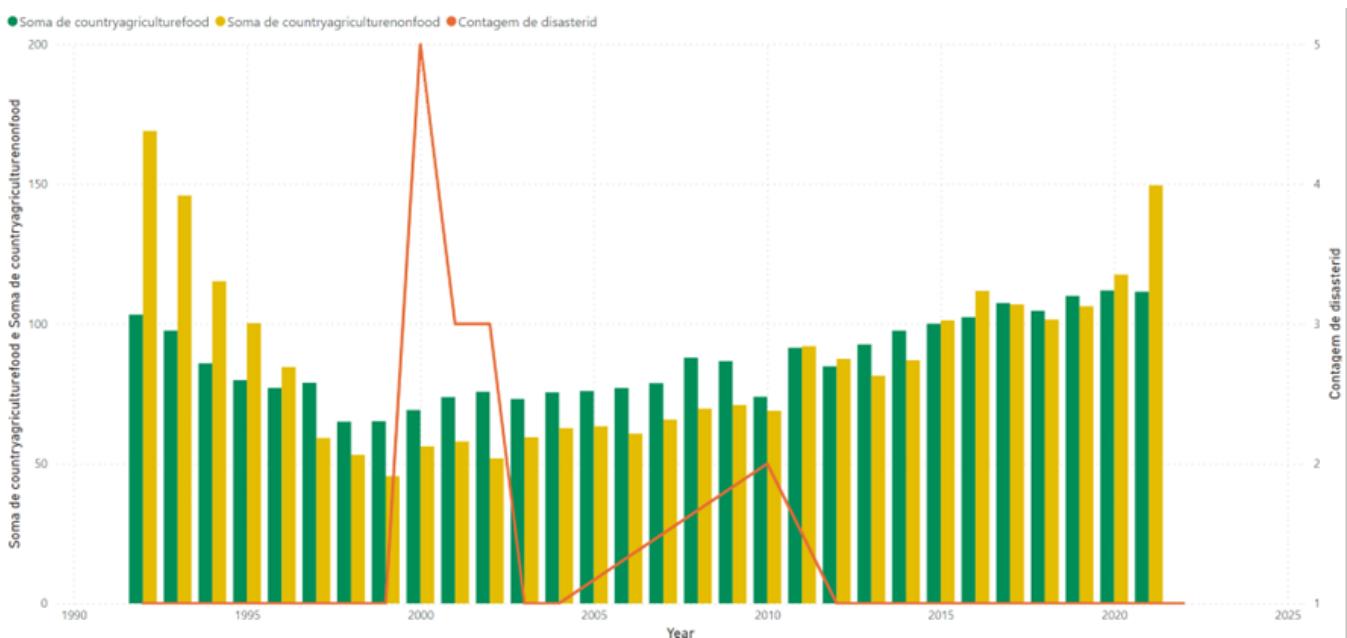


Figura 33. Produção agrícola vs incêndios florestais na Rússia.

Canadá – Incêndios florestais vs Produção agrícola:

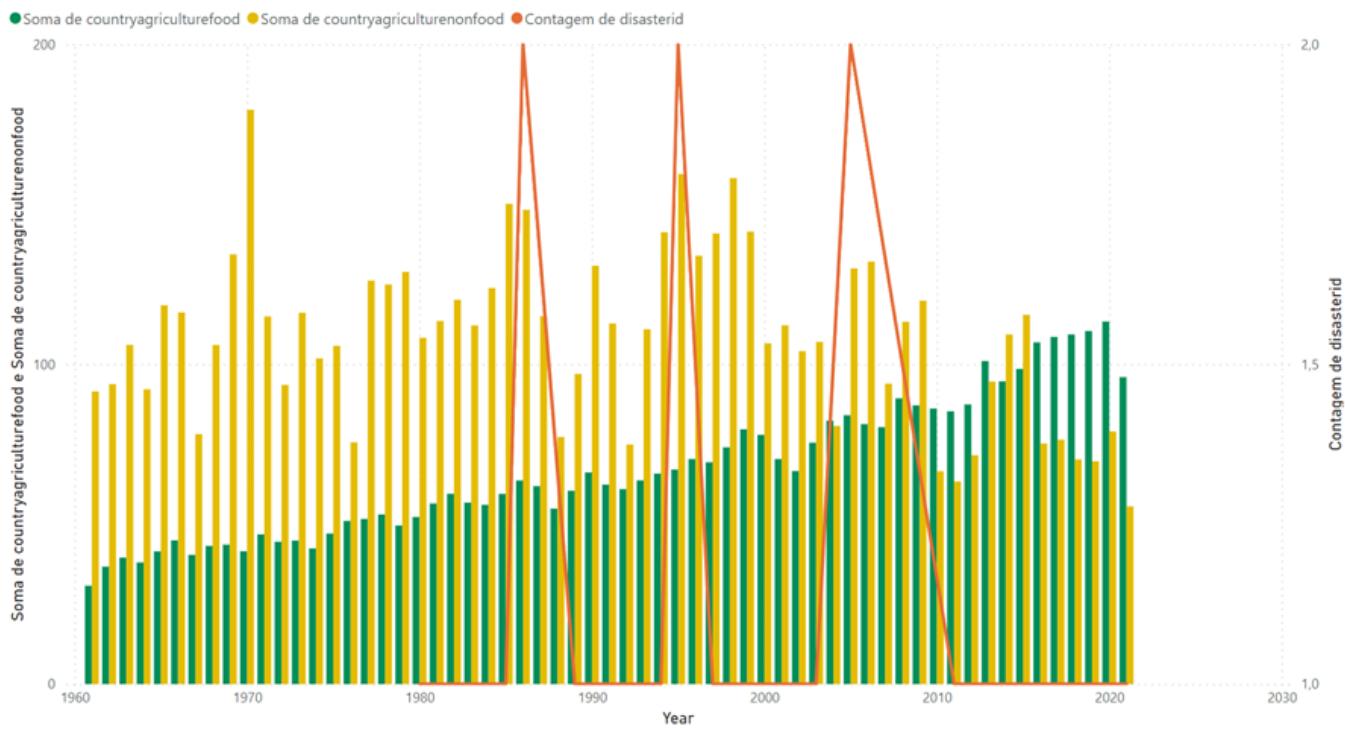


Figura 34. Produção agrícola vs Inundações no Canadá.

Produção agrícola vs Tempestades a nível global:

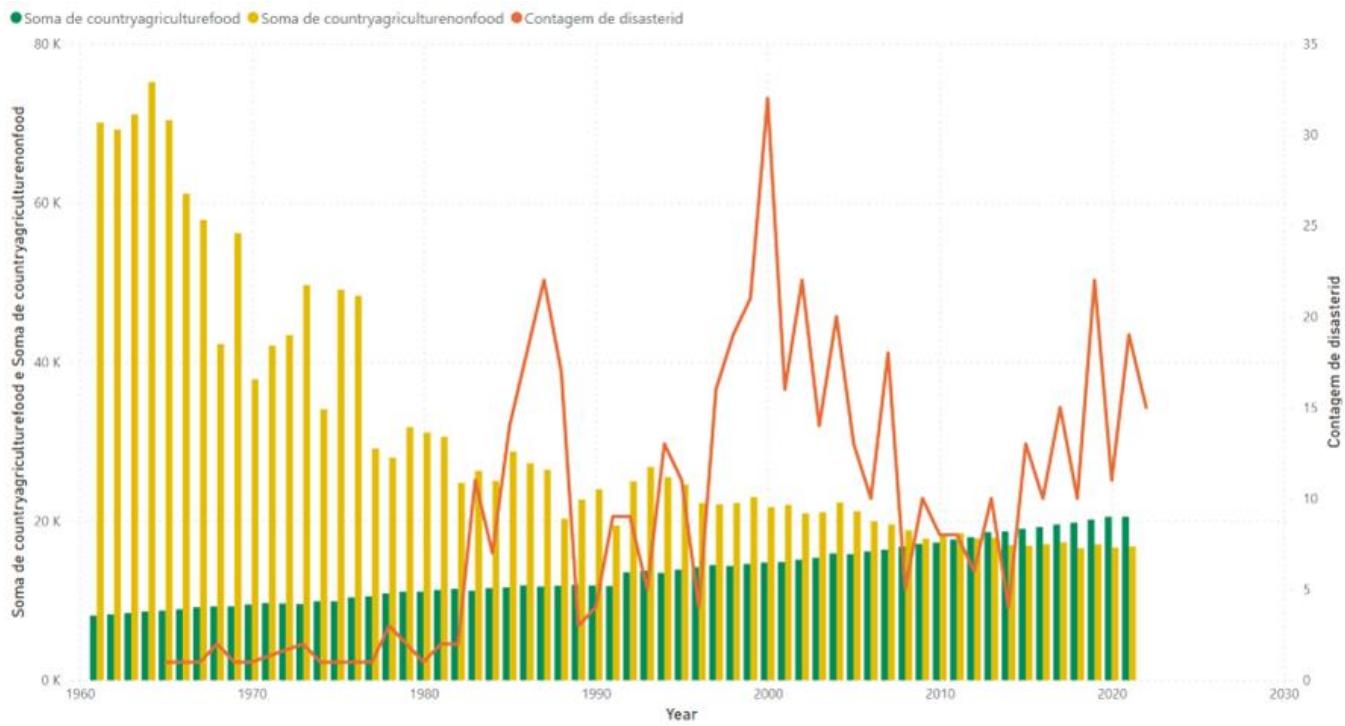


Figura 35. Produção agrícola vs tempestades no Mundo.

À semelhança dos desastres naturais mencionados anteriormente, os incêndios florestais podem ter impacto na produção agrícola nos países acima, no entanto, não conseguimos tirar nenhuma conclusão inequívoca entre a ocorrência destes e a produção agrícola. Para se poder estudar esta questão mais a fundo em cada país, teríamos de ver com mais pormenor os vários produtos que podem ser afetados.

Contudo, podemos ver que em todos os países, a produção agrícola para produtos não alimentares diminuiu no ano e/ou nos anos imediatamente a seguir. Nos Estados Unidos, é possível ver que a produção agrícola não alimentar diminui com o crescimento do número de incêndios florestais, principalmente nos últimos anos. Na Austrália, a produção agrícola foi fortemente afetada pelos incêndios desde 2006, pois tem valores baixos de produção num período de tempo em que está a ocorrer bastantes incêndios florestais. Esta observação também é aplicável ao Canadá, onde podemos ver uma baixa da produção agrícola após um incêndio florestal. Relativamente à Rússia, apenas temos dados desde 1990 e não conseguimos tirar nenhuma conclusão. Neste caso, teria de se ir mais ao pormenor, como por exemplo, estudar a produção durante o ano e ver o impacto nos meses a seguir.

Na Figura 35, podemos ver que a nível global, o nº de incêndios florestais tem aumentado nos últimos anos e houve um decréscimo da produção não alimentar. Os incêndios florestais podem ter influência neste tipo de produto, no entanto, com o avanço tecnológico, este tipo de produção agrícola pode ter decrescido.

10.2 Segunda Pergunta Analítica

“Existe alguma relação entre a afluência de turistas de um país e a frequência de ocorrência de catástrofes naturais? Se sim, como essas catástrofes afetam o desenvolvimento do setor turístico e a economia do país? Quais os tipos de desastres naturais que influenciam mais o turismo?”

A segunda pergunta analítica proposta inicialmente pretende perceber o possível impacto das catástrofes naturais no turismo dos países. Numa fase preliminar, esta questão foi efetuada de uma forma geral para todos os países, com o objetivo de comparar a afluência de turistas com a frequência de ocorrência de catástrofes naturais nos diferentes países. No entanto, nesta fase de exploração da pergunta, concluiu-se, em primeiro lugar, que para se perceber se existe algum impacto dos desastres, era necessário fazer uma avaliação temporal da afluência dos turistas antes, durante e depois do desastre. Deste modo, realizar esta análise para todos os países é impensável, sendo necessário escolher apenas alguns para realizar a questão. Em segundo lugar, utilizar a frequência de ocorrências de desastres naturais não é a melhor abordagem para este caso, tendo em conta que existem desastres com menos impacto, como deslizamento de terras ou tempestades de pouca intensidade, que iriam contar para frequência, mas que não foram noticiados em outros países (não produzindo um impacto no turismo). Assim, considerando o referido, a abordagem estabelecida partiu da intensidade dos desastres que, à partida, foram divulgados. A intensidade foi medida separadamente através de duas métricas da tabela de factos: o número de óbitos e o dano causado pelo desastre. Estas duas métricas são bastante relevantes para avaliar o impacto de um desastre, sendo as duas essenciais para o problema, uma vez que um desastre pode ter um grande número de mortes, mas não possuir qualquer tipo de danos, como numa pandemia, ou vice-versa, como num sismo que destrói completamente uma cidade, mas possui um número de mortes mais reduzido.

Como primeiro passo da análise, foram construídas duas tabelas com as métricas de forma a ordenar os desastres por ordem decrescente de impacto, sendo os primeiros 4 países diferentes mais afetados para cada métrica (8 no total) avaliados em termos de afluência de turistas. Os dados do turismo contêm informação a partir de 1995, tendo sido aplicado um filtro à tabela para que esta apresente desastres somente a partir desse ano (Tabela 26 e 27)

Tabela 26 – 4 desastres com maior número de mortes desde 1995 e respetivos países, tipo de desastre, ano e mês de ocorrência.

ID Desastre	País	Mortes	Tipo Desastre	Danos	Mês	Ano
2010-0017-HTI	Haiti	222570	Earthquake	8000000	January	2010
2004-0659-IDN	Indonesia	165708	Earthquake	4451600	December	2004
2008-0184-MMR	Myanmar	138366	Storm	4000000	May	2008
2008-0192-CHN	China	87476	Earthquake	85000000	May	2008

Tabela 27 – 4 desastres com maior número de danos de infraestruturas desde 1995 e respetivos países, tipo de desastre, ano e mês de ocorrência.

ID Desastre	País	Mortes	Tipo Desastre	Danos	Mês	Ano
2011-0082-JPN	Japan	19846	Earthquake	210000000	March	2011
2005-0467-USA	United States	1833	Storm	125000000	August	2005
2021-0411-DEU	Germany	197	Flood	40000000	July	2021
2011-0326-THA	Thailand	813	Flood	40000000	August	2011

Os países escolhidos para a análise do impacto dos desastres no turismo foram o Haiti, Indonésia, Birmânia, China, Japão, Estados Unidos, Alemanha e Tailândia. A análise foi então efetuada tendo em conta todos os anos desde 1995 até 2023 para cada país, mostrando a evolução do turismo no país ao longo do tempo e em simultâneo, o número de mortes ou os custos em infraestruturas dependendo do país em análise (Figuras 36, 37, 38, 39, 40, 41, 42 e 43). O comportamento procurado neste tipo de gráficos é um vale na altura do desastre, no entanto, este vale tanto pode ser observado no mesmo ano ou no ano posterior ao acontecimento. Isto deve-se ao facto do efeito não ser instantâneo, estado a altura do ano que se dá a catástrofe relacionada com a forma como se analisa o gráfico. Se esta acontecer no início do ano, o seu impacto é evidente no próprio ano, no entanto, se esta acontecer no final do ano, o impacto só é perceptível, provavelmente, no ano seguinte. Os primeiros 4 países, derivados do maior número de óbitos, estão representados nas figuras 36, 37, 38 e 39.

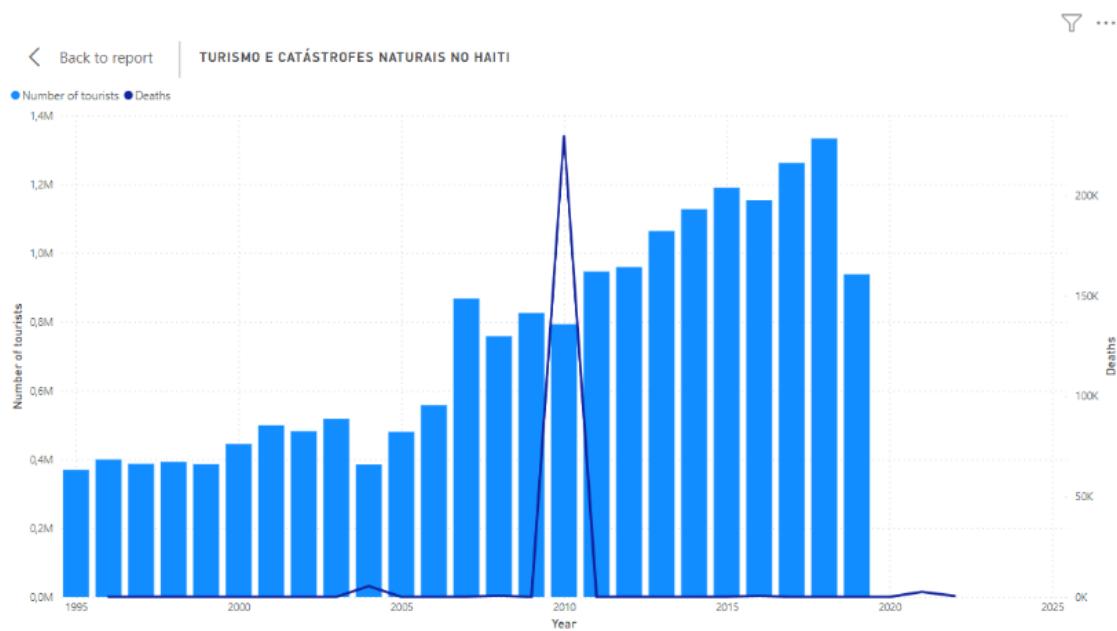


Figura 36- Número de turistas e óbitos devido a catástrofes naturais de 1995 a 2023 no Haiti.

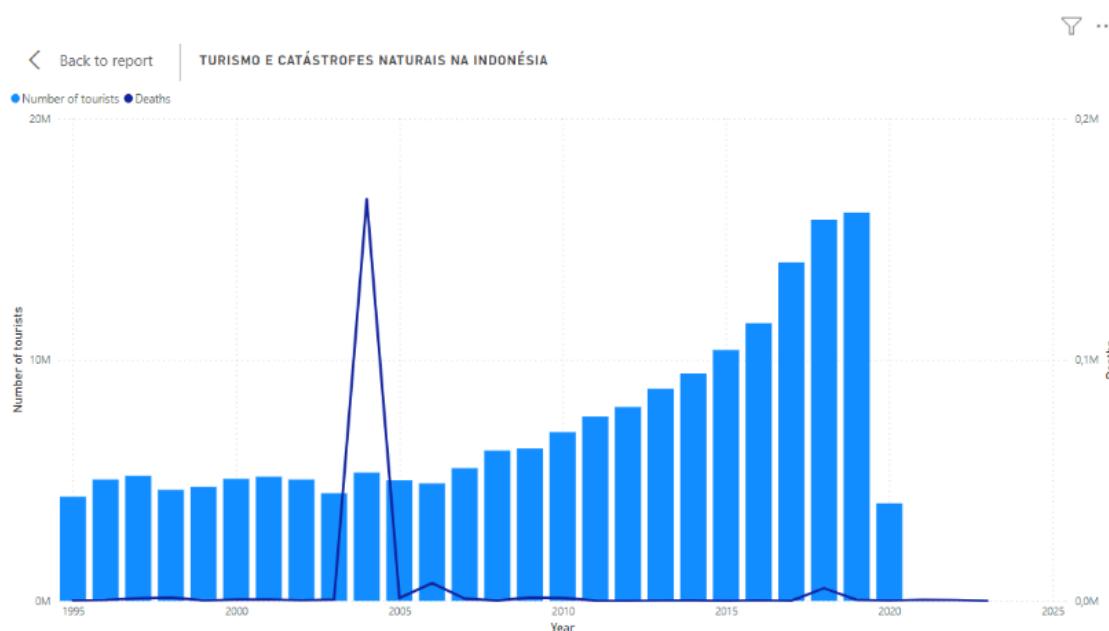


Figura 37- Número de turistas e óbitos devido a catástrofes naturais de 1995 a 2023 na Indonésia.

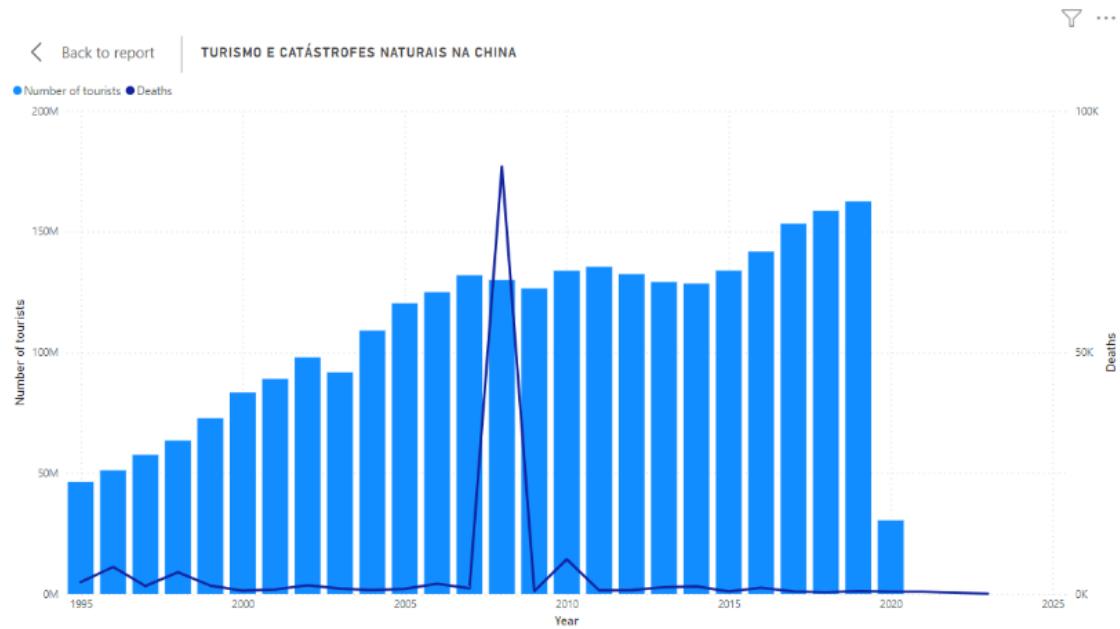


Figura 38 - Número de turistas e óbitos devido a catástrofes naturais de 1995 a 2023 na China.

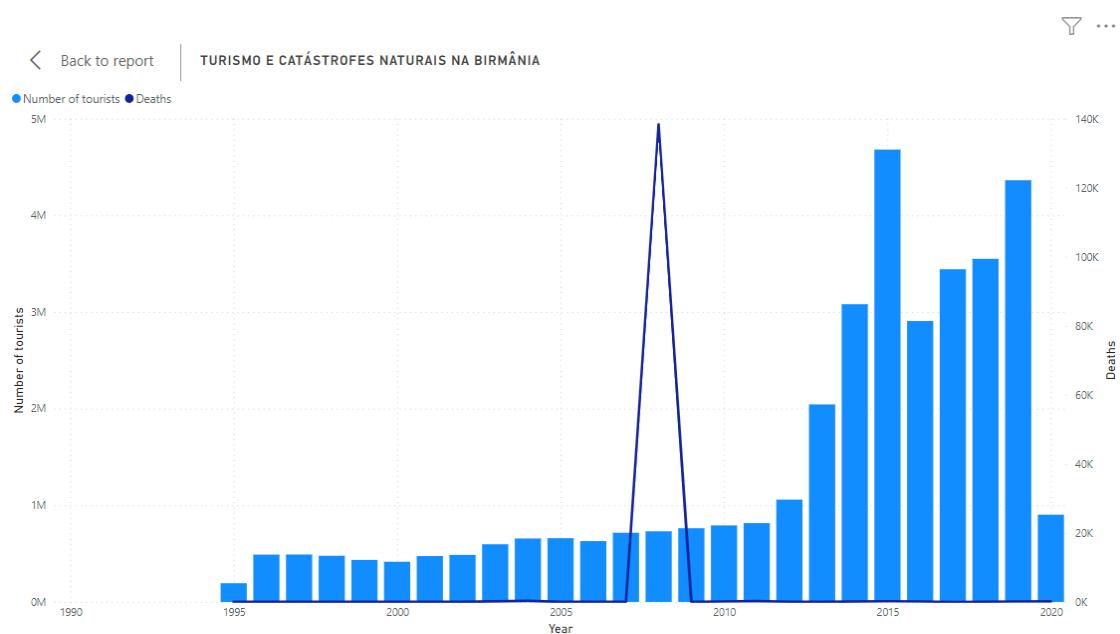


Figura 39- Número de turistas e óbitos devido a catástrofes naturais de 1995 a 2023 na Birmânia.

É possível verificar nas quatro figuras um pico na linha azul-escura que representa o número de óbitos derivados das catástrofes naturais ocorridas nesse ano. Nestes casos, os picos referem-se aos desastres específicos representados na tabela 26. Em todos os países é possível verificar um crescimento positivo do turismo desde 1995, o que mostra um desenvolvimento no setor e provavelmente um impacto positivo na economia do país. Em nenhum caso o desastre foi grande o suficiente para causar um decréscimo prolongado no turismo, sendo que se este existiu, é apenas perceptível no próprio ano ou no seguinte. No caso do Haiti (Figura 36), Indonésia (Figura 37) e China (Figura 38), é possível verificar um ligeiro decréscimo no ano ou no ano seguinte do número de turistas, no entanto, tendo em conta que existe uma variância natural todos os anos, não é possível tirar nenhuma conclusão sobre o impacto dos desastres naturais. No caso da Birmânia (Figura 39), o número de turistas continuou a crescer de forma estável, não tendo sido afetado de todo. Uma possível explicação é também o tipo de desastre que ocorreu (tempestade), que, apesar de ter um elevado número de mortes, não produz tanto choque nos turistas como um terramoto. Para ser possível o desenvolvimento de uma conclusão, era necessário verificar várias vezes o mesmo padrão e este teria de ser mais acentuado para que não fosse facilmente confundível com a variabilidade intrínseca dos dados, uma vez que, naturalmente, existem mais anos com turismo do que outros, seja devido à economia, meteorologia, etc.

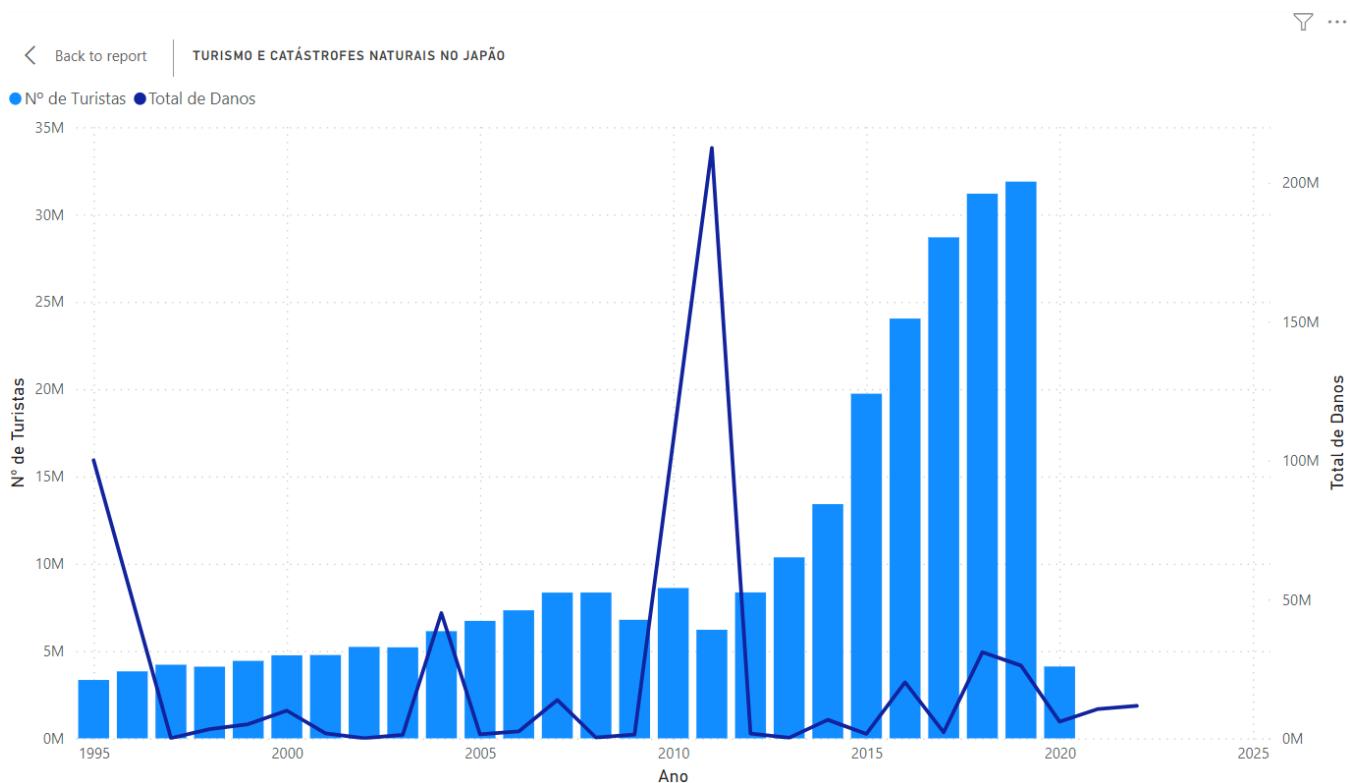


Figura 40 - Número de turistas e dano em infraestruturas devido a catástrofes naturais de 1995 a 2023 no Japão.

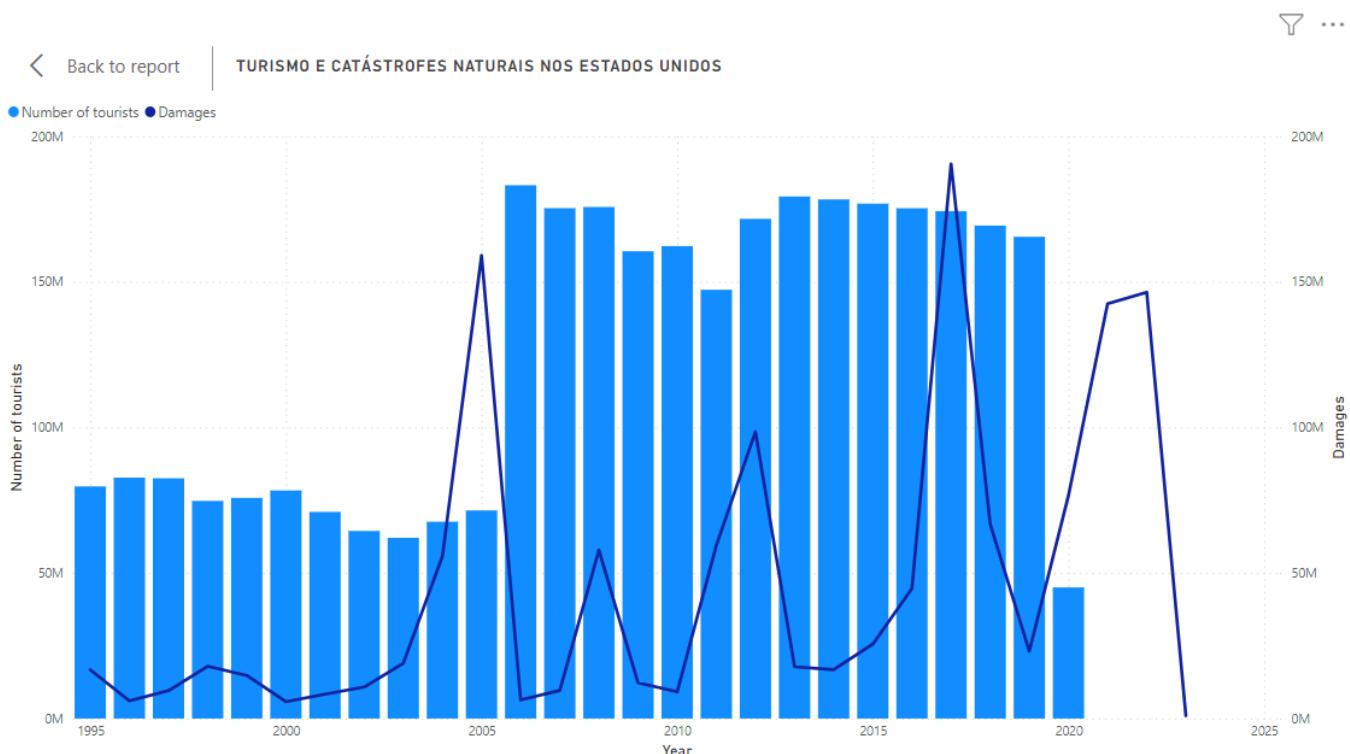


Figura 41 - Número de turistas e dano em infraestruturas de 1995 a 2023 nos Estados Unidos.

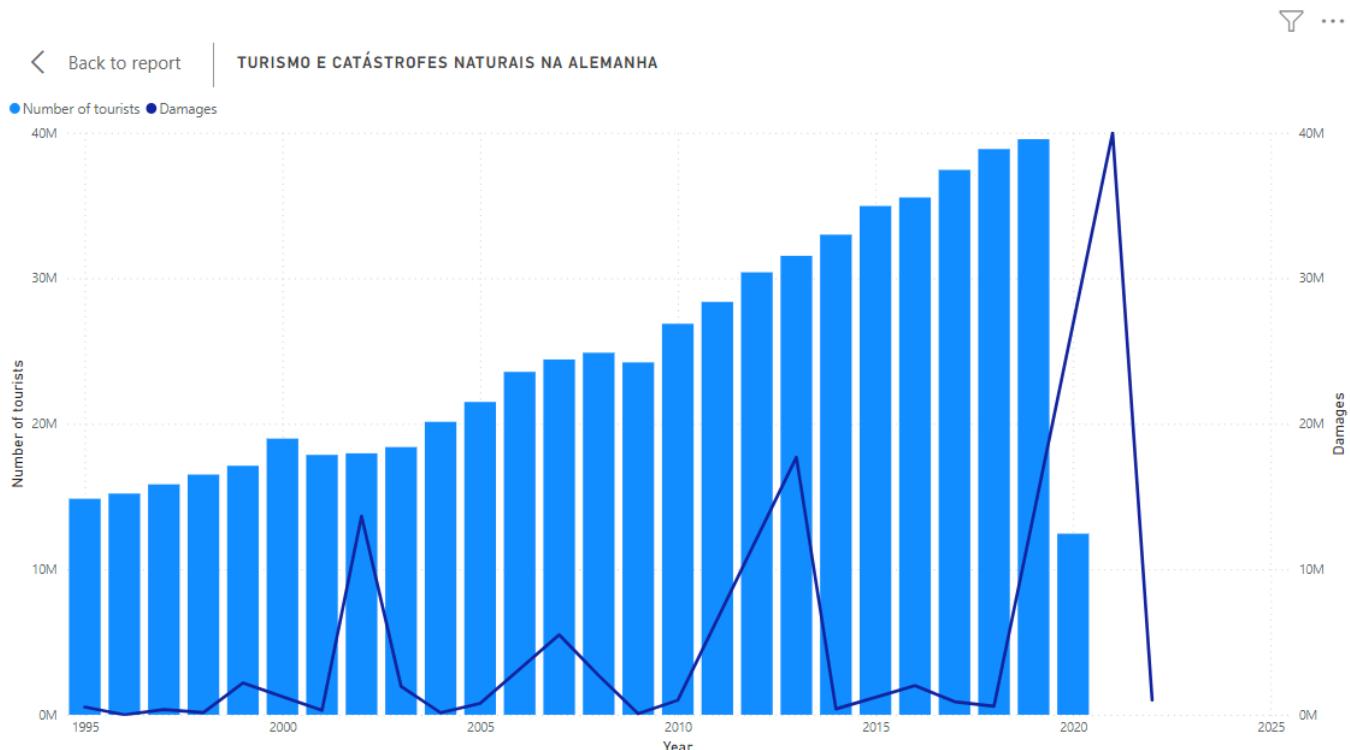


Figura 42 - Número de turistas e dano em infraestruturas de 1995 a 2023 na Alemanha.

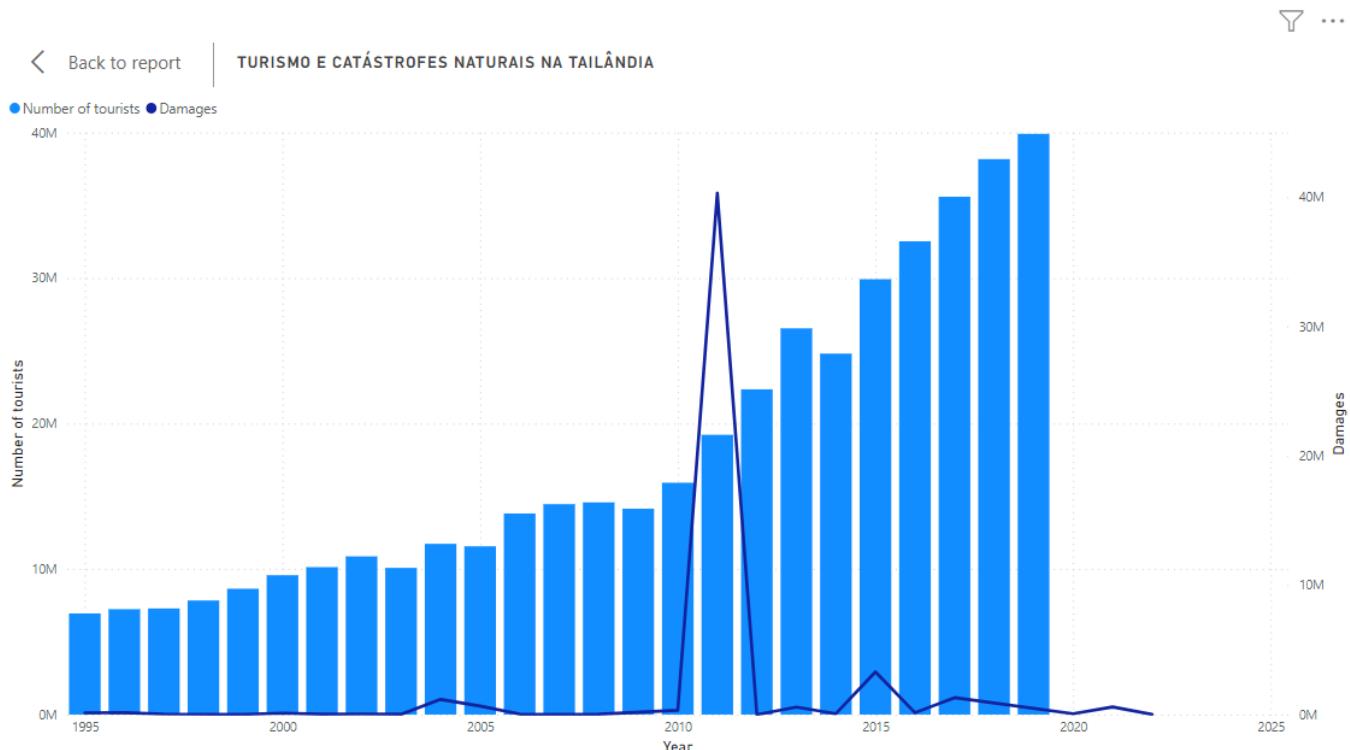


Figura 43 - Número de turistas e dano em infraestruturas devido a catástrofes naturais de 1995 a 2023 na Tailândia.

Os restantes países representados nas figuras 40, 41, 42 e 43, são referentes aos desastres com maiores danos nas infraestruturas. Numa análise geral e de forma semelhante aos casos anteriores, os picos de danos não apresentam qualquer relação com a afluência de turistas, sendo que em todos os casos, o número de turistas aumentam de ano para ano. É possível verificar no Japão (Figura 40), Estados Unidos (Figura 41), Alemanha (Figura 42), Indonésia (Figura 37), China (Figura 38) e Birmânia (Figura 39) que no ano de 2020, existe uma queda dramática no número de turistas. Este decréscimo está associado à pandemia de COVID-19 que, atualmente, não se encontra contabilizada no dataset dos desastres naturais. No entanto, esta pandemia é um exemplo de um desastre que teve um impacto no setor turístico bastante acentuado e, por conseguinte, na economia dos países. Para avaliar o tipo de catástrofe natural que mais influencia o turismo, teria de ser feita uma análise com mais países e que abrangesse todo o tipo de desastres, o que seria desproporcional para este relatório. Assim, não foi possível responder à última parte da segunda questão analítica, no entanto, é possível entender que as pandemias são provavelmente as maiores catástrofes e que têm um impacto inigualável no turismo. Por outro lado, catástrofes que têm magnitude suficiente para serem divulgadas na comunicação social, também podem ter um impacto no turismo, mas seria algo pontual e muito complicado de se medir (como verificada na análise efetuada). Uma possível explicação é que este tipo de eventos, como terramoto ou cheias, acontecem tipicamente em regiões mais específicas, como vilas ou cidades, sendo que esta análise foi efetuada para os países inteiros, diluindo o impacto dos desastres.

10.3 Terceira Pergunta Analítica

“Qual é o impacto das catástrofes naturais no índice de desenvolvimento de crescimento de um país e como fica afetada a taxa de desemprego e a produtividade da população? A inflação tem um valor de expressão maior nos países não desenvolvidos?”

A terceira pergunta analítica relaciona-se com o impacto dos desastres naturais no desenvolvimento de um país e na sua taxa de desemprego. Inicialmente, foi proposta a análise da produtividade da população e a análise da inflação, mas o data warehouse não possui os dados indicados para responder a esta parte da pergunta analítica. Assim, pretende-se nesta fase explorar qual o impacto das catástrofes naturais no produto interno bruto e na taxa de desemprego. A hipótese inicial seria que, catástrofes de grande dimensão iriam ter um impacto percetível no rendimento de um país que, por sua vez, colocaria pressão nos postos de trabalho e iria aumentar o desemprego.

De forma semelhante à segunda pergunta analítica e pelos mesmos motivos apresentados, foram escolhidos 8 países com base nos desastres mais graves, avaliados tanto a nível de óbitos como a nível de danos nas infraestruturas. Tendo em conta que existem dados a partir de 1990 para o PIB e para a taxa de desemprego, escolheu-se aplicar esse filtro aos países. As tabelas 28 e 29 apresentam os 4 países para cada situação supracitada. Estes 8 países foram avaliados ao nível do produto interno bruto e da taxa de desemprego, verificando se existiam decréscimos destes valores aquando da ocorrência de um desastre mais grave.

Tabela 28 – 4 desastres com maior número óbitos desde 1990 e respetivos países, tipo de desastre, ano e mês de ocorrência.

ID Desastre	País	Mortes	Tipo Desastre	Danos	Mês	Ano
2010-0017-HTI	Haiti	222570	Sismo	8000000	Janeiro	2010
2004-0659-IDN	Indonésia	165708	Sismo	4451600	Dezembro	2004
1991-0120-BGD	Bangladesh	138866	Tempestade	1780000	Abril	1991
2008-0184-MMR	Birmânia	138366	Tempestade	4000000	Maio	2008

Tabela 29 – 4 desastres com maior número de danos de infraestruturas desde 1990 e respetivos países, tipo de desastre, ano e mês de ocorrência.

ID Desastre	País	Mortes	Tipo Desastre	Danos	Mês	Ano
2011-0082-JPN	Japão	19846	Sismo	210000000	Março	2011
2005-0467-USA	Estados Unidos	1833	Tempestade	125000000	Agosto	2005
2008-0192-CHN	China	87476	Sismo	85000000	Maio	2008
2011-0326-THA	Tailândia	813	Cheia	40000000	Agosto	2011

Os países escolhidos para a análise do impacto dos desastres no desenvolvimento económico e desemprego foram o Haiti, Indonésia, Bangladesh, Birmânia, Japão, Estados Unidos, China e Tailândia. A análise foi então efetuada tendo em conta todos os anos desde 1990 até 2023 para cada país, mostrando a evolução do PIB e da taxa de desemprego no país ao longo do tempo e em simultâneo, o número de mortes ou os custos em infraestruturas dependendo do país em análise (Figuras 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58 e 59). De forma semelhante à pergunta anterior com o turismo, o comportamento procurado neste tipo de gráficos é um vale na altura do desastre para o PIB e um pico para o desemprego, no entanto, este vale/pico tanto pode ser observado no mesmo ano ou no ano posterior ao acontecimento, dependendo da altura do ano em que ocorre o desastre. Os primeiros 4 países, derivados do maior número de óbitos, tanto para o PIB como para o Desemprego, estão representados nas figuras 44, 45, 46, 47, 48, 49, 50 e 51.

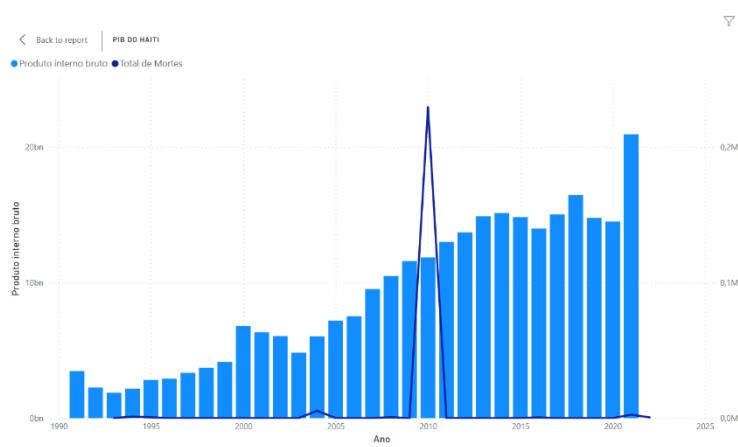


Figura 44. PIB e óbitos devido a catástrofes naturais de 1990 a 2023 no Haiti.

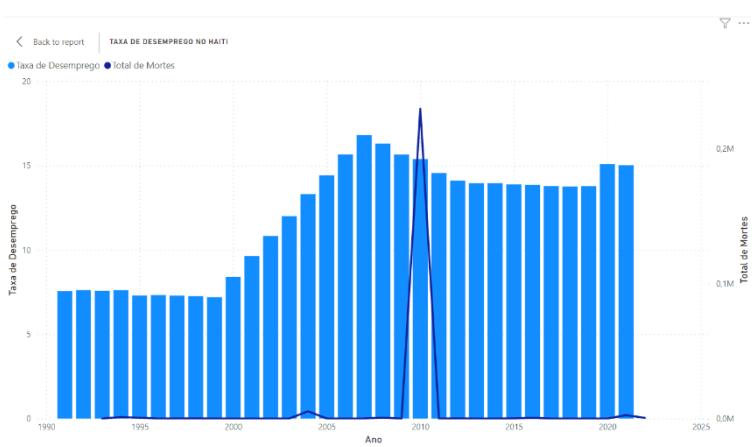


Figura 45. Taxa de desemprego e óbitos devido a catástrofes naturais de 1990 a 2023 no Haiti.

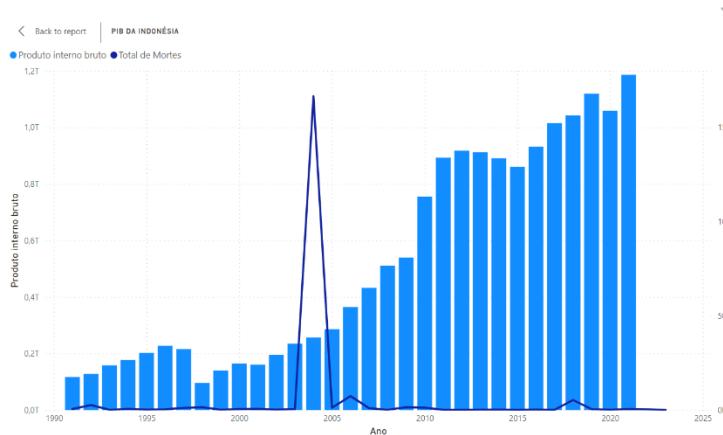


Figura 46. PIB e óbitos devido a catástrofes naturais de 1990 a 2023 na Indonésia.

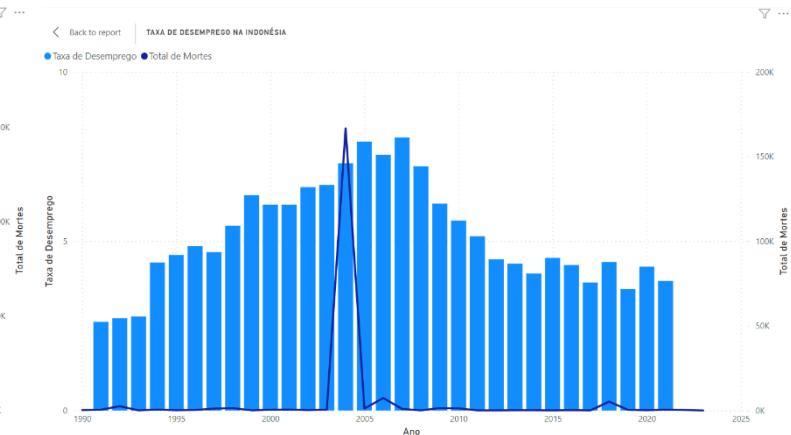


Figura 47. Taxa de desemprego e óbitos devido a catástrofes naturais de 1990 a 2023 na Indonésia.

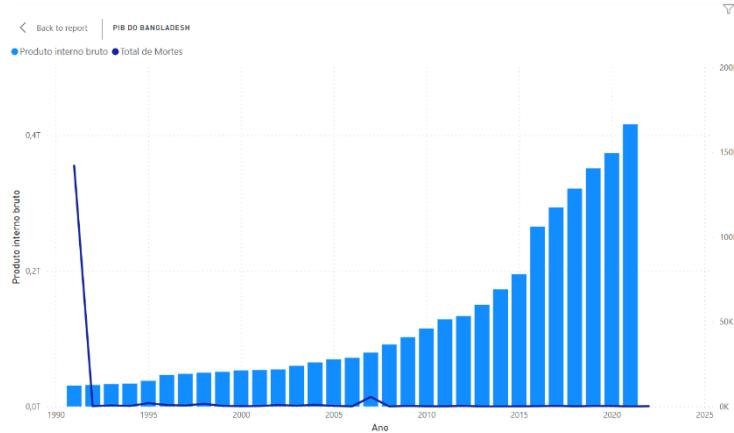


Figura 48. PIB e óbitos devido a catástrofes naturais de 1990 a 2023 no Bangladesh.

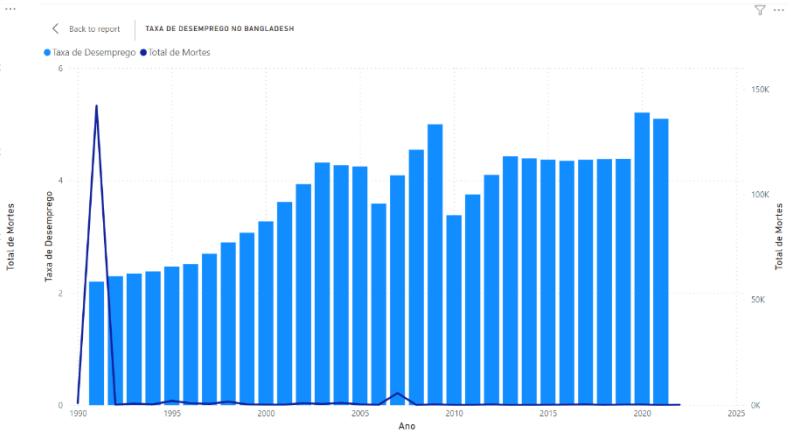


Figura 49. Taxa de desemprego e óbitos devido a catástrofes naturais de 1990 a 2023 do Bangladesh

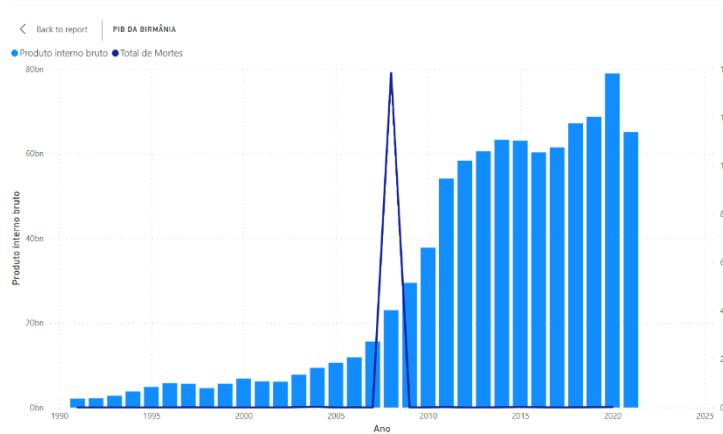


Figura 50. PIB e óbitos devido a catástrofes naturais de 1990 a 2023 na Birmânia.

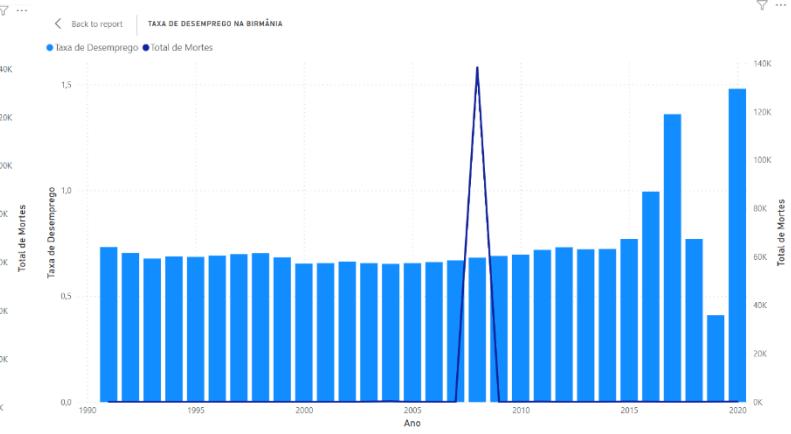


Figura 51. Taxa de desemprego e óbitos devido a catástrofes naturais de 1990 a 2023 na Birmânia

Em todos os gráficos apresentados acima relativos ao impacto dos desastres baseados nos óbitos em relação ao PIB dos países, é possível verificar os picos onde ocorreram os desastres. A nível do PIB, é possível verificar uma tendência crescente em todos os países. Quando verificadas as regiões onde ocorrem os desastres, não se encontra nenhum tipo de vale, sugerindo que as catástrofes não possuem magnitude suficiente para causar um impacto direto no PIB. Além da magnitude, o facto da maioria dos desastres ser local, implica que a existir um impacto, seria na localização e não na média global do país. No caso do desemprego, verifica-se a mesma situação, sendo que os altos e baixos da taxa de desemprego são derivados de uma evolução contínua dessa taxa e não de um impacto dos desastres naturais.

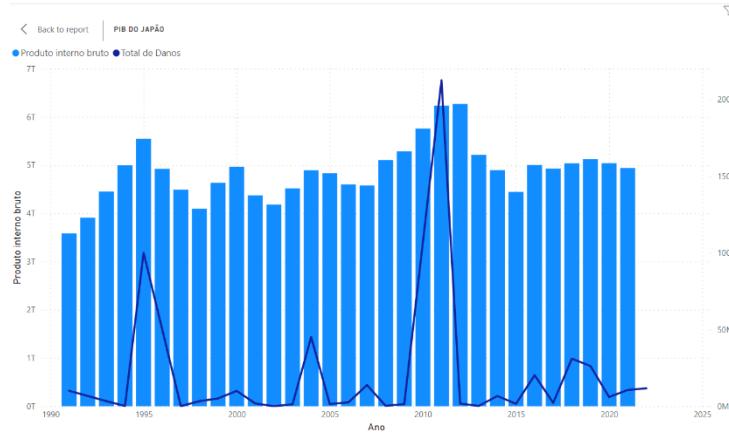


Figura 52. PIB e danos em infraestruturas devido a catástrofes naturais de 1990 a 2023 no Japão.

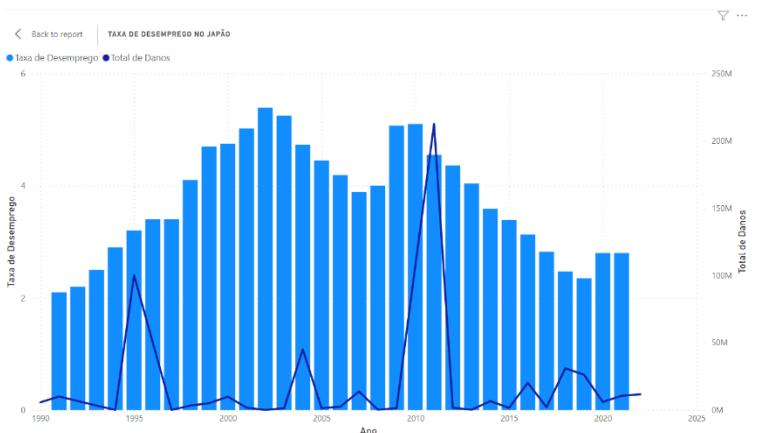


Figura 53. Taxa de desemprego e danos em infraestruturas devido a catástrofes naturais de 1990 a 2023 no Japão.

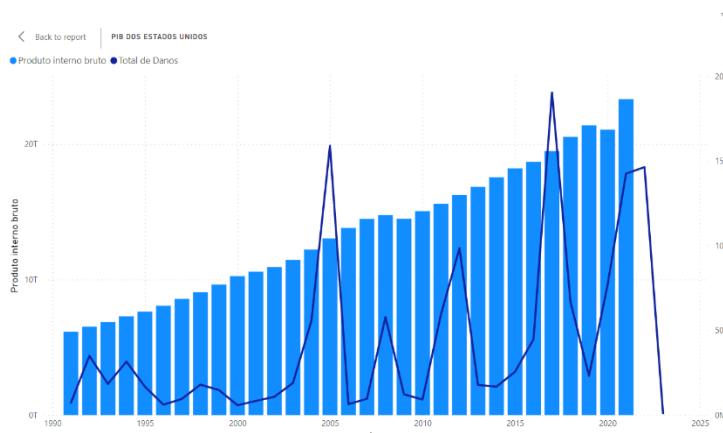


Figura 54. PIB e danos em infraestruturas devido a catástrofes naturais de 1990 a 2023 nos Estados Unidos.

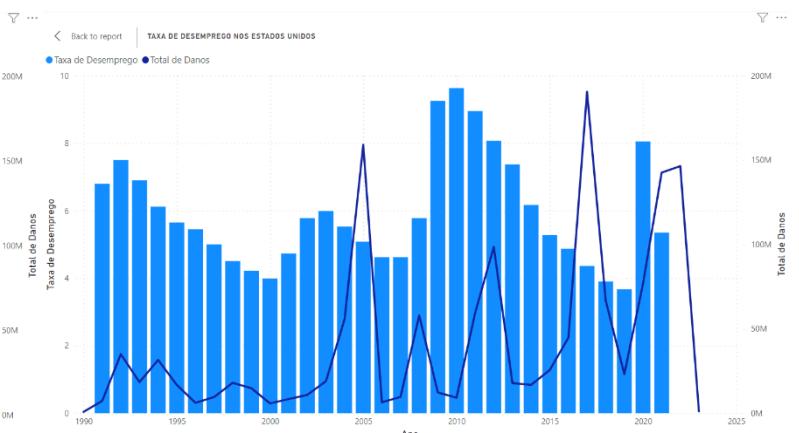


Figura 55. Taxa de desemprego e danos em infraestruturas devido a catástrofes naturais de 1990 a 2023 nos Estados Unidos.



Figura 56. PIB e danos em infraestruturas devido a catástrofes naturais de 1990 a 2023 na China.

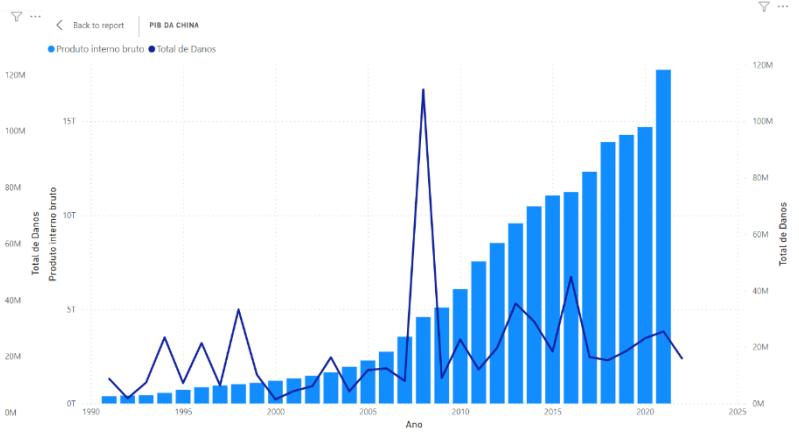


Figura 57. Taxa de desemprego e danos em infraestruturas devido a catástrofes naturais de 1990 a 2023 na China.

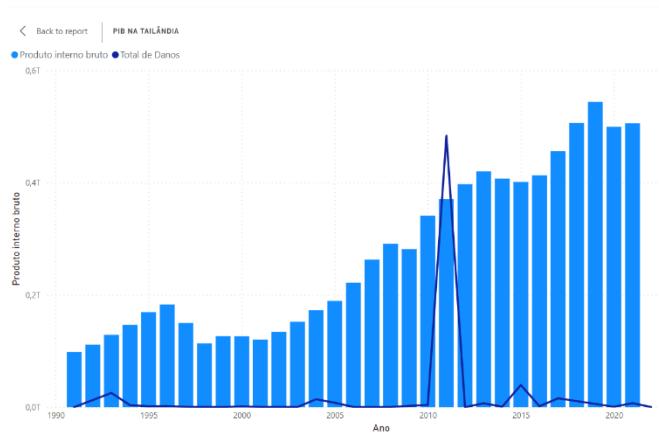


Figura 58. PIB e danos em infraestruturas devido a catástrofes naturais de 1990 a 2023 na Tailândia.

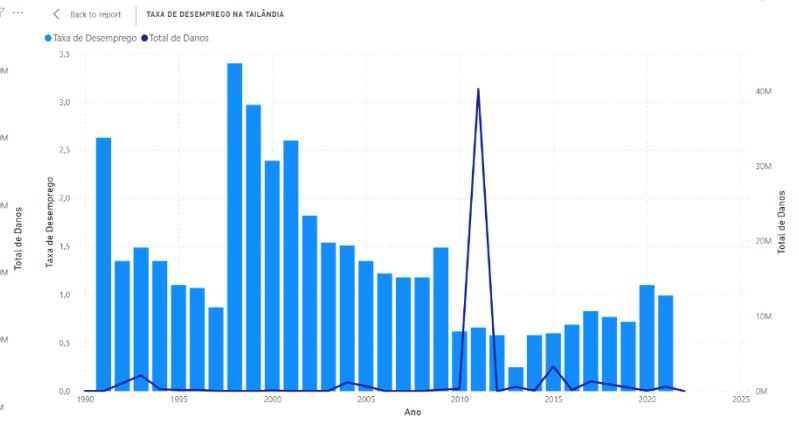


Figura 59. Taxa de desemprego e danos em infraestruturas devido a catástrofes naturais de 1990 a 2023 na Tailândia.

Os restantes países representados nas figuras 52 a 59, são referentes aos desastres com maiores danos nas infraestruturas. Numa análise geral e de forma semelhante aos casos anteriores, os picos que representam os danos não parecem possuir uma relação com o produto interno bruto e com a taxa de desemprego. Por exemplo, no caso do Japão para a taxa de desemprego (figura 53), é possível verificar um decréscimo da taxa de desemprego logo após o desastre, sendo este comportamento o oposto do expectável. No caso da China, existe um aumento do desemprego, mas este já tinha vindo a crescer desde 1990, o que uma vez mais não se traduz num padrão que se possa retirar conclusões. No caso do PIB, não se encontrou nenhum vale aquando da ocorrência dos desastres, permitindo concluir que os desastres naturais não têm qualquer tipo de impacto nos casos apresentados. Como referido nos exemplos anteriores, as catástrofes locais, que poderiam ter um impacto na taxa de desemprego da região afetada e obviamente no PIB, acabam por ser diluídos pelas restantes regiões, resultando numa média sem alterações percutíveis. Assim, para responder a este tipo de perguntas analíticas, seria necessário diminuir o grão da data warehouse para que esta também incluisse o PIB e a taxa de desemprego para as cidades onde ocorram os desastres, sendo assim impossível responder às questões efetuadas de uma forma mais correta.

Conclusão:

A primeira etapa incluiu a análise de fontes de dados abertas e a edição de dados existentes, como a exclusão de colunas e de observações irrelevantes para o negócio proposto. Este processo foi fundamental para aprimorar a manipulação de dados e, assim, obter informações relevantes e necessárias para a próxima etapa do projeto. Durante a primeira etapa, houve uma transformação dos dados de forma a tornar mais fácil a sua manipulação e compreensão. Além disso, nessa etapa, foi possível construir um processo de negócio e elaborar 3 questões analíticas relacionadas com o tema do projeto, cujas respostas serão elaboradas na etapa 3 do projeto.

Na segunda etapa, procedeu-se à realização da modelação dimensional para o processo de negócio criado na etapa 1 do projeto. Nesta etapa definiu-se o grão e o tipo de tabela de factos, estabeleceu-se as dimensões adequadas ao negócio, identificou-se medidas numéricas na tabela de factos e, por fim, desenhou-se o diagrama em estrela. O diagrama em estrela ajuda a visualizar as relações entre as tabelas de factos e dimensões criadas e, assim, facilita a compreensão das informações armazenadas.

Foram encontradas algumas dificuldades na construção da dimensão dimLocation, por se tratar de uma dimensão de mudança lenta, em que foi necessário ter extrema cautela na criação de novas colunas, assim como na junção de várias colunas de diversas tabelas diferentes.

Durante a terceira e última fase do projeto, foi estabelecido um processo completo de Extração, Transformação e Carregamento (ETL). Isso foi alcançado por meio da utilização dos scripts desenvolvidos na etapa 1, responsáveis pela transformação dos dados, bem como pela criação de novos scripts para a modelagem dimensional na etapa 2. Além disso, foram implementados scripts adicionais para realizar o carregamento apropriado das dimensões e da tabela de fatos na base de dados do PostgreSQL.

Posteriormente, a base de dados foi integrada na ferramenta do PowerBI de forma a gerar um cubo de dados, que permitiu realizar várias análises, que visaram responder às questões analíticas formuladas na etapa 1. As diversas questões permitiram ter uma ideia geral dos desastres mais graves da história, avaliando o seu impacto na agricultura, turismo, rendimento e desemprego dos países, permitindo explorar o processo de negócio proposto na sua totalidade e dar uso a uma ferramenta de extrema utilidade na integração dos dados de uma data warehouse.

Referências

- 1** - Negri, J. (2021). EMDAT (Emergency Events Database) - The Natural Disasters Dataset. Obtido a 18 Março de 2023, de https://www.kaggle.com/datasets/inegrini/emdat19002021?select=EMDAT_1900-2021_NatDis.csv
- 2** - World Bank. (2023). GDP (current US\$). Obtido a 23 Março de 2023, de <https://data.worldbank.org/indicator/NY.GDP.MKTP.CD>
- 3** - World Bank. (2023). Population, total. Obtido a 23 Março de 2023, de <https://data.worldbank.org/indicator/SP.POP.TOTL>
- 4** - World Bank. (2023). Unemployment, total (% of total labor force). Obtido a 23 Março de 2023, de <https://data.worldbank.org/indicator/SL.UEM.TOTL.ZS>
- 5** - World Bank. (2023). International tourism, number of arrivals. Obtido a 23 Março de 2023, de <https://data.worldbank.org/indicator/ST.INT.ARVL>
- 6** - Food and Agriculture Organization of the United Nations. (2023). FAOSTAT online database. Obtido a 24 Março de 2023, de <https://www.fao.org/faostat/en/#data/QI>