

Desenvolvimento de um Data Mart Para Apoiar a Análise Quanto ao Investimento na Implantação de Cursos Pré-Vestibular em SC

Augusto Pamplona, Eduardo Becker, Pablo Vicente, Roberto Rivelino, Rolf Zambon

Curso de Bacharelado em Sistemas de Informação - Universidade Federal de Santa Catarina (UFSC)

{augusto.pamplona, eduardo.becker, pablo.vicente, roberto.rvs, rolf.zambon}@grad.ufsc.br

Abstract. *The main objective of this Data Warehouse project is to search for the most prone and advantageous regions of Santa Catarina for the implementation of pre-university entrance exam courses, with the purpose of initiating an investor decision decision. To this end, we modeled a data modeling scheme to answer the questions: What is the socio economic profile of the city / region? What is the age group by city / region? How many times has a candidate taken an entrance exam in the city / region? Which cities / regions have the most candidates in very busy courses? To consult and answer these questions, use the data available from Coperve, a database of entrance exam information from 2008 to 2012 and the candidate / vacancy ratio for 2012 also available from Coperve. The Kettle tool will be used to perform the ETL process and PowerBi will implement the front end (dashboard).*

Resumo. *O objetivo principal deste projeto de Data Warehouse é buscar quais regiões/cidade em Santa Catarina são mais propensas e vantajosas para implantação de cursos pré-vestibular, com objetivo de apoiar a tomada de decisão de investidores. Com este objetivo, modelamos um esquema estrela (metodologia de modelagem de dados) para responder às perguntas: Qual o perfil socioeconômico da cidade/região?; Qual a faixa etária por cidade/região?; Qual a quantidade de vezes que um candidato prestou vestibular por cidade/região?; Quais cidades/regiões possuem mais candidatos em cursos muito concorridos?; Para realizar as consultas e responder tais questões, utilizamos os dados disponibilizado pela Coperve, um banco de dados com as informações dos vestibulares entre os anos 2008 a 2012 e relação candidato/vaga do ano 2012 também disponibilizado pela Coperve. Será utilizado a ferramenta Kettle para*

realizar o processo de ETL e o PowerBi para implementar o front-end (dashboard).

1. Introdução

Para alunos que cursam os anos finais do ensino médio sempre surge o desejo de realizar uma graduação. Para ingressar em uma instituição de ensino superior é necessário que os alunos realizem provas pelas quais são avaliados os conhecimentos adquiridos durante o ensino médio. A principal forma de acesso a uma universidade pública é através da prova do vestibular, em algumas instituições elas podem ocorrer duas vezes ao ano em quantos em outras são somente realizadas uma vez ao ano.

Para a realização de um vestibular a universidade divulga um edital com todas as informações necessárias para o ingresso em curso de graduação, neste edital são divulgadas as datas de abertura e fechamento do período de inscrição. A etapa de inscrição é uma das etapas mais importante do vestibular, pois os alunos terão que escolher qual curso desejarem tentar uma vaga como também para universidade que calcularam a relação de candidatos por vaga.

Após o período de inscrição são divulgadas a relação de candidatos dos cursos, este indicador é muito importante para alunos terem como mensurar quais são suas chances de conseguirem a vaga bem como o quanto de dedicação aos estudos será necessário para a prova. Estes dados são divulgados pela universidade para todos os vestibulares realizados.

Também após cada processo seletivo, são obtidos dados sobre candidatos às vagas, suas informações socioeconômicas, características, notas e aprovações. A partir dessas informações, dos anos de 2008 à 2012, disponibilizadas para desenvolvimento deste trabalho, decidimos analisar quais regiões/cidades em Santa Catarina são mais propensas e vantajosas para implantação de cursos pré-vestibular.

Para o desenvolvimento deste trabalho, primeiramente analisamos as informações contidas no banco de dados disponibilizado, verificando se tínhamos informações suficientes para realização das análises, vimos que precisávamos de outra informação que era a relação candidato/vaga, então através do site (<http://www.vestibular2012.ufsc.br/index.php>) acessamos o pdf disponibilizado com tal informação e o convertimos para tabelas criando uma faixa de concorrência. Depois foi feita a criação do esquema estrela, composto por quatro dimensões e um fato. Após definição da modelagem de dados utilizamos o Kettle para realização do processo de ETL e por fim utilizamos o power BI para gerar os resultados, apresentados por gráficos.

2. Materiais

O material utilizado neste projeto foi disponibilizado pela Coperve (Comissão Permanente do Vestibular) da UFSC. Foi utilizado o banco de dados entre os anos 2008 a 2012 para obtenção de informações que possibilitem a execução do projeto. O Programa de Ações Afirmativas criado pela Resolução Normativa No008/CUn/2007, criado em Jul/2007, estabeleceu que para o vestibular de 2008, foi destinado **30%**

(trinta por cento) das vagas do vestibular, em cada curso, que foram distribuídas da seguinte forma:

20% foram destinadas a candidatos que cursaram o ensino médio em escolas públicas;

10% foram destinadas a candidatos autodeclarados pretos ou pardos que cursaram o ensino médio em escolas públicas;

Além de serem criadas **5** vagas suplementares que foram preenchidas pelos candidatos com melhor classificados no vestibular destinadas a candidatos pertencentes a povos indígenas,

Algumas mudanças ocorreram no vestibular de 2012, a partir da nova resolução normativa N.º 26/CUn/2012, mudando a distribuição das vagas, sendo assim **30%** das vagas de cada curso foram destinadas a candidatos que cursaram o ensino médio em escolas públicas, com a seguinte distribuição:

20% foram destinadas a candidatos que cursaram o ensino médio em escolas públicas sendo esse distribuída entre:

50% para aqueles com renda familiar bruta mensal igual ou inferior a 1,5 salários mínimo per capita; Dessas, **32%** das vagas destinadas a candidatos autodeclarados pretos, pardos ou indígenas; **68%** para os demais;

50% para aqueles com renda familiar bruta mensal superior a 1,5 salário médio per capita; Dessas, **32%** das vagas destinadas a candidatos autodeclarados pretos, pardos ou indígenas; **68%** para os demais;

10% foram destinadas a candidatos autodeclarados pretos ou pardos que cursaram o ensino médio em escolas públicas;

Além de serem criadas **10** vagas suplementares que foram preenchidas pelos candidatos com melhor classificados no vestibular destinadas a candidatos pertencentes a povos indígenas,

O banco de dados conta com informações socioeconômicas dos candidatos às vagas na universidade que são preenchidas via formulário no momento da inscrição. O formulário conta com perguntas como: Estado civil; Frequentou ou frequenta curso pré-vestibular; Nível de instrução de seu pai; Nível de instrução da sua mãe; Tipo de estabelecimento onde cursou o ensino médio; Indique o principal responsável pelo sustento de sua família; Meio de transporte que mais utiliza; Possui computador em sua residência; Número de vezes que prestou vestibular para UFSC.

Também foi utilizado a relação de candidato vaga do ano 2012 também disponibilizado pela Coperve em formato de pdf, disponível no site <http://www.vestibular2012.ufsc.br/index.php>. O vestibular 2012 contou com 30.388 candidatos inscritos.

Utilizamos essas informações socioeconômicas sobre os candidatos para traçar um perfil e através disso obter quais regiões são mais propensas a receberem um candidato que possivelmente realizará uma matrícula em um cursinho pré-vestibular, assim mapeando quais regiões/cidades possuem maior prospectos.

Além destes dados, buscamos, para cada cidade, o seu PIB per capita. Os dados foram retirados do IBGE, mais precisamente da tabela de PIB dos municípios. disponível no site

<https://www.ibge.gov.br/estatisticas/economicas/contas-nacionais/9088-produto-interno-bruto-dos-municipios.html?t=resultados>. Utilizamos os dados de 2016 e dividimos o PIB em três faixas, baixo, médio e alto.

3. Métodos

Para implementação deste trabalho, foi proposta a realização de um Data Mart para a solução do problema. Data Mart é uma subdivisão ou subconjunto de um Data Warehouse. Os Data Marts são como pequenas fatias que armazenam subconjuntos de dados, normalmente organizados para um departamento ou um processo de negócio.

Data Warehouse é um repositório de dados digitais integrado, não volátil e variável em relação ao tempo que serve para armazenar informações detalhadas, criando e organizando relatórios através de históricos que são depois usados pela empresa, ou organização, para ajudar a tomar decisões importantes com base nos fatos apresentados.

Através dos dados obtidos nos materiais disponibilizados pela Coperve, além dos recolhidos externamente, foi possível criar um Data Mart, que possibilitou o recolhimento e a análise de informações específicas para o nosso problema, disponibilizando uma maior flexibilidade nas pesquisas de informações necessárias. Além de manter um histórico de informações, o Data Mart cria padrões melhorando os dados analisados, corrigindo os erros e reestruturando os dados sem afetar o sistema de operação, apresentando somente um modelo final e organizado para a análise em questão.

4. Metodologias

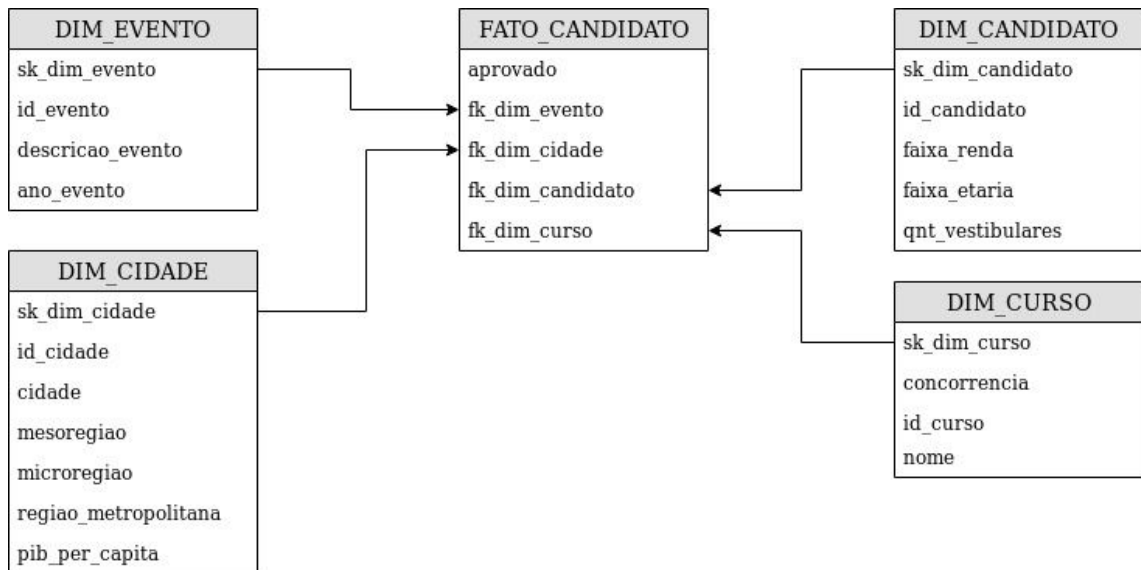


Figura 1. Esquema estrela

Para chegarmos à este esquema estrela foi necessário “passar” por quatro pontos essenciais, sendo eles:

Definir os processos: aqui desenvolvemos algumas perguntas chaves aos quais gostaríamos de conseguir responder com o produto final. Que foram, “Qual o perfil socioeconômico da cidade/região?”, “Qual a faixa etária por cidade/região?”, “Qual a quantidade de vezes que um candidato prestou vestibular por cidade/região?” e por fim “Quais cidades/regiões possuem mais candidatos em curso muito concorridos?”.

Definir o grão: Nossos dados são atualizados anualmente, ou seja, ficamos com um grão bastante grande.

Definir as dimensões:

Dimensão evento: foi aquela utilizada para definir o nosso grão, também conhecida como dimensão tempo, ou seja, relaciona cada entrada do fato com um vestibular.

O ETL da dimensão evento foi muito simples, pegamos o id_evento, ano_evento e descricao_evento, e mapeamos para nossa dimensão.



Figura 2. ETL da dimensão evento

Dimensão cidade: descreve a cidade onde o candidato mora no momento em que prestou o vestibular. Classificamos cada cidade em três hierarquias, mesorregião, microrregião e região metropolitana. Totalizam 293 chaves contendo cidades pertencentes a Santa Catarina que foram analisadas. Por fim buscamos, para cada cidade, o seu pib per capita. Os dados foram retirados do IBGE, mais precisamente da tabela de PIB dos municípios, disponível no link <https://www.ibge.gov.br/estatisticas/economicas/contas-nacionais/9088-produto-inter-no-bruto-dos-municipios.html?t=resultados>. Utilizamos os dados de 2016, e dividimos o PIB em três faixas, baixo (0 à 24.999,00), médio (25.000,00 à 39.999,99) alto (40.000,00 à 100.000,00).

Durante o processo de ETL percebemos que há inconsistências no campo do nome das cidade, sendo elas:

- S JOSÉ -> SÃO JOSÉ
- SÃO BENTO SUL e S BENTO SUL -> SÃO BENTO DO SUL
- BAL CAMBORIÚ -> BALNEÁRIO CAMBORIÚ
- RIO SUL -> RIO DO SUL

Como foram poucas inconsistências, demos alguns UPDATEs na tabela candidatos para corrigir as cidades.

Depois o ETL seguiu normalmente. Buscamos a planilha do IBGE, preparamos os dados, fizemos a faixa_pib e fizemos um LEFT JOIN com as cidades do nosso banco, e carregamos na nossa dimensão cidade.

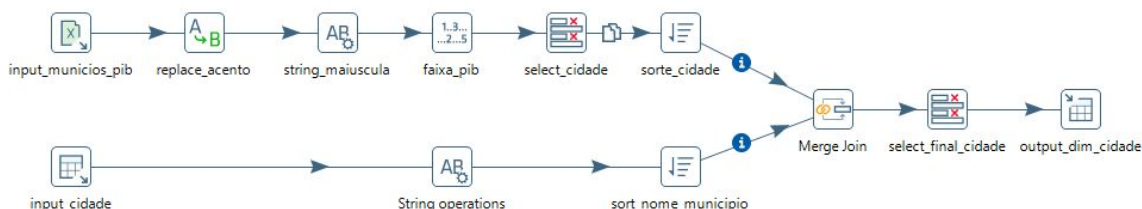


Figura 3. ETL da dimensão cidade

Dimensão curso: aqui tratamos os cursos, onde deles analisamos a concorrência, baseados na relação candidatos por vaga retirados de http://www.vestibular2012.ufsc.br/relacao_candidatosVaga.pdf, e assim criamos a faixa de concorrência, separadas em três níveis, baixa (de 0 até 3,99), média (4 até 9,99) e alta (à partir de 10). Guardamos também o id do curso, para criar as relações entre as tabelas necessárias, e o nome do curso.

Como pegamos a relação C/V dos cursos de 2012 (o mesmo ano do vestibular), o ETL ocorreu tranquilamente, trazendo as informações C/V da planilha, preparando os dados, fizemos um LEFT JOIN com os cursos do nosso banco, criamos a faixa_concorrencia e carregamos na nossa dimensão Curso.

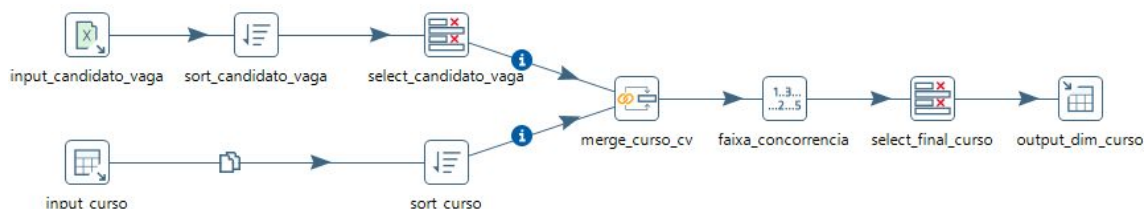


Figura 4. ETL da dimensão curso

Dimensão candidato: para cada entrada no fato há uma entrada na dimensão candidato, onde criamos faixas para caracterizá-los, sendo elas, faixa de renda que foi separada em cinco níveis, muito baixa (Até 1 salário mínimo), baixa (Acima de 1 até 5 salários mínimos), média (Acima de 5 até 10 salários mínimos), alta (Acima de 10 até 30 salários mínimos) e muito alta (Acima de 30 salários mínimos); Faixa etária, utilizamos como base a utilizada no enem (menor que 16, igual a 16, igual a 17, igual a 18, igual a 19, igual a 20, de 21 a 30, de 31 a 59 e maior ou igual a 60), e por fim a quantidade de vestibulares que o candidato já realizou também em faixas (nenhum, um, dois, três e quatro ou mais).

Para o ETL da dimensão candidato foi necessário extrair a renda e a quantidade de vestibulares do questionário_socioeconomico. Para isso analisamos as planilhas do socioeconômico, e utilizamos a função js e faixas para classificar a renda e a quantidade de vestibulares, depois calculamos a idade pegando o ano da data de nascimento do candidato e diminuindo do ano_evento, aplicamos a faixa_etaria e carregamos na dimensão candidato.

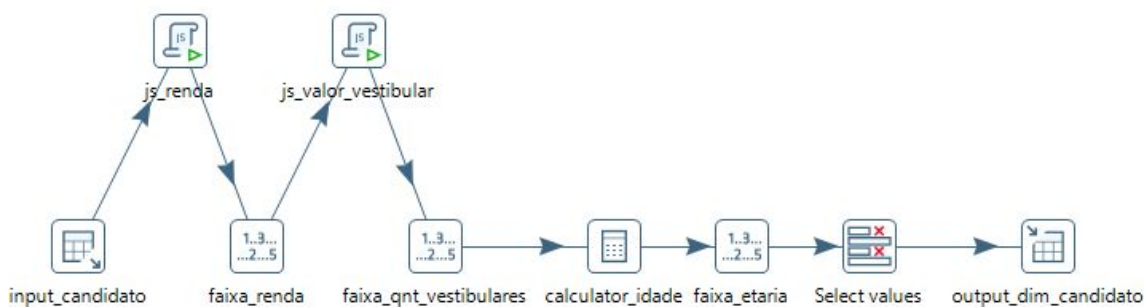


Figura 5. ETL da dimensão candidato

Definir o fato mensurável: **Fato candidato** diferente da dimensão candidato, o fato é relacionado com todas as outras dimensões, nos possibilitando assim fazer as análises desejadas entre os dados disponíveis.

Também criamos um boolean *aprovado* que tem como resposta Y para aquele candidato que foi aprovado no vestibular e N para aquele que não foi aprovado. Para inferirmos se o candidato foi aprovado ou não, criamos um step chamado js_id_curso que nos permite criar um script na linguagem javascript, nele utilizamos o campo codigo_opcao presente na tabela candidato_classificado. Caso o codigo_opcao seja diferente de nulo, o seu valor é adicionado no campo id_curso e o campo aprovado

recebe Y, caso contrário o id_curso recebe o valor do campo id_opcao1 e o campo aprovado recebe N. Dessa forma também conseguimos identificar o curso que o candidato se classificou e caso não tenha sido classificado, o curso na qual se inscreveu em primeira opção.



Figura 6. Step js_id_curso

Para relacionar as tabelas dimensão a partir de suas serial keys (sk), realizamos uma série de joins a partir dos identificadores das tabelas do database vestibular com os identificadores herdados das tabelas dimensão, ordenando-os pelos steps Sort rows que estão presentes na transformação. No final, utilizamos o step select_final_fato_candidato para retornar somente os campos necessários para popular o fato. A transformação completa pode ser vista na Figura 7.

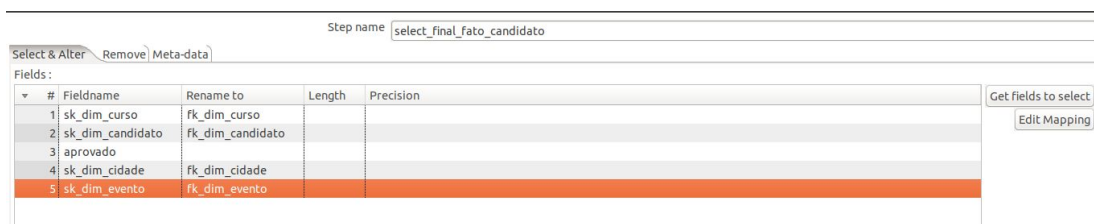


Figura 7. Step select_fato_candidato

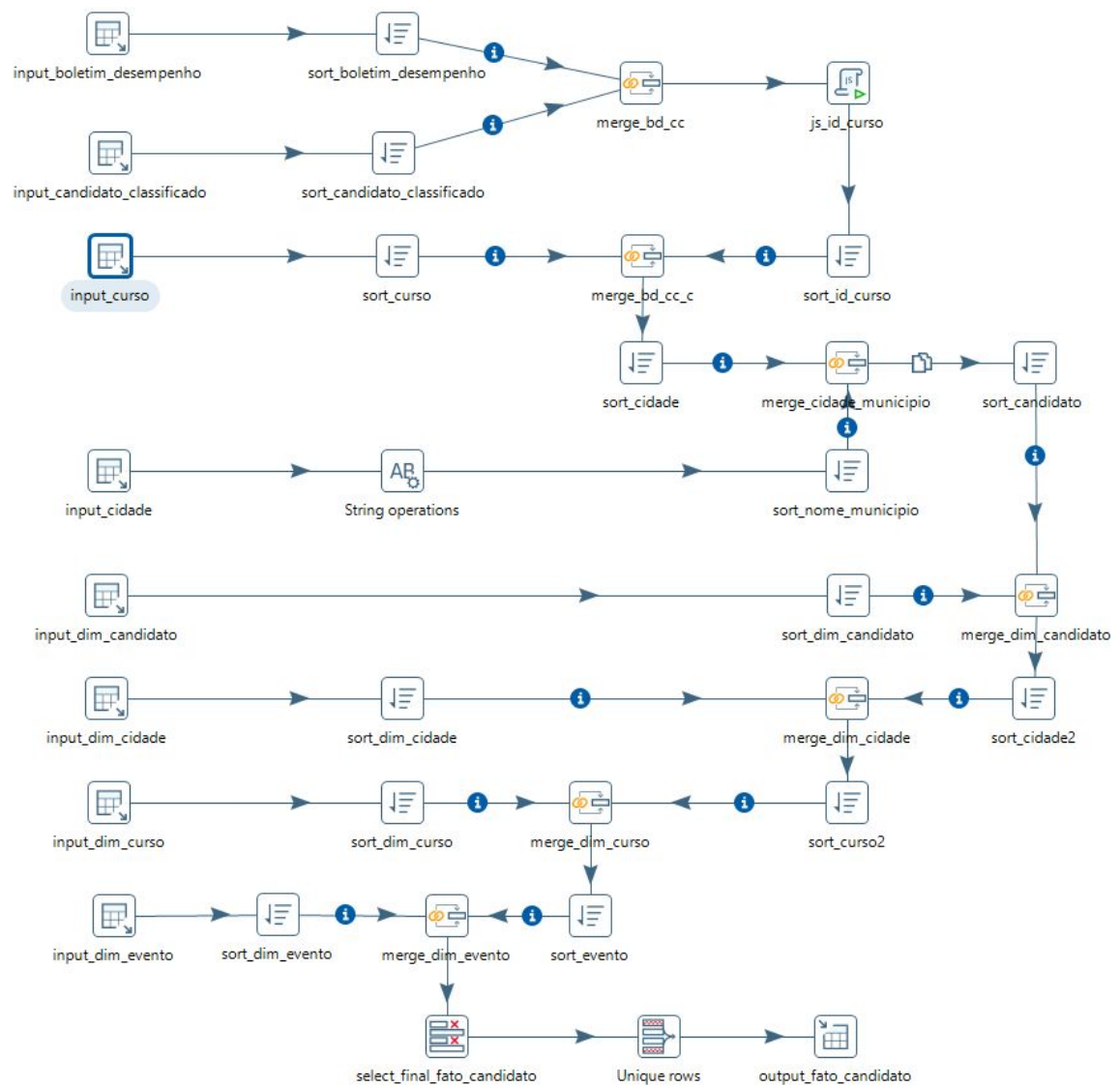


Figura 8. ETL do fato candidato

Por fim, criamos um job para executar as transformações desenvolvidas.

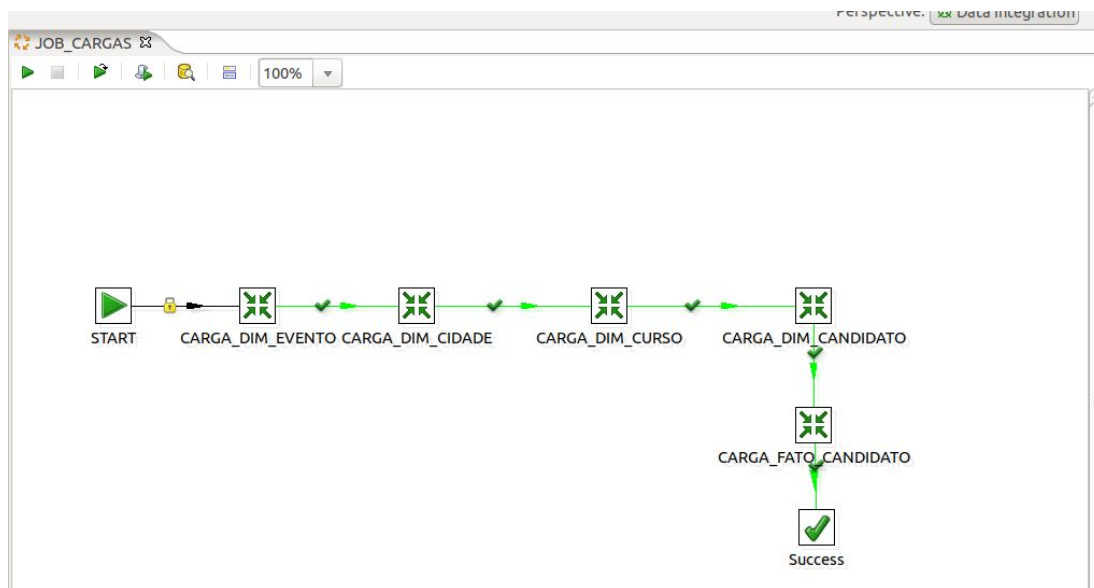


Figura 9. Job de cargas.

5. Resultados

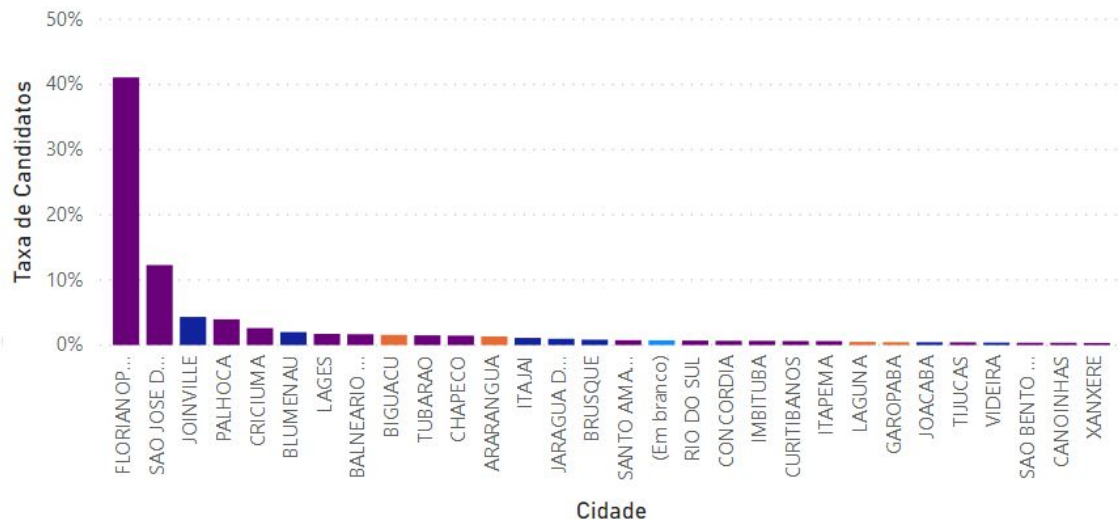
Quatro primeiras imagens foram onde conseguimos os dados para responder a nossa primeira questão:

- Qual o perfil socioeconômico da cidade/região?

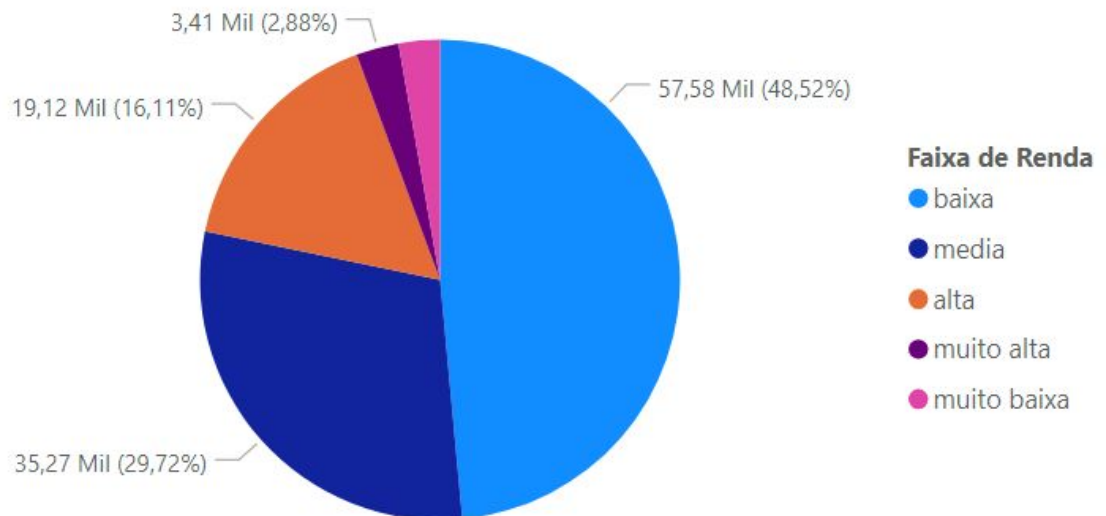
Concluimos que regiões mais “desenvolvidas” possuem um maior número de candidatos, e o que foi uma surpresa, o maior número de candidatos ao vestibular são aqueles caracterizados por baixa renda, e os menos interessados são aqueles de renda muito alta e muito baixa.

Taxa de Candidatos por Cidade e PIB

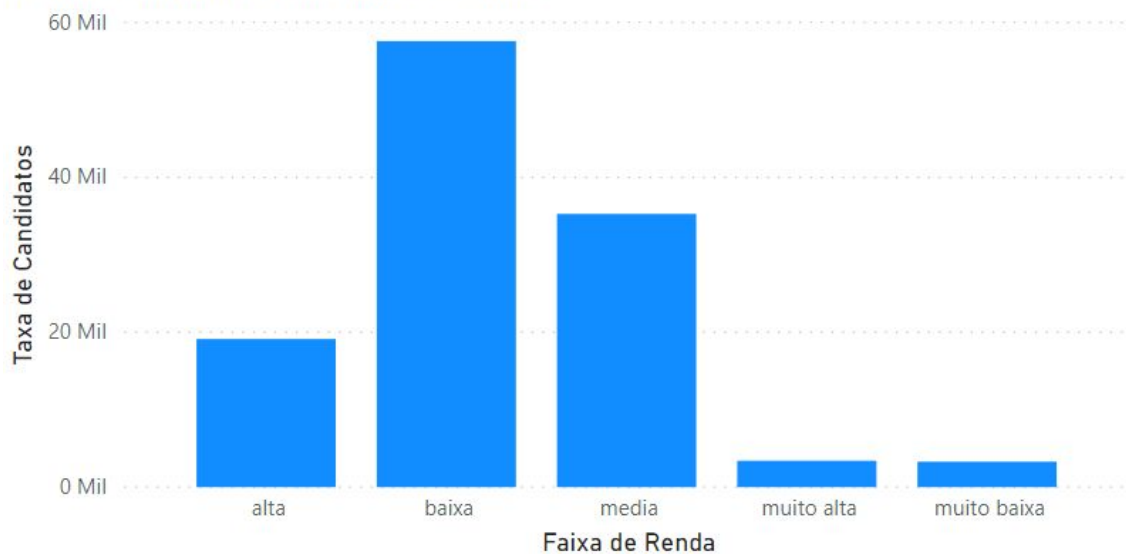
PIB ● (Em branco) ● alto ● baixo ● medio ● unknown



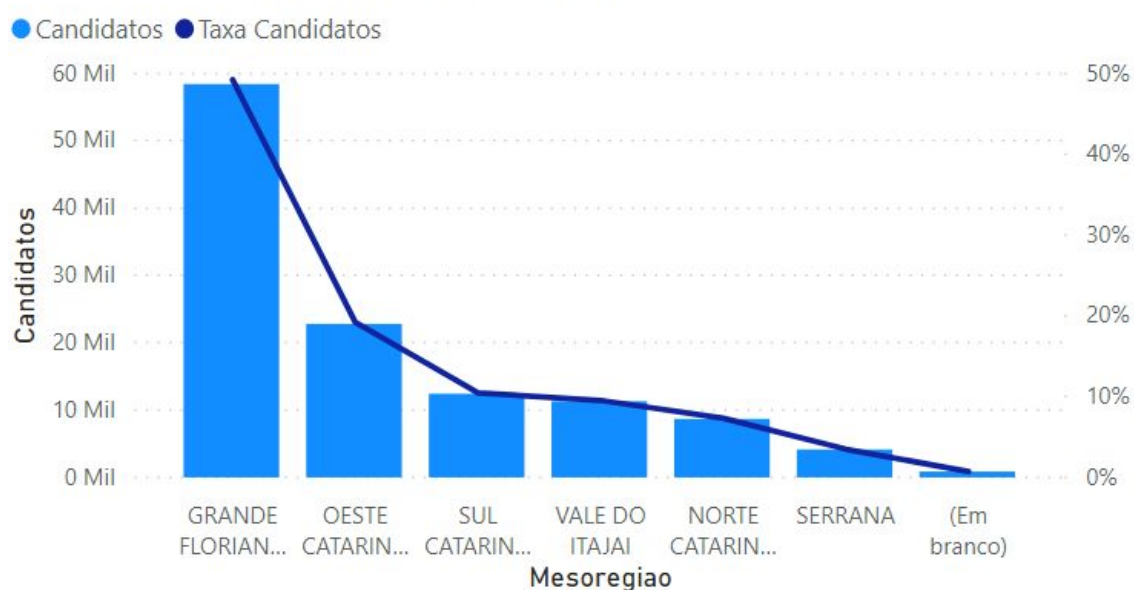
Candidatos por Faixa de Renda



Taxa de Candidatos por Faixa de Renda



Candidatos e Taxa Candidatos por Mesoregiao

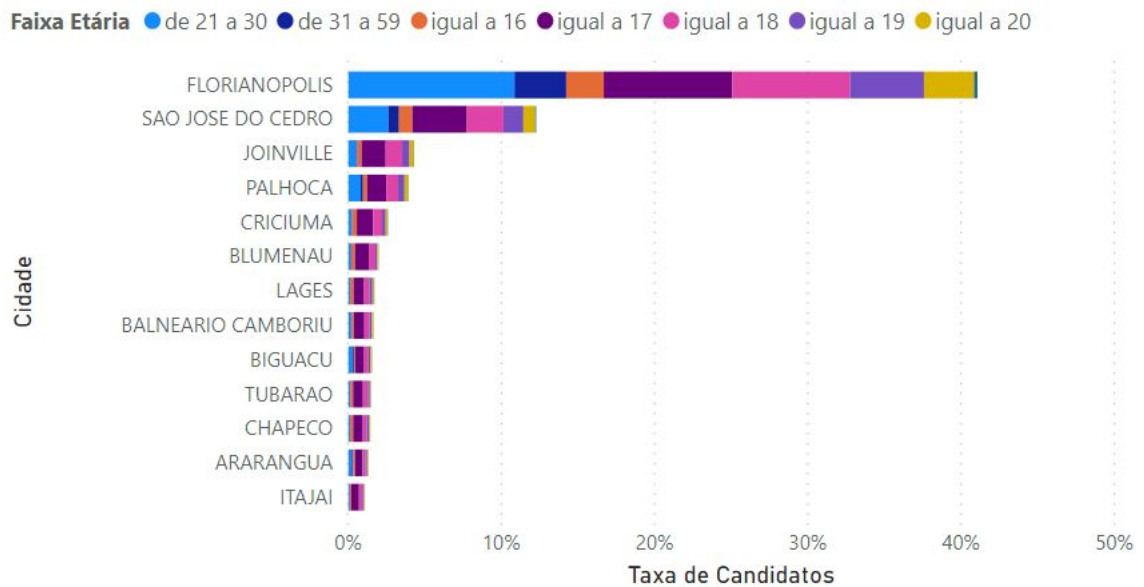


As próximas quatro imagens foram utilizadas para responder a nossa segunda questão:

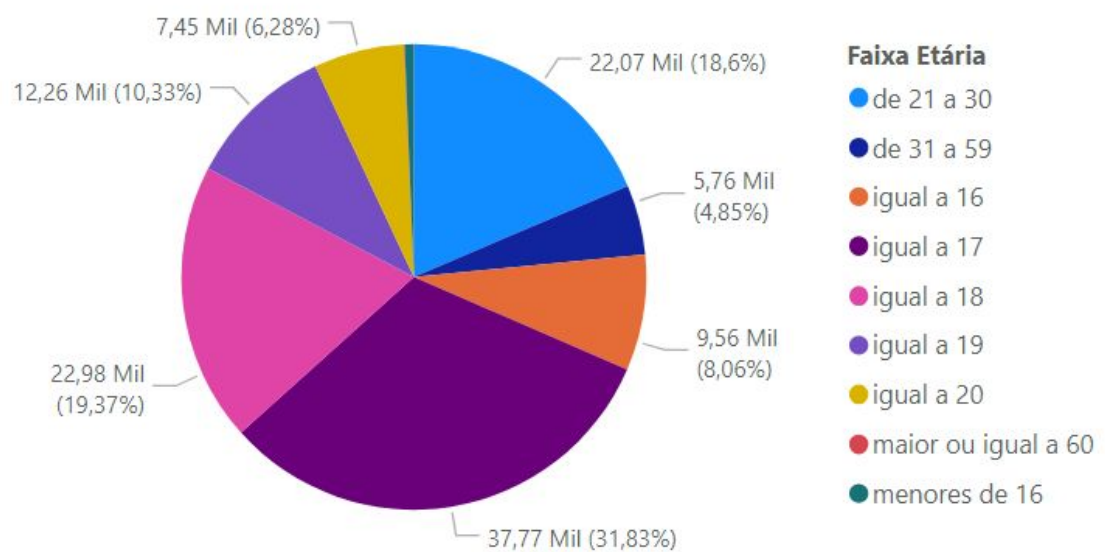
- Qual a faixa etária por cidade/região?

Concluimos que, tirando florianópolis que tivemos um resultado diferente das demais, a média de idade para os candidatos ficou muito forte para candidatos de 17 anos, em seguida a faixa entre 21 e 30 anos, e a menor, que quase não é perceptível a aparição nos resultados é a faixa de maior ou igual a 60 anos.

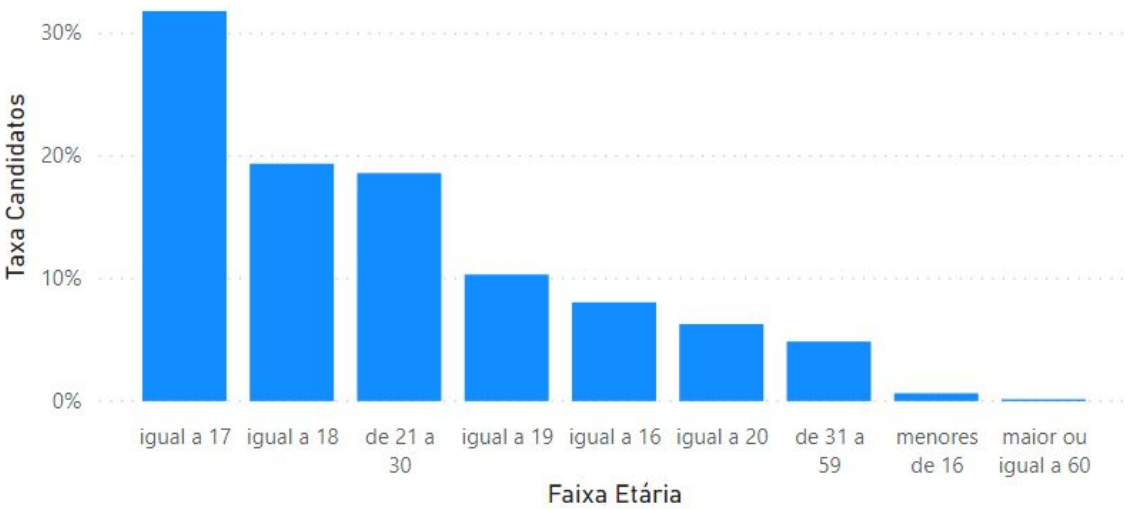
Taxa de Candidatos por Cidade e Faixa Etária



Candidatos por Faixa Etária



Taxa Candidatos por Faixa Etária



Candidatos e Taxa Candidatos por Mesoregiao

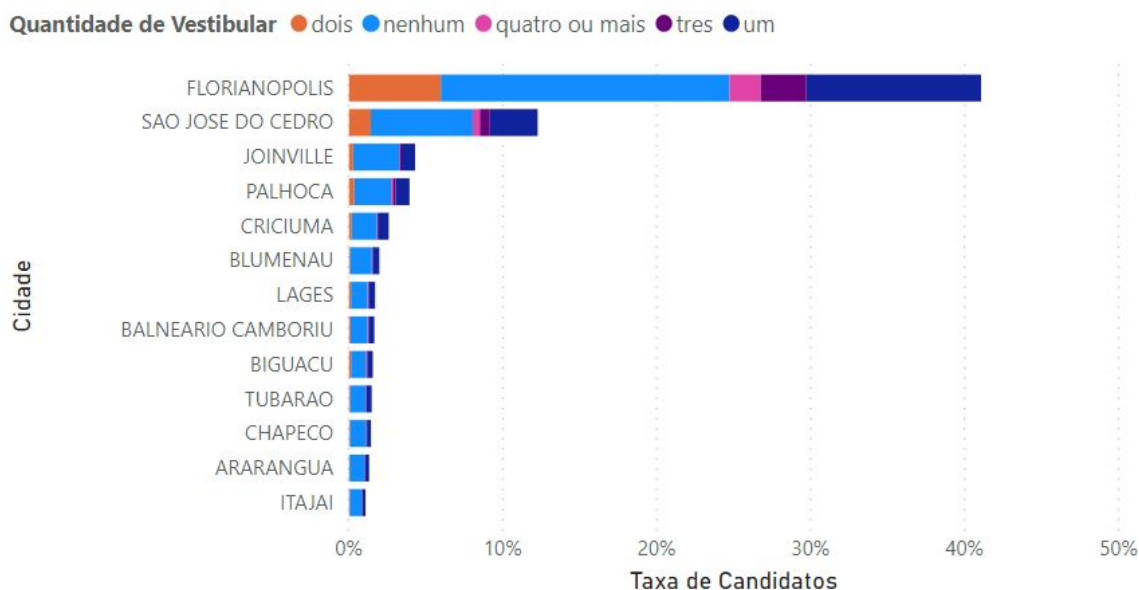


Os próximos 4 gráficos foram utilizados para responder a nossa terceira questão:

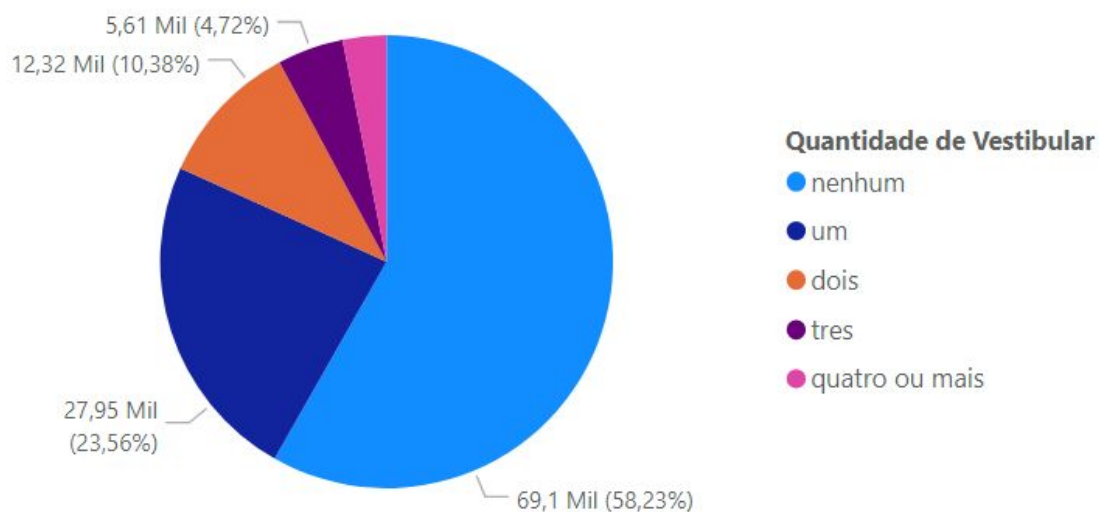
- Qual a quantidade de vezes que um candidato prestou vestibular por cidade/região?

Sem dúvida alguma, obtemos como resposta para todas as regiões que o número de candidatos que não realizaram nenhum vestibular anteriormente é disparado o maior, em seguida vêm aqueles que já realizaram um vestibular, o valor entre esses dois pontos resultantes ficou mais parecido para Florianópolis, para as outras regiões a diferença fica maior.

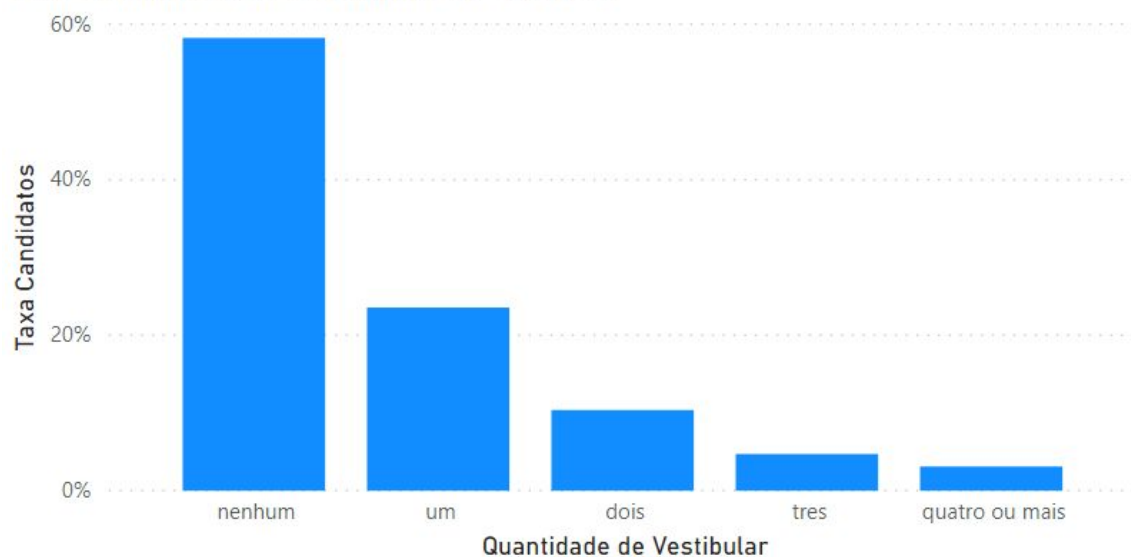
Taxa de Candidatos por Cidade e Quantidade de Vestibular



Candidatos por Quantidade de Vestibular



Taxa Candidatos por Quantidade de Vestibular



Candidatos e Taxa Candidatos por Mesoregiao



Os 3 próximos gráficos foram utilizados para responder a nossa última questão:

- Quais cidades/regiões possuem mais candidatos em cursos muito concorridos?

Percebemos que baseado no número de candidatos, as cidades com mais candidatos obtiveram valores parecidos para os três níveis de concorrência, para as demais cidades há um valor mais semelhante para os cursos de alta e média concorrência. Também tivemos como resultado o valor de aprovados para cada cidade, para este

ponto, pegamos aqueles aprovados por segunda opção de curso, que mais tarde isso será analisado com maior foco, notamos que para toda cidade analisada, o número de aprovados é muito inferior ao número de reprovados. Por fim, analisamos os cursos com seus candidatos e o nível de concorrência atribuídos a eles, aqui o dados dos não aprovados foi para a primeira opção de curso selecionada na inscrição, e para os aprovados, o curso para qual foram, não importante qual das opções era, nos mostrando que os cursos mais concorridos são aqueles que têm o maior número de concorrência, fugindo desse resultado apenas em casos muitos específicos, onde um curso bastante concorrido tem baixa relação de candidato por vaga.

Taxa de Candidados por Cidade e Concorrência

Concorrência (Em branco) alta baixa media



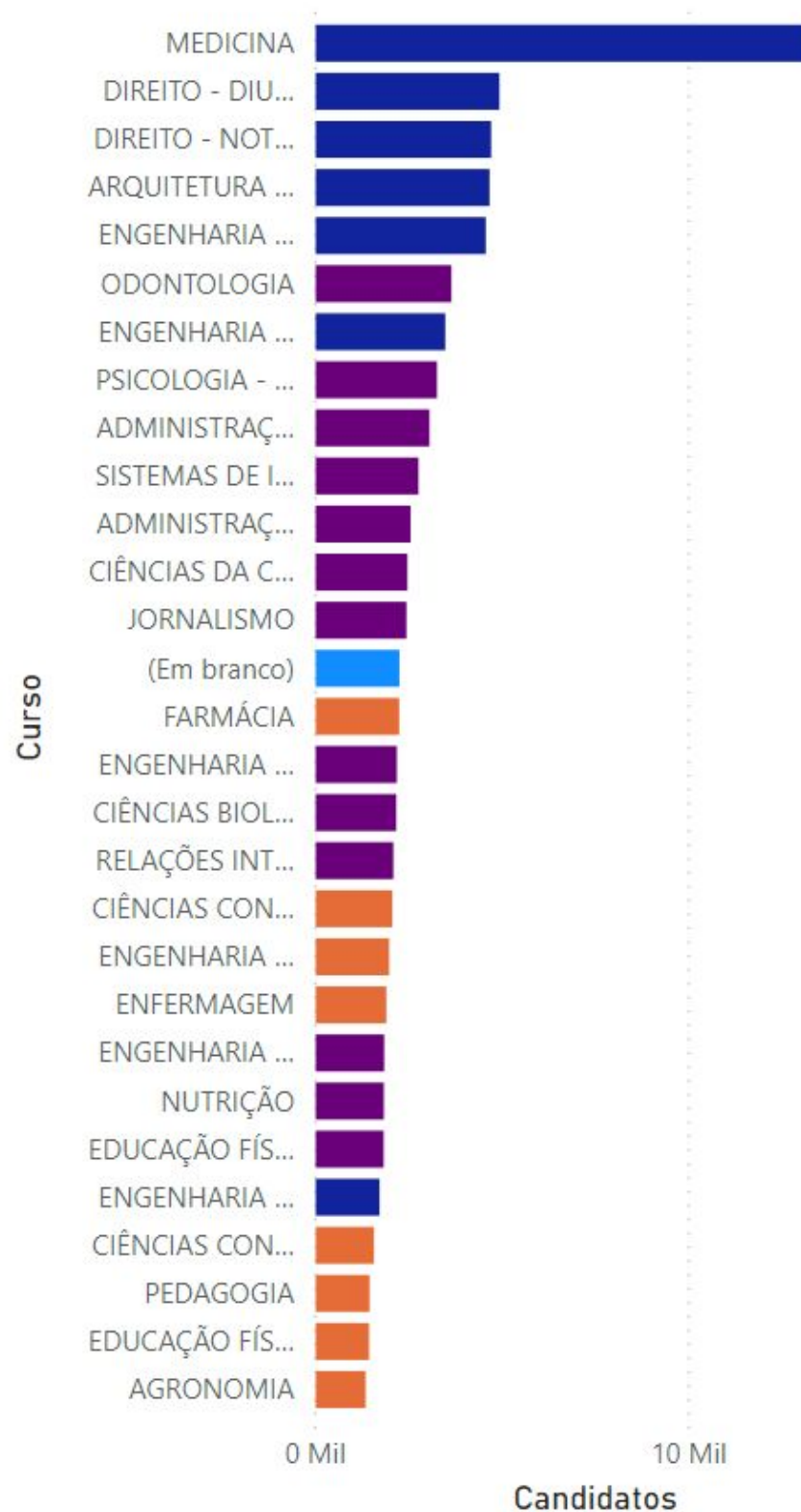
Taxa de Candidatos por Cidade e Aprovação

Aprovação ● N ● Y



Candidatos por Curso e Concorrência

Concorrência ● (Em branco) ● alta ● baixa ● media



6. Conclusões e trabalhos futuros

Através dos resultados obtidos nesta análise e estudo de caso com informações de vestibular UFSC, podemos concluir que nos anos de 2008 à 2012, a cidade de Florianópolis mostrou ser a cidade mais promissora para investir-se em um curso pré-vestibular, tendo em vista os resultados apresentados no tópico anterior, tais resultados refletem em se tratar da capital do estado e ser residência do maior campus pertencente a UFSC. Não muito atrás, a cidade São José também se destacou nos resultados, que pode ser atribuído à sua proximidade com Florianópolis e por ser sua cidade satélite com maior densidade demográfica. Outras cidades que também mostraram bons indicadores foram Joinville, Blumenau, Lages e Balneário Camboriú, que podem ter mostrado menor número de candidatos comparado com as outras já citadas, porém mostraram uma grande proporção de candidatos que estão tentando cursos que classificamos como concorrência “alta”, sendo esse um dos fatores que avaliamos ser o de maior peso quanto a possível ingresso do candidato a um curso pré-vestibular.

Ressaltando que este projeto ficou restrito apenas a uma análise de mercado em busca de possíveis prospectos, sendo avaliado o mercado baseado, através dos dados dos candidatos, perfil socioeconômico das regiões, faixa etária dos inscritos por região, quantidade de vezes que um candidato já prestou vestibular e a tendência de candidatos que prestam cursos que são mais concorridos por região. Não foi realizado, neste trabalho, uma análise a fundo sobre a economia regional e incentivos ou não para implantação de um curso pré-vestibular, como também não analisamos o aluguel médio das regiões.

Como este projeto ficou restrito aos dados da Coperve, somente aos vestibulares ocorridos entre os anos de 2008 a 2012. Futuramente o projeto pode ser expandido para um estudo maior de anos, tendo em vista que o último vestibular analisado foi de 7 anos atrás, isso trará uma análise mais próxima da realidade atual e uma maior gama de dados que podem ser analisados.

A análise foi feita levando em consideração apenas o vestibular da UFSC e candidatos pertencentes ao estado de Santa Catarina. Não vemos nenhuma restrição em expandir nossa análise englobando outras universidades e faculdades, isso trará uma análise mais complexa e detalhada, podendo assim comparar instituições umas com as outras e analisar quais, segundo os indicadores, possuem maior número de possíveis prospectos a cursos pré-vestibular. Podendo assim elaborar a criação de um Data Warehouse unificado nacionalmente com o perfil dos estudantes.

7. Referências bibliográficas

| | | | |
|---|-----|-----|--------------|
| Para | uso | dos | indicadores: |
| https://moodle.ufsc.br/mod/resource/view.php?id=1920024 . | | | Acessado em |
| Novembro. | | | |

Para o uso do material dos vestibulares:
<https://moodle.ufsc.br/mod/url/view.php?id=1920036>. Acessado em Novembro.

Excel usados para o consumo dos dados desejados:
<https://moodle.ufsc.br/mod/url/view.php?id=1920040>. Acessado em Novembro.

Consumo dos dados desejados para relação candidato/vaga 2012:
<http://www.vestibular2012.ufsc.br/index.php>. Acessado em Novembro.

Para análise de vestibulares anteriores:
<https://coperve.ufsc.br/vestibulares-anteriores/>. Acessado em Novembro.

Consumo dos dados desejados para PIB e regiões dos municípios:
<https://www.ibge.gov.br/estatisticas/economicas/contas-nacionais/9088-produto-inter-no-bruto-dos-municipios.html?t=resultados>. Acessado em Novembro.

Base para formação do artigo:
<http://www.sbc.org.br/documentos-da-sbc/category/169-templates-para-artigos-e-capitulos-de-livros>. Acessado em Novembro.

Além das aulas administradas pelo Prof. José Leomar Todesco, seu material de apoio disponibilizado para matéria INE5643-07238 - Data Warehouse para apoio do embasamento teórico: <https://moodle.ufsc.br/course/view.php?id=108814> Acessado em Novembro.