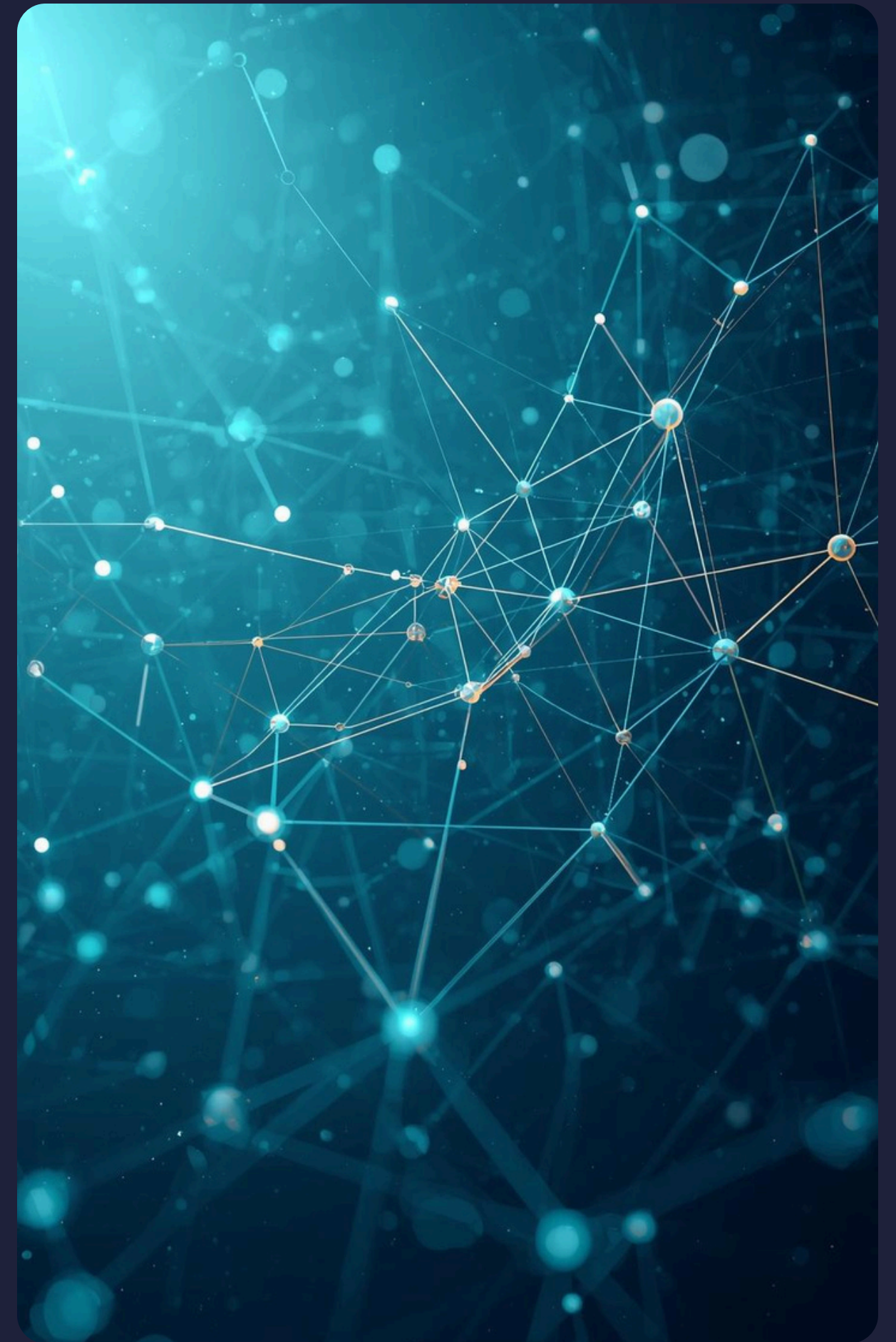


Privacidade Diferencial

João Gabriel Borges Monteiro
Felipe dos Reis



Pseudo-código do kNN Diferencial

Entrada: Conjunto de Treino (X_{train}, Y_{train}), Instância de Teste X_{query} , Hiperparâmetro k , Orçamento de privacidade ϵ .

Saída: Classe prevista \hat{y} .

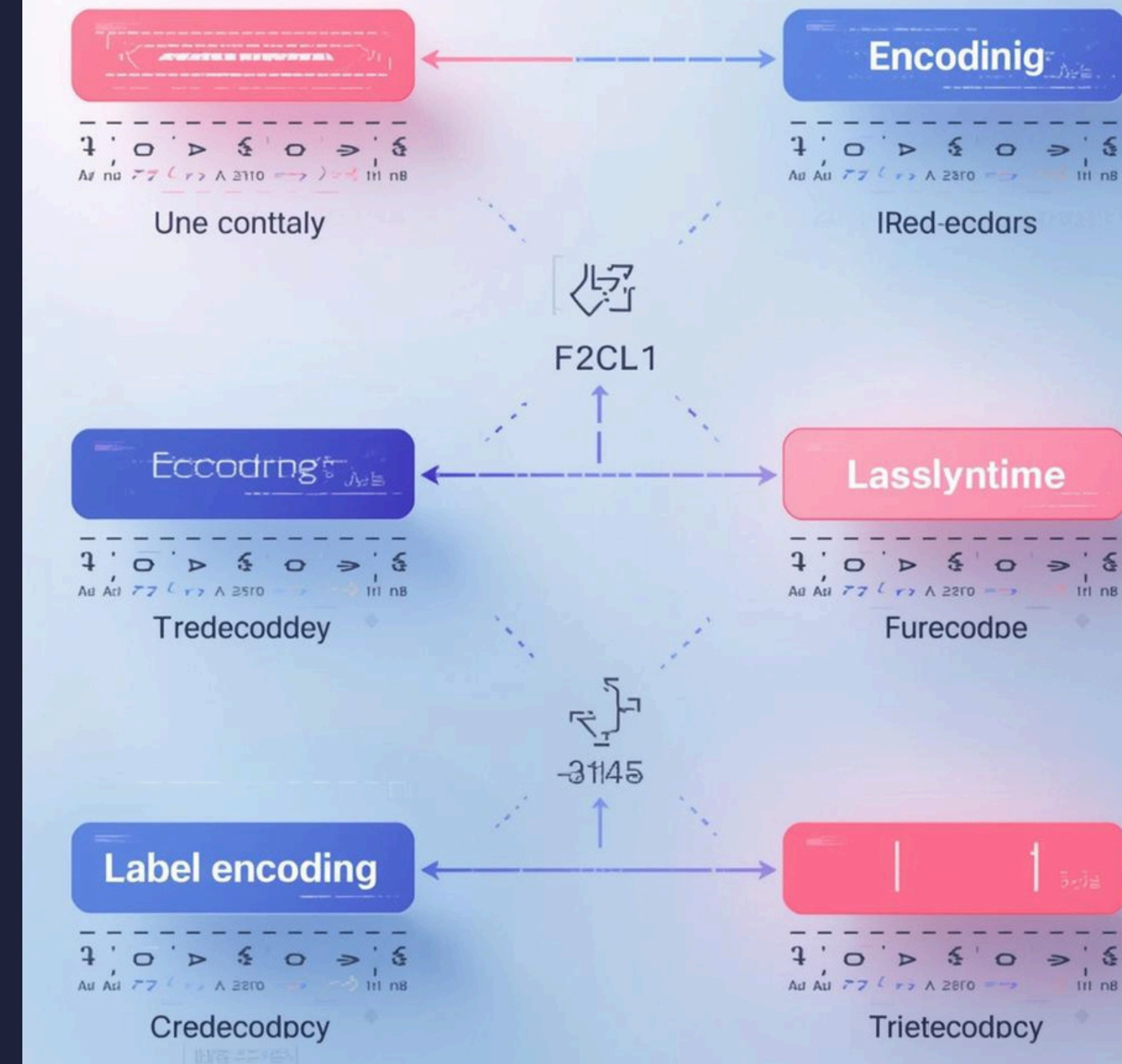
1. Carregar dados: Armazenar X_{train} e Y_{train} na memória.
2. Para cada instância de teste x_{query} :
3. Calcular Distâncias:
 - Calcular a distância Euclidiana entre x_{query} e todos os pontos em X_{train} .
 - $d(x_{query}, x_i) = \sum (x_{query} - x_i)^2$
4. Identificar Vizinhos:
 - Ordenar as distâncias em ordem crescente.
 - Selecionar os índices das k menores distâncias.
 - Recuperar os rótulos (classes) Y desses k vizinhos.
5. Contabilizar Votos (Contagem Real):
 - Para cada classe única c , calcular $count(c)$: quantos vizinhos pertencem à classe c .
6. Se kNN Tradicional:
 - Retornar a classe com o maior valor de $count(c)$.
7. Se kNN com Privacidade (Mecanismo de Laplace):
 - Definir sensibilidade $\Delta f = 1$.
 - Calcular o parâmetro de escala $b = \Delta f / \epsilon$
 - Para cada classe c :
 - Gerar ruído $ruido \sim \text{Laplace}(0, b)$
 - Calcular $voto_ruidoso(c) = count(c) + ruido$
 - Retornar a $\max(voto_ruidoso(c))$.

Codificação de Atributos Não Numéricos

No código apresentado, a conversão é feita explicitamente na linha:

```
'''  
for col in df.columns:  
    if df[col].dtype == 'object' or pd.api.types.is_categorical_dtype(df[col]):  
        df[col], _ = pd.factorize(df[col])  
'''
```

Isso fez com que os dados de categorização ficassem no intervalo entre 1, ...n, onde n é o número de classes categóricas única da coluna codificada.

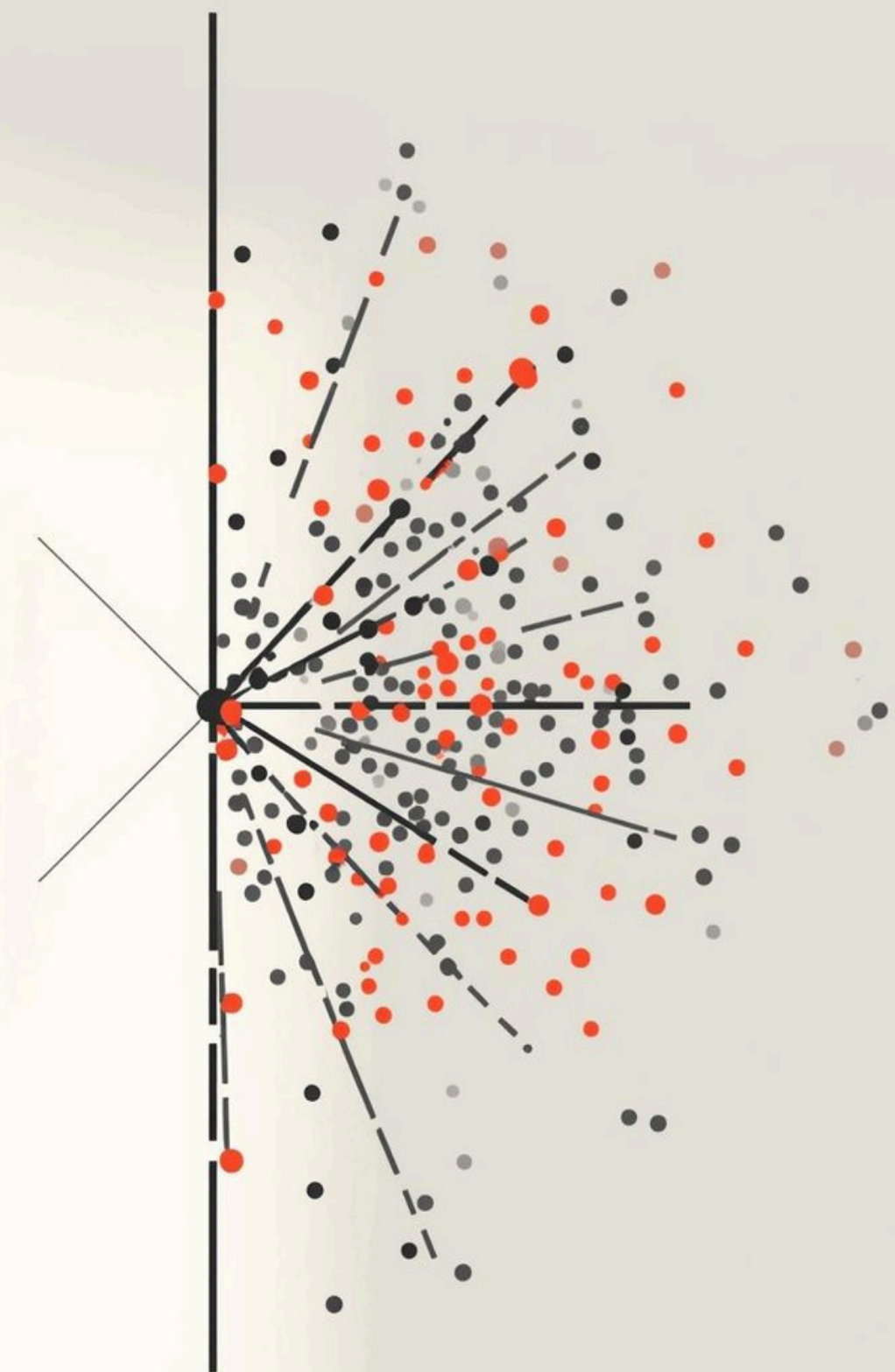


Aplicação do Mecanismo de Laplace

"O mecanismo de Laplace foi aplicado sobre a contagem de votos das classes vizinhas (histograma de votos). Em vez de simplesmente escolher a classe majoritária entre os k vizinhos, adicionou-se um ruído aleatório extraído de uma distribuição de Laplace, centrado em 0 com escala $b = 1/\epsilon$, a cada contagem de classe. A classe predita foi aquela que obteve a maior contagem após a adição do ruído. Isso possivelmente mascara a influência exata de qualquer vizinho individual na decisão final do algoritmo."

Pega a contagem \rightarrow adiciona ruído \rightarrow escolhe quem ganhou \rightarrow retorna a classe

EPLACE Mechaniser



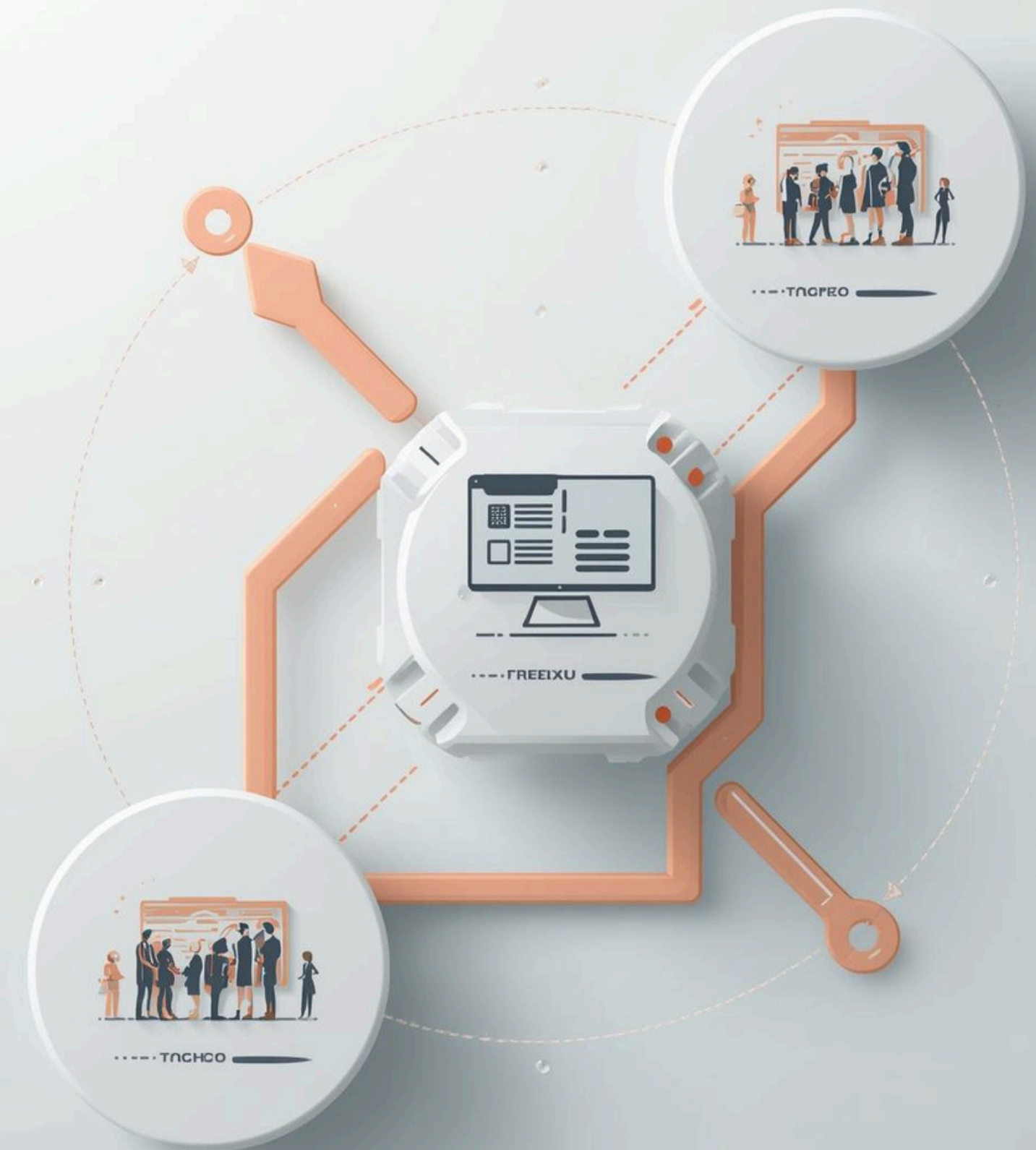
MAETHONISAN ◀ TERVATION = PIRCHCY

Orçamento do Mecanismo de Laplace

O orçamento usado nos mecanismo foi do intervalo de $\epsilon = [0.5, 1, 1.5, 2.0, 2.5, 3.0]$

Motivações :

Sugestão da especificação + Construção de Curva de Convergência de Acurácia



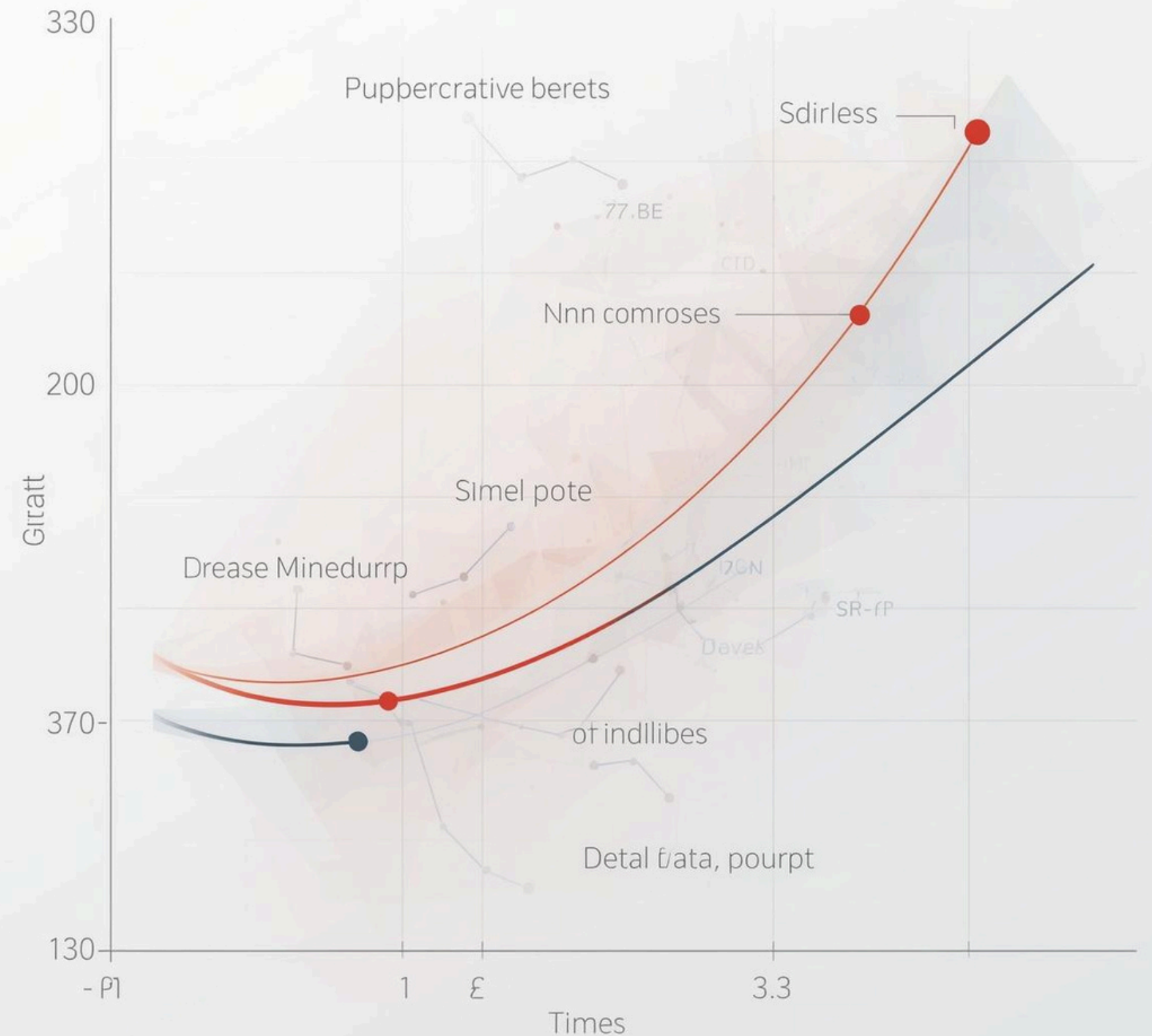
Sensibilidade

O valor da sensibilidade é $\Delta f=1$

Isso ocorre porque a função de consulta é uma contagem de votos. Ao modificar, adicionar ou remover um único registro da base de dados de treino, o conjunto de vizinhos em R pode ser alterado.

No pior caso, um vizinho da classe A é substituído por um vizinho da classe B. Isso altera a contagem da classe A em -1 ou da classe B em +1. Portanto, a magnitude máxima da mudança na contagem de qualquer classe dada a alteração de um único indivíduo é 1.

Global Sensitivity



Acurácia x Orçamento(ϵ)

