

# ENEM na Rede: a social, online and free environment to support Brazilian students in knowledge creation.

Marcos Arrais  
Programa de Pós-graduação em Informática  
Universidade Federal do Rio de Janeiro  
Rio de Janeiro, RJ  
marcos.arrais@gmail.com

Jonice Oliveira  
Programa de Pós-graduação em Informática  
Universidade Federal do Rio de Janeiro  
Rio de Janeiro, RJ  
jonice@gmail.com

## ABSTRACT

This paper presents the structure and functionalities of the “ENEM na Rede” tool, a free online platform that helps students from Brazil in preparing for the National High School Exam (Exame Nacional do Ensino Médio – ENEM). The application uses social networking and gamification concepts to build an engaging environment conducive to knowledge building. During the use of the platform students are introduced to the exam content through an automated method that detects deficiencies and proficiency in the exam areas.

## KEYWORDS

Information retrieval, Data mining, Social Network, Web Applications, Gamification.

## 1 CENÁRIO E CONTEXTO

No Brasil, a escola pública muitas vezes sofre pela falta de investimentos, infraestrutura e corpo docente, podendo apresentar uma disparidade de até 120% em indicadores de qualidade de ensino quando comparadas com escolas privadas [1]. O grande problema é que os alunos formados na escola pública e na privada passam pelo mesmo processo de ingresso no ensino superior.

Existem inúmeras iniciativas livres na internet que apoiam o ensino formal com técnicas de autoaprendizagem, sejam através de plataformas de simulados, videoaulas ou mesmo discussões em um grupo aberto de uma rede social on-line. A grande dificuldade dessas plataformas é que não existe um método avaliativo que possa aferir o crescimento, erro, acerto e carência de parte do conteúdo programático. Muitas vezes quando o aluno inicia um processo de estudo informal ele não sabe por onde começar e salta unidades formativas que tem pré-requisito. Esse processo ocorre em mecanismos de aprendizagem sem nenhuma mediação e avaliação.

Essa carência da rede pública de ensino deve ser suprida para garantir uma igualdade de acesso ao ensino superior. Estudantes

de baixa renda muitas vezes não podem pagar por um pré-vestibular e tem na internet uma fonte livre para complementar seus estudos.

Nesse contexto, este estudo objetiva a construção de uma ferramenta que incorporasse um método de autoaprendizagem, focado na preparação de estudantes para o Exame Nacional do Ensino Médio – ENEM. O estudo para construção do artefato foi realizado como requisito parcial para validação do método proposto na tese de doutorado do autor.

A ferramenta proposta utiliza dados públicos da rede social online Facebook, para classificar os conteúdos a serem estudados com base nas proficiências e deficiências do aluno. O sistema possui uma camada de gamificação para envolver e engajar [2] o aluno, que precisa realizar interações sociais para que a ferramenta colete os dados necessários para inferir os conteúdos a serem estudados.

A importância e relevância desse estudo está em fornecer uma ferramenta que possa ser apropriada por pesquisadores e autores para a construção de plataformas livres, que se apoiem no poder da internet para ajudar alunos no processo de construção do conhecimento.

Este método foi validado com quarenta estudantes da rede pública de ensino de uma escola estadual da cidade de Belo Horizonte – MG e apresenta dados promissores para o uso em cenários genéricos em trabalhos derivados.

## 2 DESENVOLVIMENTO DA FERRAMENTA

Para garantir significativos ganhos de aprendizagem dos conteúdos das quatro grandes áreas do ENEM, sendo elas: linguagens e códigos, matemática, ciências humanas e da natureza, foi elencada uma estrutura de funcionamento da ferramenta que está definida em: gerar testes com questões no padrão ENEM; avaliar os erros/acertos dos testes e definir proficiências e deficiências de conteúdos; montar um plano de estudos individualizado com as necessidades de conteúdos do aluno; oferecer conteúdos de repositórios externos para que o aluno utilize como apoio no processo de estudo; reiniciar esse processo de forma contínua. A Figura 1 representa essa estrutura:



Figura 1: Ciclo de funcionamento da ferramenta proposta.

A primeira fase do desenvolvimento foi à **mineração de questões de domínio público do ENEM** de anos anteriores e de vestibulares de todo o Brasil para compor a base de perguntas que seriam utilizados para validar os conhecimentos do aluno. Essa foi uma das etapas mais importantes porque todo conteúdo adicionado ao banco de dados necessitava de uma adequação para encaixar com a matriz de habilidades do ENEM. Ao todo foram mapeadas 32 mil questões divididas nas quatro grandes áreas do ENEM. Com a interface gráfica (Figura 2), tangível a operação do usuário deu-se início a construção dos **robôs autônomos** que funcionam como serviços da ferramenta. Tais processos autônomos são responsáveis por coletar os conteúdos utilizados na plataforma e realizar rotinas autônomas em intervalos de tempo definidos, como por exemplo, calcular as proficiências e deficiências de alunos e extrair dados relevantes da rede social do aluno. O funcionamento dos robôs será apresentado a seguir.

## 2.1 Funcionalidades da Ferramenta

O funcionamento da ferramenta está dividido em quatro requisitos, são eles:

- **Batalha:** principal área da plataforma. Aqui são apresentadas as questões no formato do ENEM e a partir dos erros e acertos dos alunos o sistema calcula as proficiências e deficiências para construir um plano de estudo do aluno. Essa área da plataforma é executada em duplas, e um estudante convida outro para realizar uma rodada de perguntas. Cada rodada é composta por 10 perguntas que fazem parte do escopo de deficiências do desafiador e do oponente;
- **Sugestão de Estudos:** é o campo onde são apresentados os planos de estudos do aluno. O sistema seleciona alguns assuntos relacionados aos que o aluno apresentou deficiência durante as batalhas e por isso devem ser revisados. Após a revisão, o conteúdo é inserido novamente como batalha para que seja avaliado se houve o aprendizado;
- **Grupo de Estudos:** coleta as interações sociais dos estudantes em um grupo público de Facebook e analisa seu conteúdo. Caso os assuntos debatidos sejam pertinentes ao ENEM, a plataforma atribui pontos ao estudante no sistema de gamificação;
- **Ranking:** é onde o estudante pode ver seus pontos e sua colocação no sistema. Todas as atividades da ferramenta são pontuadas. Essa mecânica é responsável pelo maior envolvimento do aluno uma vez que suas conquistas são compartilhadas com outros usuários da rede.

## 2.2 Arquitetura e Tecnologias da Plataforma

A Figura 3 demonstra a arquitetura conceitual da plataforma e mostra como os serviços e abordagens web estão interconectados para a construção de um artefato único. É importante observar que toda a interação do aluno acontece em um *front-end web* que usa os dados de usuário fornecidos pela plataforma Facebook.

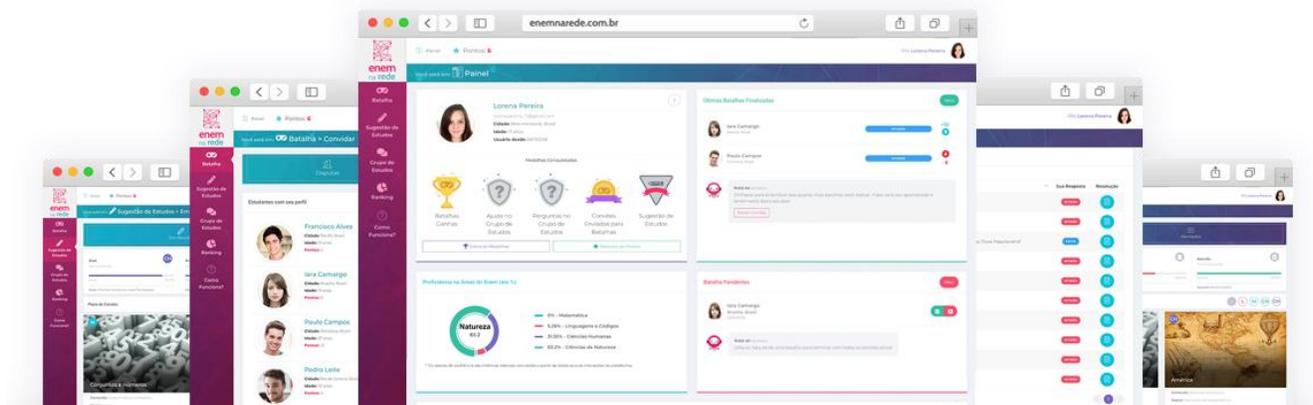
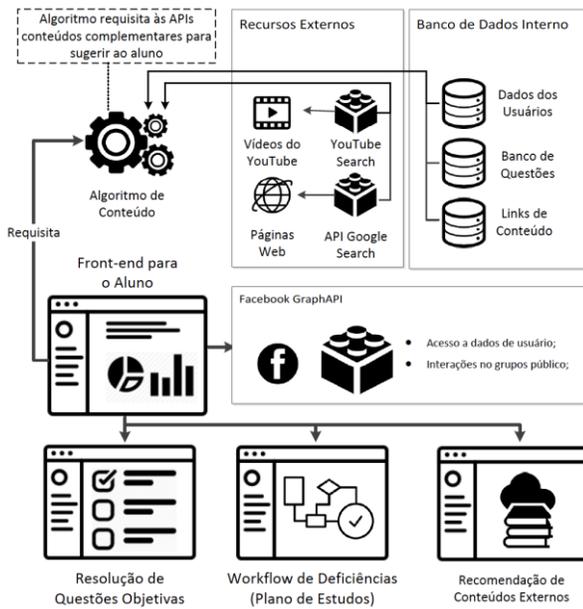
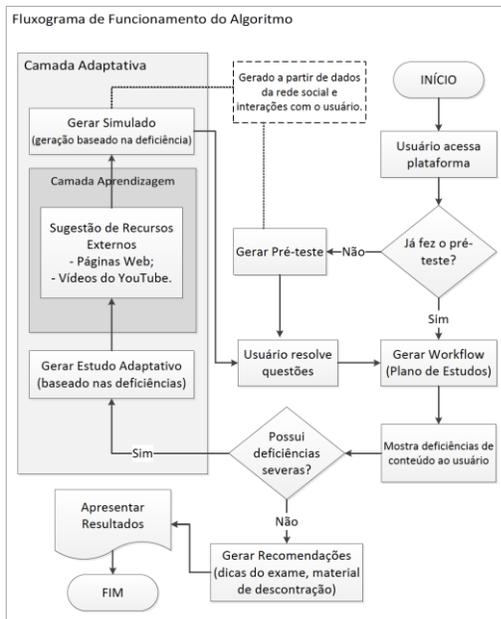


Figura 2: Interface gráfica da plataforma ENEM na Rede.



**Figura 3: Proposta de Arquitetura Conceitual da Plataforma.**

O funcionamento do “**Algoritmo de Conteúdo**”, proposto na Figura 3, está estruturado em um conjunto de passos demonstrados na Figura 4.



**Figura 4: Funcionamento do Algoritmo de Conteúdo.**

Para o *front-end* do usuário foi criado um website responsivo que utilizou as tecnologias:

- Framework de programação: PHP/Laravel;

- Motor de recuperação da informação: Elastic Search;
- Robôs autônomos: Python;
- Servidor: Apache 2 + Linux Ubuntu;
- APIs: Facebook Graph API 3.0, Google Custom Search API e Youtube Search API.

### 2.3 Processador de Conteúdos

É responsável por percorrer todo o banco de questões e para cada questão minerar através da Google Search API e do Youtube Search API conteúdos confiáveis para os estudantes. O robô escolhe os conteúdos que irão compor a base em função da reputação atribuída pelas *search engines*. No caso da Youtube Search também são consideradas as curtidas e inscrições que o vídeo e o canal possuem. É importante ressaltar que somente conteúdos públicos da internet são indexados pelos robôs.

A abordagem para coleta de websites aproveitava a inteligência do algoritmo de *pagerank* do Google Search. Sendo assim, o robô utilizando a API pública Google Search API, extraia os resultados e classificava quais iriam para a base. O termo de consulta era obtido através da extração das palavras chaves da questão, que era fornecido por um atributo já pré-processado para toda a base mais a combinação do assunto, tópico e subtópico da questão. A extração da palavra chave utilizou o seguinte procedimento: 1º - remoção de *stop words*; 2º - *stemming* e *lowercase*; 3º - contagem da frequência do termo; 4º - seleção dos três primeiros termos. Para consolidar a *query* de busca eram utilizados operadores booleanos e aspas para enclausurar expressões. Seguindo a estrutura: “*assunto*” AND “*tópico*” AND “*subtopico*” AND (“*palavra-chave 1*” OR “*palavra-chave N*”).

Durante o processo de extração de dados via API o robô executa também o processo de classificação dos itens. Para que os endereços de website, entregues pela Google Search API sejam selecionados eles devem atender os seguintes critérios:

- O conteúdo deve ter mais de um ano que figura nos resultados de busca;
- Devem estar em português;
- Obrigatório conter o assunto e tópico do corpo do texto;

Para os vídeos selecionados da plataforma Youtube:

- O saldo de curtidas (*likes*) deve ser maior que o de não-curtidas (*unlikes*);
- Deve existir comentários no vídeo;
- Deve ter um canal com outros vídeos associados.

É importante ressaltar que em ambas as consultas, são respeitados o ranking que o sistema de busca classificou. Sempre que o número mínimo de resultados (2 websites e 2 vídeos) não é atingido o robô refaz a busca removendo uma palavra-chave (das 3 utilizadas). Se depois de remover todas as palavras chaves as quantidades mínimas não forem cumpridas, o algoritmo cadastra a quantidade encontrada e envia um e-mail para o administrador do sistema com a *query* de busca e o código da questão para que o processo seja realizado manualmente.

## 2.4 Processador de Proficiências e Deficiências

Esse robô é responsável por processar os conteúdos que o aluno tem deficiências e sugerir um plano de estudos a ele. Esse robô entra em execução toda vez que um aluno realiza um teste dentro da plataforma. Baseado nos erros/acertos do teste o robô utiliza o algoritmo de proficiência/deficiência para aferir quais conteúdos demandam atenção.

Quando um aluno acerta uma questão, é necessário invocar os graus de proficiência históricos (Valores de Proficiência -  $V_p$ ), gravados em uma base. Quando o aluno é novo e não tem dados históricos, seu valor é 1. Esse valor será dividido pelos graus históricos incorretos (Valores de Deficiência -  $V_d$ ), e aplicado um logaritmo. A parte fundamental do processo está em somar ao final o grau de dificuldade do item ( $D_i$ ) [3] estabelecido por:

$$D_i = \frac{C_i - \frac{E_i}{K_i - 1}}{N_i} [3] \quad (1)$$

Onde a dificuldade de um item ( $D_i$ ) é representada pelo número de indivíduos que responderam corretamente um item ( $C_i$ ) menos a razão entre o número de alunos que erraram o item ( $E_i$ ) pelo número de respostas do item ( $K_i$ ) menos um, dividido pelo número de indivíduos que realizaram o teste ( $N_i$ ).

Sendo assim para a equação de proficiência temos:

$$\text{Proficiência} = \log\left(\frac{\sum_{i=0}^n V_p + 1}{\sum_{i=0}^n V_d + 1}\right) + D_i \quad (2)$$

Na equação de deficiência o grau de dificuldade de questão tem um peso maior do que na proficiência. Essa adequação faz com que o estudante tenha que acertar mais que uma questão para que seu grau seja convertido de deficiência em proficiência.

$$\text{Deficiência} = \log\left(\frac{\sum_{i=0}^n V_d + 1}{\sum_{i=0}^n V_p + 1}\right) + 1 + D_i \quad (3)$$

Quando uma questão não possui a dificuldade ( $D_i$ ) calculada, por ainda não ter uma base de respondentes de 100 usuários o coeficiente de dificuldade é alterado para o valor 1 de deficiência e 0.5 quando se tratar de proficiência. O coeficiente de dificuldade ( $D_i$ ) é recalculado sempre que um teste é processado, mantendo o equilíbrio da base, visto os novos usuários e testes que são realizados durante a execução do artefato.

Os valores de proficiência e deficiência são utilizados para definir quais conteúdos devem ser apresentados no *workflow*. Quanto maior a deficiência, maior a prioridade no *ranking* de conteúdos. Quando o saldo de proficiência é maior que o de deficiência o conteúdo é transferido para uma área de itens revisados.

## 2.5 Processador de Interações Sociais

Responsável por quantificar a participação do aluno no grupo de estudos da plataforma. O grupo de estudos opera dentro da

rede social online Facebook e diariamente, a cada 2 horas, o robô visita a rede, extrai as participações do aluno e verifica se elas têm correlação com algum conteúdo da base. Caso positivo o aluno é bonificado com uma quantidade de pontos no sistema de gamificação da plataforma.

Para processar as interações dos alunos no grupo, diariamente a cada duas horas um robô conecta-se a Graph API do Facebook e lê todas as postagens realizadas no grupo. As postagens são divididas em dois grupos: perguntas e respostas. As perguntas são postagens realizadas no fluxo principal do grupo, conhecido como *feed* e as respostas são comentários adicionados a perguntas do *feed*. Para cada post realizado, o algoritmo verifica qual o autor do conteúdo e processa o texto para identificar se está relacionado a algum conteúdo da matriz de competência do ENEM, caso positivo, ele salva o conteúdo na área de “Grupo de Estudos” do aluno e atribui pontos pela publicação.

O processamento do conteúdo é realizado por técnicas de recuperação da informação, que extrai do texto do aluno as palavras-chaves e compara com a base de dados de questões do artefato. A primeira fase de construção do algoritmo foi a indexação de todas as palavras-chaves das questões, agrupadas por subtópico e tópico. Sempre que um novo conteúdo era recuperado, suas palavras-chaves eram comparadas as palavras-chaves da base indexada, esse processo retorna um coeficiente de correlação, quanto mais alto, maior a probabilidade de os conteúdos serem relacionados. Para a postagem do aluno, os tratamentos realizados nas *strings* foram: 1° - remoção de *stop words*; 2° - *stemming e lowercase*; 3° - contagem da frequência do termo. Já para o processo de indexação, que é usado como a base de comparação para a primeira fase foram considerados os seguintes passos:

- Junção de todas as palavras-chaves das questões por tópico;
- Remoção dos *strings* duplicados para o atributo palavras-chaves, *stop words, stemming e lowercase*;
- Junção de todos os subtópicos das questões por tópico;
- Remoção dos *strings* duplicados para o atributo subtópico, *stop words, stemming e lowercase*;
- Agrupamento por tópico e assunto;
- Indexação dos termos utilizando o método *Term frequency—inverse document frequency selection* (TF-IDF);

Após o processo de indexação, o robô comparava se o texto digitado pelo aluno tinha correlação com os textos indexados na base, retornando um coeficiente.

Nesse momento percebeu-se que somente aplicar o método TF-IDF [4] não geraria uma correlação funcional para essa pesquisa. Foi então que o algoritmo foi adaptado para funcionar com pesos individuais em cada atributo (tópico, subtópico e palavra-chave). Além de acrescentar o peso individual foi adicionado um fator de normalização do TF que utiliza o tamanho do documento como um critério de importância. Essa normalização garantiu melhorias significativas da correlação, visto que alguns documentos indexados tinham palavras-chaves composta por um baixo número de questões da base, e outros documentos tinham muitas questões relacionadas, o que resultava um grande número de palavras-chaves e

consequentemente o aumento do tamanho do documento. O processo de normalização permitiu que os documentos com tamanhos diferentes fossem comparados sem que o documento maior expressasse uma significativa discrepância contra um documento pequeno.

Estruturando novamente a equação temos então:

$$TF - IDF \text{ Normalizado} = TF * IDF$$

$$TF = \sqrt{TF_{x,y}} * \frac{1}{\sqrt{l}} * w$$

$$IDF = \left( \log \frac{N+1}{df+1} + 1 \right) * w$$
(4)

Onde  $TF_{x,y}$  é a frequência de uma palavra-chave  $X$  em um documento  $Y$ ,  $N$  é o total de documentos na base, e  $df$  é o número de documentos que contém a palavra-chave  $X$ , o  $w$  é o peso que será variável por atributo, trazendo a possibilidade de definir uma maior importância, por exemplo ao subtópico. O  $l$  é o tamanho do documento em palavras.

O modelo proposto traz a pesquisa as seguintes propriedades:

- Pesos individuais para os atributos;
- Normalização baseado no tamanho do documento;
- Valores sempre positivos;

Com o modelo proposto os documentos foram indexados e seus atributos receberam os seguintes fatores de importância:  $topico_w = 6$ ;  $subtopico_w = 4$ ;  $palavras - chaves_w = 0.6$ . Esse coeficiente foi ajustado realizando inúmeras buscas para processar a eficiência. É possível observar que se uma *string* está presente no tópico do documento ele ganha um alto fator de importância, o que faz sua nota de correlação crescer.

Durante a etapa de testes uma outra característica foi adicionada ao algoritmo para melhorar sua eficiência nas consultas. Essa característica foi chamada de **Correspondência**, e definia que para uma *string* ter correlação com a base de documento, 30% das palavras-chaves deveriam estar presentes em todos os atributos – tópico, subtópico e palavras-chaves. Essa característica removeu muitos falsos positivos nas consultas, porque garantia que o que estava sendo procurado existia no tópico, subtópico e palavras-chave.

O funcionamento desse algoritmo foi essencial para que o artefato trabalhasse uma nova dimensão do aprendizado e utilizasse as interações sociais como fator de engajamento de alunos.

### 3 PERSPECTIVAS E LICENCIAMENTO

A ferramenta proposta abre novas perspectivas para o desenvolvimento de aplicações para construção do conhecimento. Os resultados demonstram o potencial promissor da aplicação e a capacidade de uso em outros contextos e ambientes. Para continuar aprimorando a ferramenta e trazendo

a possibilidade da geração de produtos ou serviços derivados ela continuará com acesso gratuito para a comunidade acadêmica e sociedade. A tese de doutorado e artigos que detalham o funcionamento da aplicação podem ser utilizados como ponto de partida para outros pesquisadores que desejam criar aplicativos correlatos.

O uso da plataforma é livre e não é regida por nenhum modelo de comercialização ou uso de software (como EULA, GNU GPL ou MIT License), seus utilizadores não podem comercializar, divulgar ou se apropriar das informações e dados da ferramenta.

### 4 CONCLUSÃO

Essa pesquisa objetivou o desenvolvimento de uma ferramenta de autoaprendizagem que pudesse ser combinada ao estudo formal de alunos do ensino médio da rede pública de ensino no Brasil para ajudar na preparação para o ENEM. O método partia da premissa de identificar as proficiências e deficiências de conteúdo dos estudantes e criar um ambiente de colaboração social com atividades para corrigir os déficits de aprendizagem.

A fase de avaliação da plataforma, que gerou o conjunto de resultados que demonstram ganhos de aprendizagem de quando o grupo de alunos iniciou o processo até a final da pesquisa. Foi possível estabelecer que quanto maior o uso das atividades propostas na ferramenta, melhor o rendimento nos testes que aferem a proficiência. Os experimentos demonstram que o método pode promover a aprendizagem de forma eficiente, o que garante a efetividade da ferramenta.

### 5 DEMONSTRAÇÃO DA FERRAMENTA

Foi desenvolvido um vídeo tutorial da ferramenta proposta que apresenta todas as áreas de sistema, esse vídeo pode ser acessado pelo link:

<https://www.youtube.com/watch?v=4pyhTwEFjMY>

### REFERÊNCIAS

- [1] A. Moraes and W. Belluzzo, "O diferencial de desempenho escolar entre escolas públicas e privadas no Brasil", Nova Economia. vol.24 no.2 Belo Horizonte Mai/Ago. 2014.
- [2] G. Zichermann and C. Cunningham, "Gamification by Design: Implementing Game Mechanics in Web and Mobile Apps", Canada: O'Reilly Media, 2011. [3] A. Dresch, D. P. Lacerda and J. A. Junior, "Design Science Research". Bookman, 2015.
- [3] L. Pasquali, "Psicometria: teoria dos testes na psicologia e na educação", Petrópolis, Vozes, 2003. 397 p.
- [4] M. Leginus and C. Zhai, DOLOG, Peter. "Personalized generation of word clouds from tweets," J. Assoc. Inf. Sci. Technol., p. n/a–n/a, 2015.