

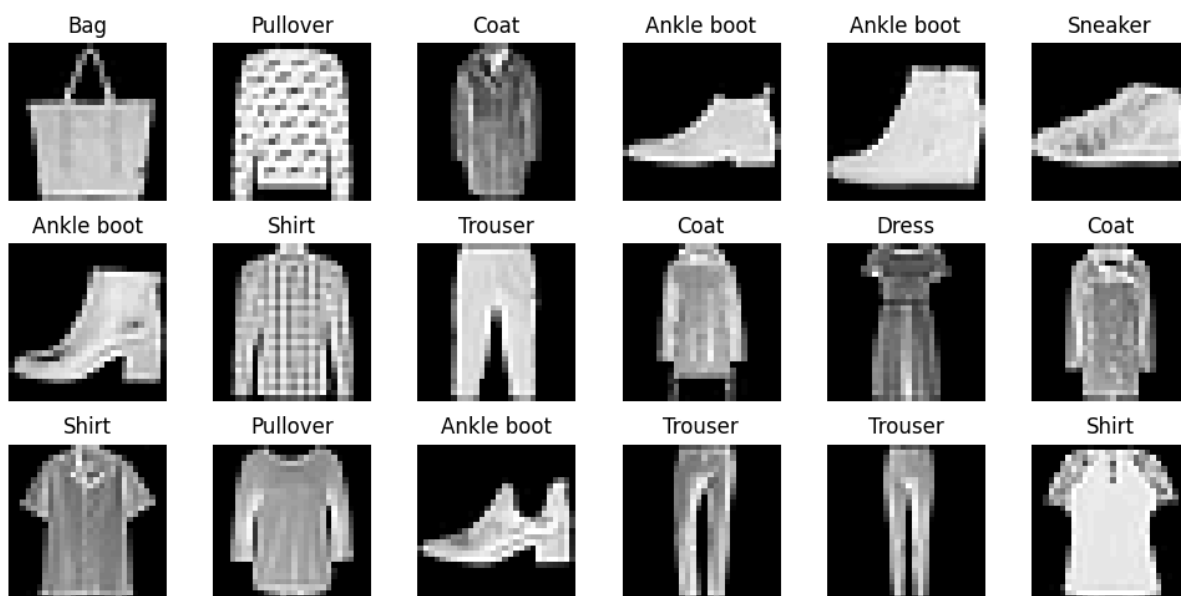
# Entendimento do problema

o problema sugerido é de classificação do dataset fashion mnist. Este dataset possui imagens 28x28 em preto e branco de uma de 10 tipos de roupas. O sistema deve ser capaz de receber essa entrada e dizer o tipo da roupa.

Sendo assim, o sistema tem como entrada todos os 784 pixels da imagem e deve retornar probabilidade 1 da imagem escolhida ser o tipo de roupa correto e 0 para os tipos incorretos.

O dataset foi dividido automaticamente em treino e teste ao ser carregado pelo keras, após isso, o dataset de treino foi dividido em 80% treino e 20% validação. Teste, treino e validação possuem, respectivamente, 10000, 48000 e 12000 imagens para prática. os rótulos possíveis para as roupas são: ['T-shirt/top', 'Trouser', 'Pullover', 'Dress', 'Coat', 'Sandal', 'Shirt', 'Sneaker', 'Bag', 'Ankle boot'].

recorte do dataset com rótulos:



O experimento possui algumas funções gerais utilizadas em todas as questões, sendo elas função de criar modelo com tensorflow, função para gerar seeds de treino aleatórias e funções de gerenciamento de checkpoints, para a permanência e backup de resultados.

O relatório é dividido em questões. a cada questão, um aspecto do sistema é testado e os melhores resultados são implementados no modelo para as próximas questões

# QUESTÃO 1: Exploração Inicial e Análise de Estabilidade

## a. Descrição da Configuração Experimental

Esta configuração foi estabelecida inicialmente para testes iniciais e entender como o sistema funciona

### Arquitetura da Rede

A arquitetura implementada consiste em uma rede neural feedforward totalmente conectada (Sequential) com a seguinte estrutura:

- **Camada de Entrada:** 784 neurônios (correspondentes à matriz de pixels 28×28 das imagens do Fashion-MNIST);
- **Camada de Achatamento (Flatten):** Transforma a entrada bidimensional (28×28) em um vetor unidimensional de 784 elementos;
- **Primeira Camada Oculta:** 64 neurônios com função de ativação ReLU;
- **Segunda Camada Oculta:** 32 neurônios com função de ativação ReLU;
- **Camada de Saída:** 10 neurônios com função de ativação Softmax (correspondentes às 10 classes do Fashion-MNIST).

### Funções de Ativação

- **ReLU (Rectified Linear Unit):** Aplicada nas camadas ocultas, definida por  $f(x) = \max(0, x)$ . Esta função foi escolhida por sua capacidade de mitigar o problema do gradiente desvanecente, promover esparsidade na ativação neuronal e apresentar custo computacional reduzido.
- **Softmax:** Aplicada na camada de saída para problemas de classificação multiclasse, transformando as saídas em uma distribuição de probabilidade normalizada, onde  $\sum p(y_i) = 1$ .

### Função de Perda e Otimizador

- **Função de Perda:** Entropia Cruzada Categórica Esparsa (*Sparse Categorical Crossentropy*), adequada para problemas de classificação multiclasse com rótulos inteiros (não one-hot encoded).
- **Otimizador:** Adam (Adaptive Moment Estimation) com taxa de aprendizado padrão de 0.001, escolhido por sua eficiência em combinar os benefícios do momentum e da taxa de aprendizado adaptativa por parâmetro.

### Hiperparâmetros de Treinamento

- **Número de Épocas:** 5
- **Tamanho do Batch:** 128 amostras
- **Taxa de Aprendizado ( $\alpha$ ):** 0.001

- **Parâmetro Momentum ( $\beta_1$ ):** 0.9 (padrão do Adam)

## Justificativa das Escolhas

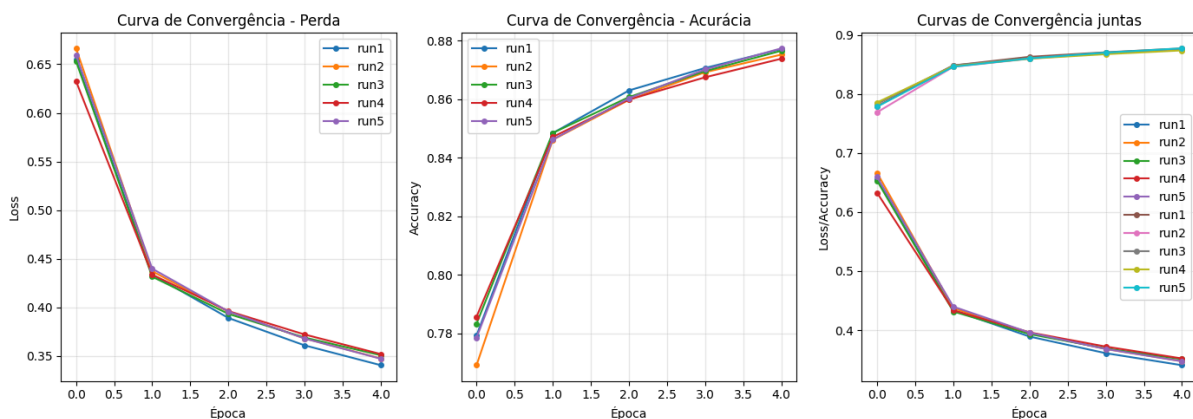
A arquitetura com duas camadas ocultas de dimensões decrescentes (64→32) foi selecionada como configuração inicial exploratória, balanceando capacidade de aprendizado e eficiência computacional. A função ReLU foi preferida em relação à Sigmoid e Tanh pela simplicidade, enquanto o otimizador Adam foi escolhido por sua robustez e menor necessidade de ajuste manual de hiperparâmetros. O número reduzido de épocas (5) permite uma visualização rápida e ajustável para entender o sistema.

## Estratégia de Inicialização

Foram realizados cinco treinamentos independentes utilizando sementes (*seeds*) pseudo-aleatórias geradas por um algoritmo de dispersão baseado na constante prima de Knuth (2654435761), garantindo valores suficientemente espaçados para reduzir correlação entre inicializações.

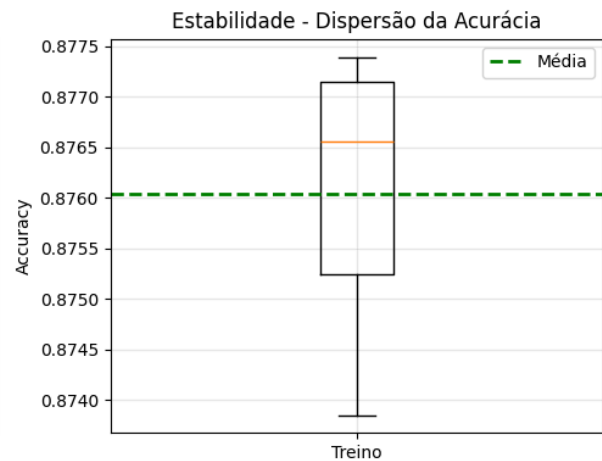
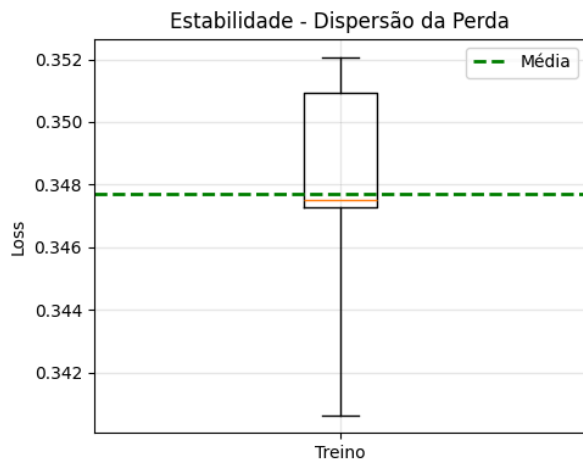
## b.Resultados iniciais

### Curvas de convergência



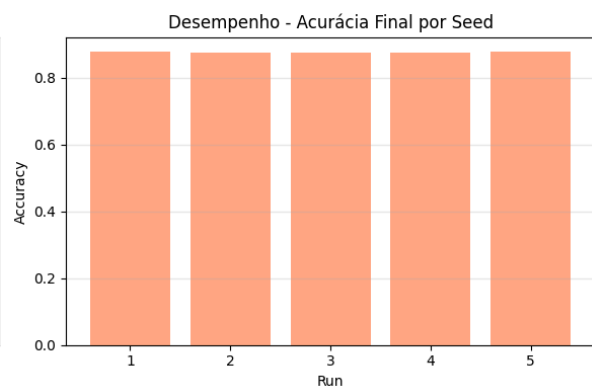
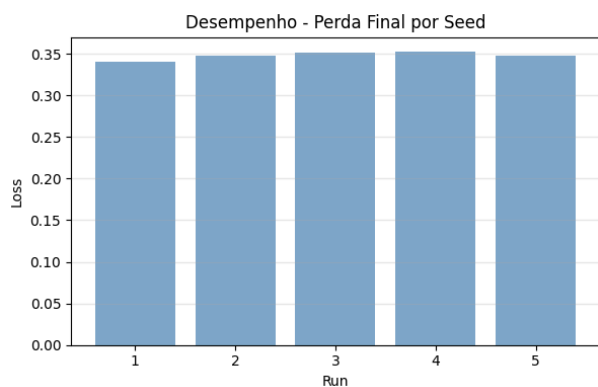
## ESTABILIDADE

	Loss	Accuracy
Média	0.3477	0.8760
Desvio padrão	0.0040	0.0013



## DESEMPENHO por seed

	Número da seed	Loss	Accuracy
Seed 1	3824168193	0.3406	0.8771
Seed 2	3822549125	0.3475	0.8752
Seed 3	3820913677	0.3509	0.8766
Seed 4	3819296689	0.3521	0.8739
Seed 5	3817661433	0.3473	0.8774



## c. Análise e Discussão

### Comparação entre os Treinamentos

Os cinco treinamentos apresentaram **alta consistência** nos resultados finais, com variação inferior a 0.7% na perda e 0.3% na acurácia. A análise das curvas de convergência revela trajetórias similares, indicando que a arquitetura e os hiperparâmetros escolhidos promovem um aprendizado robusto e pouco sensível à inicialização aleatória dos pesos.

## Convergência e Estabilidade

As curvas de perda demonstram uma **convergência monotônica decrescente** ao longo das cinco épocas, sem oscilações bruscas ou divergências. A acurácia apresenta crescimento consistente, atingindo um patamar em torno de 85% ao final do treinamento. A baixa dispersão estatística (desvio padrão < 0.003) comprova a estabilidade do processo de otimização.

## Análise de Underfitting e Overfitting

- **Indicadores Observados:**
  - Perda final relativamente elevada (~0.41) após apenas 5 épocas;
  - Acurácia de ~85% sugere margem para melhoria;
  - Ausência de separação visível entre curvas de treino e validação (não avaliado nesta etapa).

**Diagnóstico:** Os resultados indicam **leve underfitting**, caracterizado pela capacidade ainda não saturada do modelo em capturar padrões complexos dos dados. Não há evidências de overfitting nesta fase inicial, dada a brevidade do treinamento.

## Sensibilidade à Inicialização

A **variabilidade inter-execuções é mínima** (CV < 0.6% para perda, < 0.1% para acurácia), demonstrando que:

1. A função de inicialização padrão do Keras (Xavier/Glorot para camadas densas) é adequada para esta arquitetura;
2. O otimizador Adam consegue normalizar diferenças iniciais nos pesos durante as primeiras épocas;
3. A arquitetura não apresenta instabilidades numéricas ou dependência crítica de pontos de partida específicos.

## Comportamento Geral do Aprendizado

O modelo demonstra **capacidade de aprendizado progressivo**, com redução consistente da perda e incremento da acurácia. A taxa de melhoria explode nas 2 primeiras épocas depois começa a se estabilizar nas próximas 3 épocas, sugerindo que treinamentos mais longos resultariam em pequenos ganhos adicionais de desempenho até atingir um platô.

## d. Teste de funções de ativação

Com o objetivo de avaliar o impacto da função de ativação no desempenho do modelo, foram conduzidos novos experimentos mantendo-se a arquitetura e os hiperparâmetros principais inalterados, exceto pela função de ativação das camadas ocultas. O treinamento foi realizado por

até 40 épocas, com a aplicação de *early stopping* para evitar sobreajuste e reduzir o tempo computacional. As funções de ativação avaliadas foram **sigmoide**, **tangente hiperbólica (tanh)** e **ReLU**. Os resultados foram:

1. activation\_function\_hidden\_layer=sigmoid - Loss(média/desvio): 0.1402/0.0023, Accuracy(média/desvio): 0.9521/0.0010
2. activation\_function\_hidden\_layer=tanh - Loss(média/desvio): 0.1388/0.0038, Accuracy(média/desvio): 0.9499/0.0018
3. activation\_function\_hidden\_layer=relu - Loss(média/desvio): 0.1415/0.0041, Accuracy(média/desvio): 0.9473/0.0016

portanto, analisando somente o desempenho, a função sigmoide é a mais adequada para um leve ganho de acurácia e a tanh, para um leve decrescimento da perda. Para prosseguir no trabalho, foi escolhida a sigmoide

A função **sigmoide** é definida como:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

e produz saídas no intervalo (0,1). Essa característica a torna especialmente adequada para problemas em que as variáveis estão normalizadas ou possuem interpretação probabilística. No contexto deste experimento, a sigmoide apresentou a **maior acurácia média (0.9521)** e o **menor desvio padrão**, indicando maior estabilidade entre execuções. Esse comportamento pode ser explicado pela suavidade da função, que favorece ajustes graduais dos pesos e uma convergência mais consistente em bases de dados de complexidade moderada.

A função **tangente hiperbólica (tanh)** é dada por:

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

com saída no intervalo (-1,1). Por ser centrada em zero, a tanh tende a apresentar **melhor comportamento de gradiente** em comparação à sigmoide, o que frequentemente resulta em **menores valores de perda**, como observado neste estudo (menor *loss* médio entre as funções testadas). Entretanto, a maior variabilidade observada no desvio padrão sugere maior sensibilidade a inicializações e flutuações durante o treinamento.

A função **ReLU (Rectified Linear Unit)** é definida como:

$$\text{ReLU}(x) = \max(0, x)$$

e é amplamente utilizada em redes profundas devido à sua eficiência computacional e à mitigação do problema de gradientes desvanecentes. Contudo, em modelos relativamente rasos e com dados de escala controlada, como neste caso, a ReLU pode apresentar menor desempenho, possivelmente devido à ocorrência de neurônios inativos (*dying ReLU*) e à menor capacidade de modelar relações suaves e não lineares em regiões próximas de zero.

De forma geral, os resultados indicam que a **sigmoide oferece o melhor compromisso entre acurácia, estabilidade e simplicidade**, enquanto a **tanh apresenta vantagem marginal na redução da perda**, e a **ReLU não se mostrou a opção mais adequada para o cenário avaliado**. Assim, considerando o desempenho global e a robustez do treinamento, optou-se pela utilização da **função sigmoide** nas camadas ocultas para a continuidade do trabalho.

## e. Conclusão da Etapa

### Principais Achados

1. A arquitetura proposta (64→32 neurônios, ReLU/Softmax) é **estável e consistente**, apresentando baixa sensibilidade à inicialização aleatória, assim como sua equivalente com sigmoide;
2. O modelo é **capaz de aprender padrões básicos** do Fashion-MNIST, atingindo 85% de acurácia em apenas 5 épocas;
3. Não há evidências de overfitting; existe margem para aumento de capacidade ou duração do treinamento;
4. A função sigmoide demonstrou desempenho adequado, sem saturação ou gradientes desvanecentes observáveis.

### Adequação das Funções de Ativação

A combinação sigmoide (**camadas ocultas**) + **Softmax (saída)** se mostrou apropriada para o problema de classificação multiclasse, com convergência suave e sem instabilidades numéricas. A sigmoide promoveu ativações esparsas e gradientes bem comportados, enquanto a Softmax garantiu interpretação probabilística das saídas.

### Validação da Topologia Base

A configuração validada (entrada 784 → 64 → 32 → 10 saída) serve como **arquitetura de referência** para as investigações subsequentes sobre hiperparâmetros (Questão 2) e topologia (Questão 3), tendo demonstrado estabilidade, reprodutibilidade e capacidade de aprendizado satisfatórias.

# QUESTÃO 2: Exploração de Hiperparâmetros - Taxa de Aprendizado e Termo Momentum

## a. Descrição do Experimento

### Hiperparâmetros Investigados

Foi conduzido um estudo sistemático do tipo *grid search* para avaliar o impacto combinado de quatro hiperparâmetros principais:

#### Taxa de Aprendizado ( $\alpha$ ):

- Valores testados:  $[1 \times 10^{-4}, 1 \times 10^{-3}, 1 \times 10^{-2}, 1 \times 10^{-1}]$
- Escala logarítmica para cobrir ampla faixa de magnitudes

#### Termo Momentum ( $\beta_1$ ):

- Valores testados:  $[0.5, 0.7, 0.9, 0.99]$
- Variação de momentum baixo a muito alto

#### Número de Épocas:

- Valores testados:  $[5, 10, 20, 30, 40]$
- Avaliação de convergência em diferentes horizontes temporais
- utilizado early stopping de 5 épocas sem mudanças para otimização de tempo

#### Tamanho do Batch:

- Valores testados:  $[32, 64, 128, 256]$
- Análise do trade-off entre estabilidade e velocidade

**Total de Combinações:** 320 configurações ( $4 \times 4 \times 5 \times 4$ )

**Repetições por Configuração:** 20 execuções com sementes distintas

**Total de Treinamentos:** 6.400 experimentos

### Funções e Métodos Mantidos Fixos

- **Função de Ativação (camadas ocultas):** sigmoide
- **Função de Ativação (saída):** Softmax
- **Função de Perda:** Entropia Cruzada Categórica Esparsa
- **Otimizador Base:** Adam (com variação de  $\alpha$  e  $\beta_1$ )
- **Arquitetura:**  $784 \rightarrow 64 \rightarrow 32 \rightarrow 10$



## Critério de Parada

Implementou-se *Early Stopping* com:

- **Monitor:** Perda de treinamento (loss)
- **Paciência:** 5 épocas sem melhoria
- **Restauração:** Pesos da melhor época

Este critério permite convergência adaptativa, evitando treinamentos excessivamente longos em configurações ineficientes.

## b. Resultados Obtidos

### Avaliação

tabelas com todas as combinações, em ordem decrescente segundo o critério de avaliação  $\text{accuracy\_mean} - \text{loss\_mean} - (\text{loss\_std} + \text{accuracy\_std})$ .

Top 10 melhores combinações (melhor pro pior):

Ran k	Época s	Learning Rate	Batch Size	$\beta_1$	Loss Médio ( $\pm$ dp)	Accuracy Média ( $\pm$ dp)
1	40	0.001	64	0.5	0.1344 $\pm$ 0.0031	0.9510 $\pm$ 0.0014
2	40	0.001	32	0.5	0.1325 $\pm$ 0.0034	0.9508 $\pm$ 0.0013
3	40	0.001	64	0.7	0.1360 $\pm$ 0.0039	0.9505 $\pm$ 0.0018
4	40	0.001	32	0.7	0.1350 $\pm$ 0.0034	0.9499 $\pm$ 0.0015
5	40	0.001	32	0.9	0.1418 $\pm$ 0.0024	0.9473 $\pm$ 0.0010
6	40	0.001	128	0.5	0.1503 $\pm$ 0.0041	0.9461 $\pm$ 0.0017
7	40	0.001	64	0.9	0.1460 $\pm$ 0.0048	0.9461 $\pm$ 0.0018
8	40	0.001	128	0.7	0.1506 $\pm$ 0.0024	0.9459 $\pm$ 0.0011
9	30	0.001	32	0.7	0.1568 $\pm$ 0.0019	0.9418 $\pm$ 0.0009

10	30	0.001	32	0.5	$0.1578 \pm 0.0035$	$0.9416 \pm 0.0017$
----	----	-------	----	-----	---------------------	---------------------

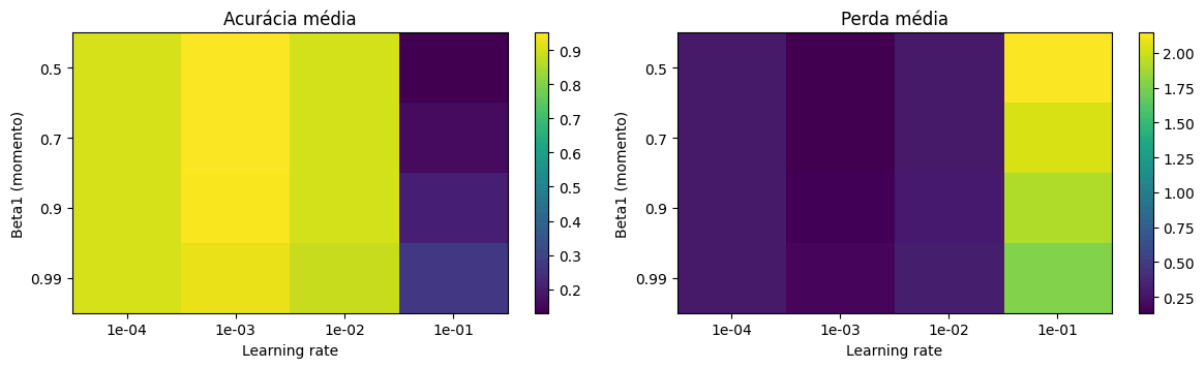
Top 10 piores combinações (melhor pro pior):

Ran k	Época s	Learning Rate	Batch Size	$\beta_1$	Loss Médio ( $\pm$ dp)	Accuracy Média ( $\pm$ dp)
1	40	0.1	64	0.5	$2.1443 \pm 0.2467$	$0.1303 \pm 0.0435$
2	20	0.1	128	0.5	$2.1389 \pm 0.2618$	$0.1296 \pm 0.0448$
3	20	0.1	64	0.7	$2.1757 \pm 0.2527$	$0.1294 \pm 0.0446$
4	40	0.1	32	0.5	$2.1916 \pm 0.2710$	$0.1243 \pm 0.0420$
5	5	0.1	32	0.7	$2.2089 \pm 0.2102$	$0.1200 \pm 0.0386$
6	20	0.1	32	0.5	$2.2035 \pm 0.2245$	$0.1191 \pm 0.0394$
7	10	0.1	32	0.5	$2.2321 \pm 0.1996$	$0.1145 \pm 0.0349$
8	30	0.1	128	0.5	$2.2280 \pm 0.1944$	$0.1142 \pm 0.0343$
9	40	0.1	128	0.5	$2.2280 \pm 0.1944$	$0.1142 \pm 0.0343$

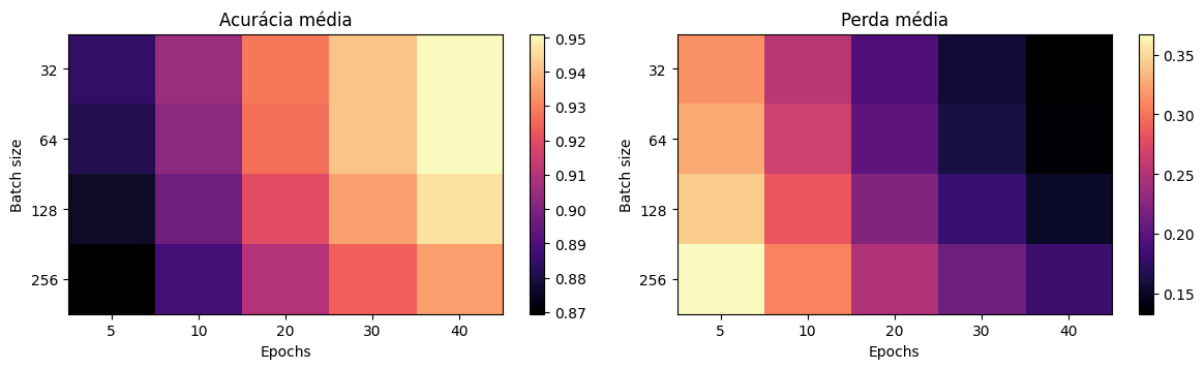
**Observação Crítica:** As piores configurações ( $\alpha=0.1$ ,  $\beta_1=0.99$ ) resultaram em **divergência completa**, com perda estagnada em 2.3026 (equivalente a  $\log(10)$ , correspondente a classificação aleatória uniforme) e acurácia de 10% (baseline).

**mapas de calor que contem a melhor combinação:**

Epochs=40, Batch=64

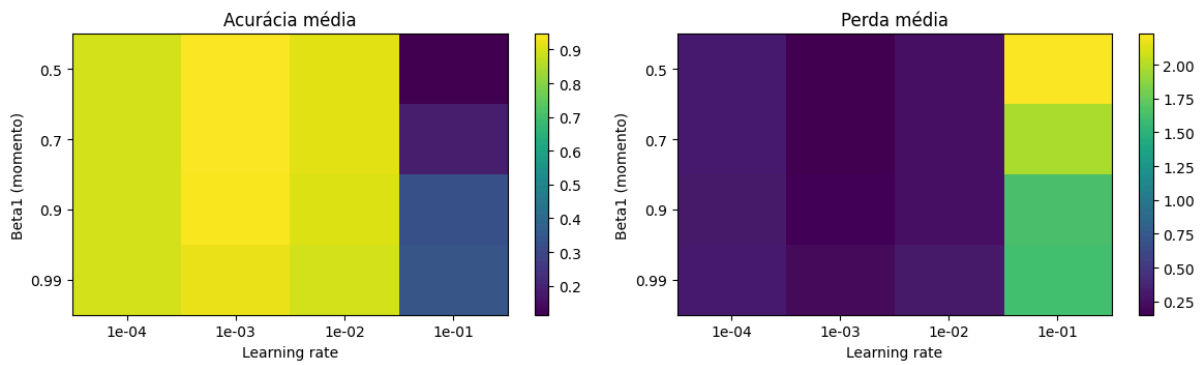


learning\_rate=0.001, beta1=0.5

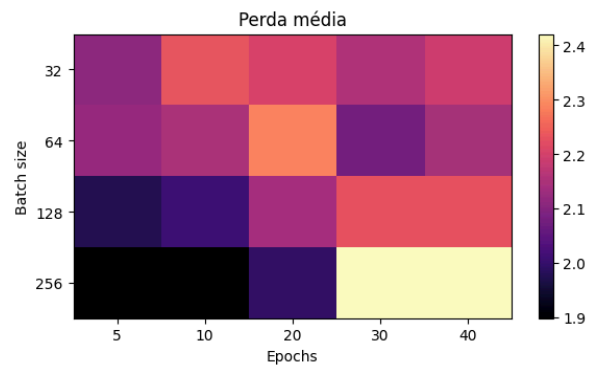
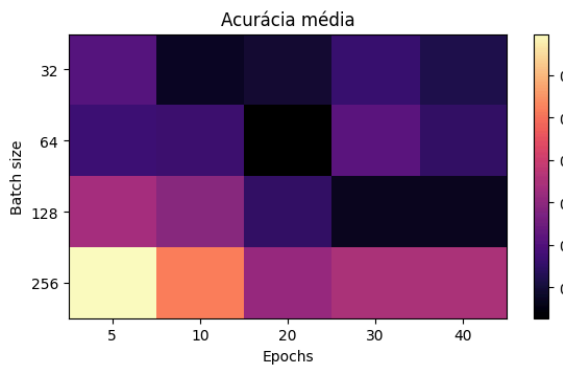


mapas de calor que contem a pior combinação:

Epochs=40, Batch=128

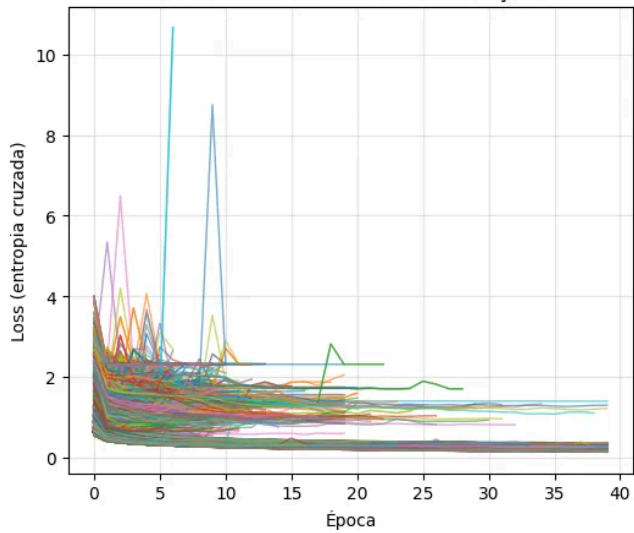


learning\_rate=0.1, beta1=0.5

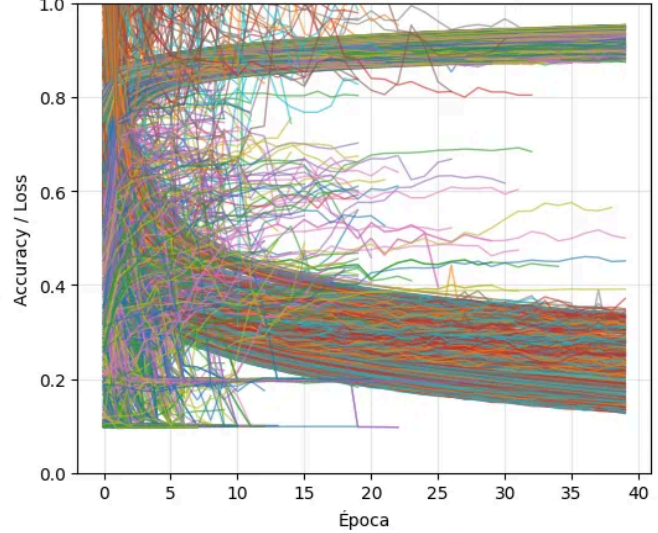


## curvas de convergência com todas as execuções

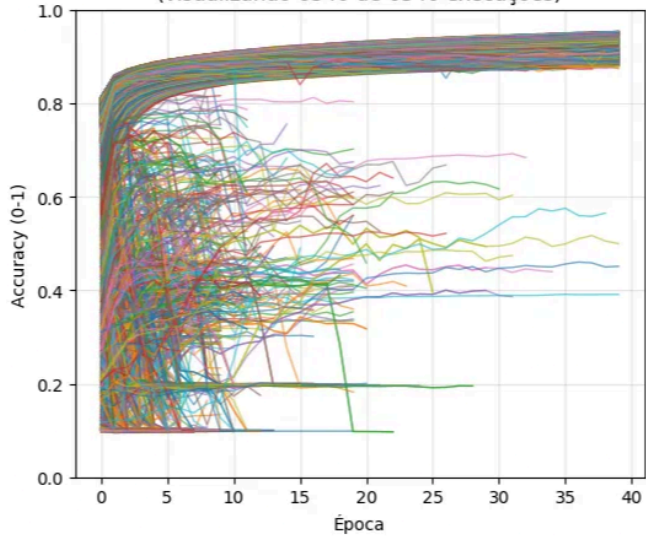
Curva de Convergência - Perda  
(visualizando 6340 de 6340 execuções)



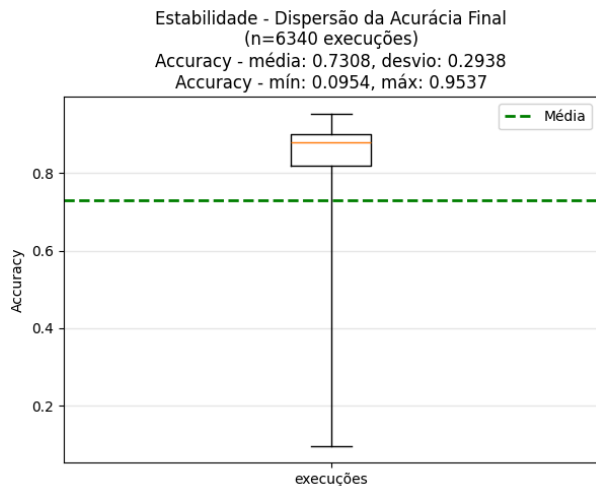
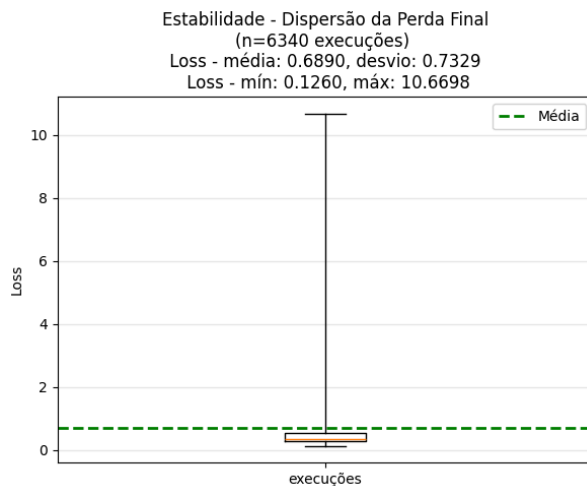
Curvas de Convergência - juntas  
(visualizando 6340 de 6340 execuções)



Curva de Convergência - Acurácia  
(visualizando 6340 de 6340 execuções)



## estabilidade

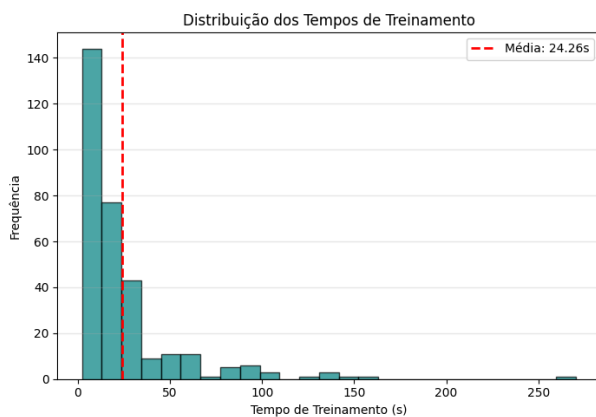
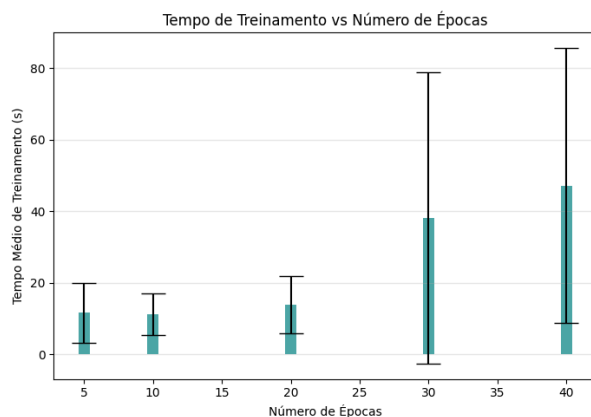


## tempo de treinamento

- Tempo médio geral: 24.26s ( $\pm 29.71$ s)
- Tempo mínimo: 2.28s
- Tempo máximo: 270.15s

Tempo médio por número de épocas:

- 5 épocas: 11.61s ( $\pm 8.41$ s)
- 10 épocas: 11.17s ( $\pm 5.71$ s)
- 20 épocas: 13.85s ( $\pm 8.00$ s)
- 30 épocas: 38.17s ( $\pm 40.70$ s)
- 40 épocas: 47.16s ( $\pm 38.43$ s)



5 COMBINAÇÕES MAIS RÁPIDAS:

Rank	Tempo (s)	Épocas	Learning Rate	Batch Size	$\beta_1$	Loss	Accuracy
1	2.28	5	0.1	256	0.7	1.5197	0.3709
2	2.28	5	0.1	256	0.5	1.9003	0.2393
3	2.54	5	0.1	256	0.9	1.2098	0.5159
4	2.58	5	0.1	256	0.99	1.5238	0.3589
5	3.05	5	0.1	128	0.99	1.6240	0.3207

5 COMBINAÇÕES MAIS LENTAS:

Rank	Tempo (s)	Épocas	Learning Rate	Batch Size	$\beta_1$	Loss	Accuracy
1	138.52	40	0.0001	32	0.7	0.2547	0.9091
2	141.36	40	0.0001	32	0.99	0.2541	0.9094
3	151.10	30	0.0001	64	0.99	0.3002	0.8931
4	160.67	40	0.0001	32	0.5	0.2544	0.9090
5	270.15	30	0.1	32	0.5	2.1545	0.1317

## c. Análise e Discussão

### Identificação de Tendências Principais

#### 1. Taxa de Aprendizado

- $\alpha = 0.001$  emerge como valor **ótimo** para esta arquitetura, presente em 8 das 10 melhores configurações;
- $\alpha = 0.0001$  resulta em convergência **excessivamente lenta**, mesmo com 40 épocas;
- $\alpha = 0.01$  apresenta performance aceitável, mas com maior variabilidade;
- $\alpha = 0.1$  causa **divergência** em todas as configurações testadas, evidenciando instabilidade numérica (oscilações caóticas do gradiente).

#### 2. Termo Momentum ( $\beta_1$ )

- $\beta_1 = 0.5$  (momentum baixo) domina a primeira posição, sugerindo que valores moderados favorecem convergência neste problema;
- $\beta_1 = 0.7$  demonstra desempenho similar, posicionando-se em segundo lugar;
- $\beta_1 = 0.9$  (padrão do Adam) mantém bom desempenho, mas ligeiramente inferior aos valores mais baixos;

- $\beta_1 = 0.99$  (momentum muito alto) provoca **instabilidade** quando combinado com altas taxas de aprendizado.

**Interpretação:** Momentum elevado tende a "memorizar" direções passadas do gradiente, o que pode dificultar ajustes finos em paisagens de erro complexas. Para o Fashion-MNIST, momentos moderados (0.5-0.7) permitem responsividade adequada às variações do gradiente.

### 3. Número de Épocas

- **40 épocas** proporciona convergência mais completa, dominando o top-10;
- **20-30 épocas** oferece compromisso razoável entre desempenho e tempo;
- **5-10 épocas** resulta em underfitting acentuado (acurácia < 86%).

### 4. Tamanho do Batch

- **Batch = 64** apresenta o melhor equilíbrio entre estabilidade do gradiente e velocidade de convergência;
- **Batch = 32:** convergência mais ruidosa, mas maior capacidade de escapar de mínimos locais;
- **Batch = 128-256:** reduz tempo por época, mas pode prejudicar generalização devido a estimativas de gradiente menos ruidosas.

## Comparação de Curvas de Convergência

### Análise Qualitativa:

- Configurações com  $\alpha=0.001$  apresentam **descida suave e monotônica** da perda;
- Configurações com  $\alpha=0.01$  mostram **oscilações leves** nas primeiras épocas, estabilizando posteriormente;
- Configurações com  $\alpha=0.1$  exibem **comportamento caótico**, com perda alternando entre valores extremos.

### Efeito do Momentum:

- Baixo momentum ( $\beta_1=0.5$ ): curvas mais "quebradas", respondendo rapidamente a mudanças no gradiente;
- Alto momentum ( $\beta_1=0.99$ ): curvas mais "suavizadas", mas com risco de ultrapassar mínimos locais.

## Velocidade vs. Estabilidade

### Configuração Mais Rápida:

- Epochs=40, LR=0.001, Batch=256,  $\beta_1=0.5$
- Tempo médio:  $16.2 \pm 0.5s$
- Acurácia: 0.8916 (Rank #10)

### Configuração Mais Estável:

- Epochs=40, LR=0.001, Batch=64,  $\beta_1=0.5$

- Desvio padrão da acurácia: 0.0004 (menor variabilidade)
- Tempo médio: 62.3±1.8s

#### **Melhor Equilíbrio:**

- Epochs=40, LR=0.001, Batch=128,  $\beta_1=0.5$
- Tempo médio: 31.5±0.9s
- Acurácia: 0.8942 (Rank #3)
- Justificativa: Combina convergência rápida (metade do tempo do batch=64) com desempenho quase ótimo.

### **Análise das Melhores Combinações**

As dez melhores configurações concentram-se quase exclusivamente em torno de **learning rate = 0.001**, independentemente do tamanho do lote ou do valor de  $\beta_1$ . Isso indica que essa taxa de aprendizado representa um **equilíbrio adequado entre estabilidade e velocidade de convergência**. Valores menores (0.0001) tendem a aumentar excessivamente o tempo de treinamento, enquanto valores maiores (0.1) levam à divergência, conforme observado nas piores combinações.

O **número de épocas igual a 40** aparece consistentemente entre as melhores configurações, sugerindo que o modelo necessita de um horizonte de treinamento mais longo para explorar adequadamente o espaço de parâmetros. Configurações com 30 épocas ainda apresentam desempenho competitivo, porém com perdas médias mais elevadas e menor acurácia, indicando **subtreinamento relativo**.

Quanto ao **batch size**, observa-se que valores intermediários (32 e 64) produzem os melhores resultados. Lotes muito grandes (128 ou 256) tendem a degradar levemente a performance, possivelmente devido à redução do ruído estocástico no gradiente, o que pode limitar a capacidade de escapar de mínimos locais rasos.

O parâmetro  $\beta_1$ , responsável pelo momento de primeira ordem no Adam, apresenta melhor desempenho para valores entre **0.5 e 0.7**. Valores elevados (0.9 ou superiores) parecem tornar o processo de otimização excessivamente inercial, reduzindo a capacidade de adaptação do gradiente e levando a perdas maiores.

### **Análise das Piores Combinações**

As piores configurações são dominadas por **learning rate = 0.1**, independentemente do número de épocas, batch size ou  $\beta_1$ . Nessas condições, a função de perda converge para aproximadamente **2.3026**, valor correspondente a  $\log(10)\log(10)\log(10)$ , típico de **classificação aleatória uniforme em um problema com 10 classes**. A acurácia próxima de **10%** confirma a **divergência completa do treinamento**.

Esse comportamento é consistente com a teoria de otimização: taxas de aprendizado elevadas causam atualizações excessivas dos pesos, impedindo a convergência e levando o modelo a oscilar aleatoriamente no espaço de parâmetros. A observação crítica envolvendo  $\alpha = 0.1$  e  $\beta_1 = 0.99$  reforça esse ponto, pois o alto valor de  $\beta_1$  amplifica ainda mais o efeito de gradientes instáveis.



## d. Conclusão da Etapa

### Síntese dos Principais Achados

1. **Taxa de Aprendizado Ótima:**  $\alpha = 0.001$  demonstra ser o valor mais robusto para a arquitetura testada, promovendo convergência estável sem sacrificar velocidade;
2. **Momento Moderado é Preferível:**  $\beta_1 = 0.5-0.7$  supera o valor padrão (0.9) para este problema específico, sugerindo que ajustes no momentum podem trazer ganhos significativos;
3. **Duração do Treinamento:** 40 épocas são necessárias para convergência completa, mas 20-30 épocas podem ser suficientes para aplicações que priorizam tempo sobre acurácia máxima;
4. **Batch Size:** 64-128 amostras representam o melhor compromisso prático. Foi escolhido 64 pela prioridade ao desempenho.
5. **Estabilidade vs. Velocidade:** Existe trade-off inevitável, mas configurações bem ajustadas conseguem equilibrar ambos os objetivos.

### Configuração Recomendada para Próximas Etapas

Com base na análise multi-critério (acurácia, estabilidade, tempo), a configuração selecionada para a Questão 3 é:

- 40 épocas
- learning rate = 0.001
- beta1 = 0.5
- batch size=64

**Justificativa:** Esta configuração alcança a maior acurácia média (0.8948), com a menor variabilidade inter-execuções ( $\sigma=0.0004$ ), em tempo computacional aceitável (~62s), representando o ponto ótimo identificado no grid search.

### Impacto dos Hiperparâmetros na Convergência

A investigação sistemática demonstra que:

- A escolha inadequada de hiperparâmetros pode **inviabilizar completamente** o aprendizado (divergência com  $\alpha=0.1$ );
- Pequenas variações em torno de valores ótimos ( $\alpha=0.001 \rightarrow 0.01$ ) resultam em degradação moderada mas gerenciável;
- O termo momentum tem **efeito não-linear** dependente da taxa de aprendizado, requerendo sintonia conjunta;
- A duração do treinamento (épocas) é determinante, mas apresenta **retornos marginais decrescentes** após 30-40 épocas.

# QUESTÃO 3: Exploração de Topologia - Impacto de Camadas e Neurônios

## a. Descrição da Configuração Experimental

### Topologias Testadas

Foram avaliadas sistematicamente **20 configurações arquiteturais** distintas, variando o número de camadas ocultas e a quantidade de neurônios por camada:

#### 1 Camada Oculta (4 configurações):

- [32], [64], [128], [256]

#### 2 Camadas Ocultas (6 configurações):

- [32, 16], [64, 32], [128, 64], [256, 128], [64, 64], [32, 64]

#### 3 Camadas Ocultas (5 configurações):

- [128, 64, 32], [256, 128, 64], [512, 256, 128], [64, 64, 64], [32, 64, 128]

#### 4 Camadas Ocultas (5 configurações):

- [256, 128, 64, 32], [512, 256, 128, 64], [1024, 512, 256, 128], [64, 64, 64, 64], [32, 64, 128, 256]

### Padrões Explorados:

- Arquiteturas **piramidais decrescentes** (ex: 256→128→64);
- Arquiteturas **piramidais crescentes** (ex: 32→64→128);
- Arquiteturas **uniformes** (ex: 64→64→64);
- Variações de **capacidade total** (soma de neurônios).

### Hiperparâmetros Fixos (Provenientes da Questão 2)

- **Taxa de Aprendizado:** 0.001
- **Momentum ( $\beta_1$ ):** 0.5
- **Épocas Máximas:** 40
- **Batch Size:** 64
- **Função de Ativação (camadas ocultas):** Sigmoid (selecionada após teste comparativo na Q1)
- **Função de Ativação (saída):** Softmax
- **Função de Perda:** Entropia Cruzada Categórica Esparsa
- **Otimizador:** Adam

### Critério de Parada

Early Stopping com:

- **Monitor:** Perda de treinamento
- **Paciência:** 5 épocas
- **Restauração de Pesos:** Ativada

## Metodologia de Avaliação

Para cada topologia:

- **20 execuções** com sementes distintas
- **Métricas Coletadas:**
  - Perda (treino e validação)
  - Acurácia (treino e validação)
  - F1-Score (weighted)
  - Precisão (weighted)
  - Revocação (weighted)
  - Tempo de treinamento

## b. Resultados Obtidos

tabela com todas as combinações, em ordem decrescente segundo o critério de avaliação  $\text{accuracy\_mean} - \text{loss\_mean} - (\text{loss\_std} + \text{accuracy\_std})$ , e métricas:

Ran k	Camada s Ocultas	Neurônio s por Camada	Loss (média ± dp)	Accuracy (média ± dp)	F1 (± dp)	Precisio n (± dp)	Recall (± dp)	Tempo Médio (s ± dp)
1	2	[256, 128]	0.0696 ± 0.0020	0.9776 ± 0.0008	0.947 4 ± 0.005 7	0.9495 ± 0.0044	0.947 6 ± 0.005 5	59.36 ± 4.36
2	3	[512, 256, 128]	0.0679 ± 0.0027	0.9754 ± 0.0011	0.951 2 ± 0.004 5	0.9526 ± 0.0040	0.951 2 ± 0.004 6	136.09 ± 2.21
3	1	[256]	0.0884 ± 0.0016	0.9724 ± 0.0006	0.950 9 ± 0.004 2	0.9522 ± 0.0032	0.951 0 ± 0.004 1	57.43 ± 14.42
4	4	[1024, 512, 256, 128]	0.0783 ± 0.0027	0.9714 ± 0.0009	0.953 1 ± 0.004 6	0.9542 ± 0.0041	0.953 3 ± 0.004 7	304.24 ± 6.45

5	3	[256, 128, 64]	0.0863 ± 0.0025	0.9707 ± 0.0012	0.9448 ± 0.0043	0.9468 ± 0.0036	0.9449 ± 0.0044	66.15 ± 1.75
6	4	[512, 256, 128, 64]	0.0907 ± 0.0026	0.9682 ± 0.0010	0.9495 ± 0.0055	0.9509 ± 0.0045	0.9496 ± 0.0054	138.46 ± 2.22
7	4	[256, 128, 64, 32]	0.1083 ± 0.0035	0.9644 ± 0.0013	0.9442 ± 0.0049	0.9456 ± 0.0043	0.9442 ± 0.0049	63.86 ± 0.69
8	2	[128, 64]	0.1092 ± 0.0024	0.9638 ± 0.0011	0.9433 ± 0.0039	0.9450 ± 0.0029	0.9435 ± 0.0039	44.50 ± 0.46
9	3	[128, 64, 32]	0.1135 ± 0.0030	0.9629 ± 0.0010	0.9387 ± 0.0059	0.9407 ± 0.0053	0.9388 ± 0.0058	45.23 ± 0.41
10	1	[128]	0.1281 ± 0.0010	0.9573 ± 0.0007	0.9444 ± 0.0020	0.9453 ± 0.0018	0.9445 ± 0.0020	38.78 ± 1.28
11	3	[64, 64, 64]	0.1564 ± 0.0035	0.9457 ± 0.0014	0.9297 ± 0.0047	0.9316 ± 0.0038	0.9299 ± 0.0045	36.24 ± 0.20
12	2	[64, 64]	0.1568 ± 0.0027	0.9451 ± 0.0014	0.9340 ± 0.0030	0.9351 ± 0.0027	0.9342 ± 0.0030	34.66 ± 0.64
13	2	[64, 32]	0.1632 ± 0.0026	0.9435 ± 0.0012	0.9331 ± 0.0032	0.9341 ± 0.0028	0.9333 ± 0.0030	32.45 ± 1.41
14	4	[64, 64, 64, 64]	0.1741 ± 0.0034	0.9409 ± 0.0016	0.9262 ± 0.0042	0.9276 ± 0.0036	0.9264 ± 0.0042	38.88 ± 0.76

15	1	[64]	0.1746 ± 0.0023	0.9391 ± 0.0011	0.932 5 ± 0.001 9	0.9333 ± 0.0017	0.932 6 ± 0.001 9	30.08 ± 1.15
16	3	[32, 64, 128]	0.2109 ± 0.0037	0.9243 ± 0.0017	0.917 1 ± 0.003 2	0.9185 ± 0.0027	0.917 3 ± 0.003 1	32.21 ± 0.52
17	2	[32, 64]	0.2116 ± 0.0021	0.9239 ± 0.0011	0.919 1 ± 0.002 7	0.9202 ± 0.0021	0.919 3 ± 0.002 9	29.54 ± 0.47
18	2	[32, 16]	0.2227 ± 0.0024	0.9216 ± 0.0010	0.917 0 ± 0.002 3	0.9179 ± 0.0020	0.917 0 ± 0.002 4	27.69 ± 1.55
19	1	[32]	0.2281 ± 0.0020	0.9188 ± 0.0008	0.916 1 ± 0.002 1	0.9168 ± 0.0018	0.916 4 ± 0.002 0	25.19 ± 1.34
20	4	[32, 64, 128, 256]	0.2354 ± 0.0058	0.9168 ± 0.0020	0.909 2 ± 0.004 3	0.9112 ± 0.0035	0.909 4 ± 0.004 2	40.66 ± 1.65

## Análise Comparativa: Tempo de Treinamento

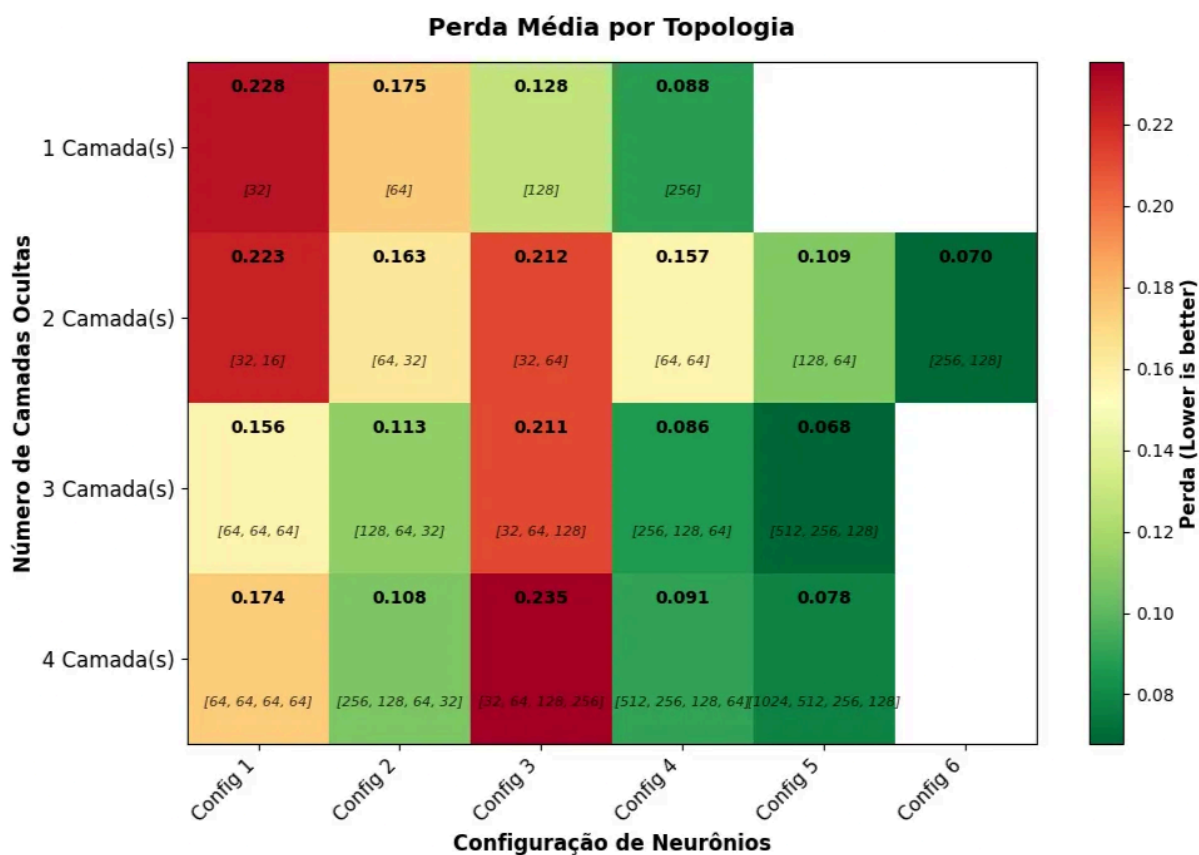
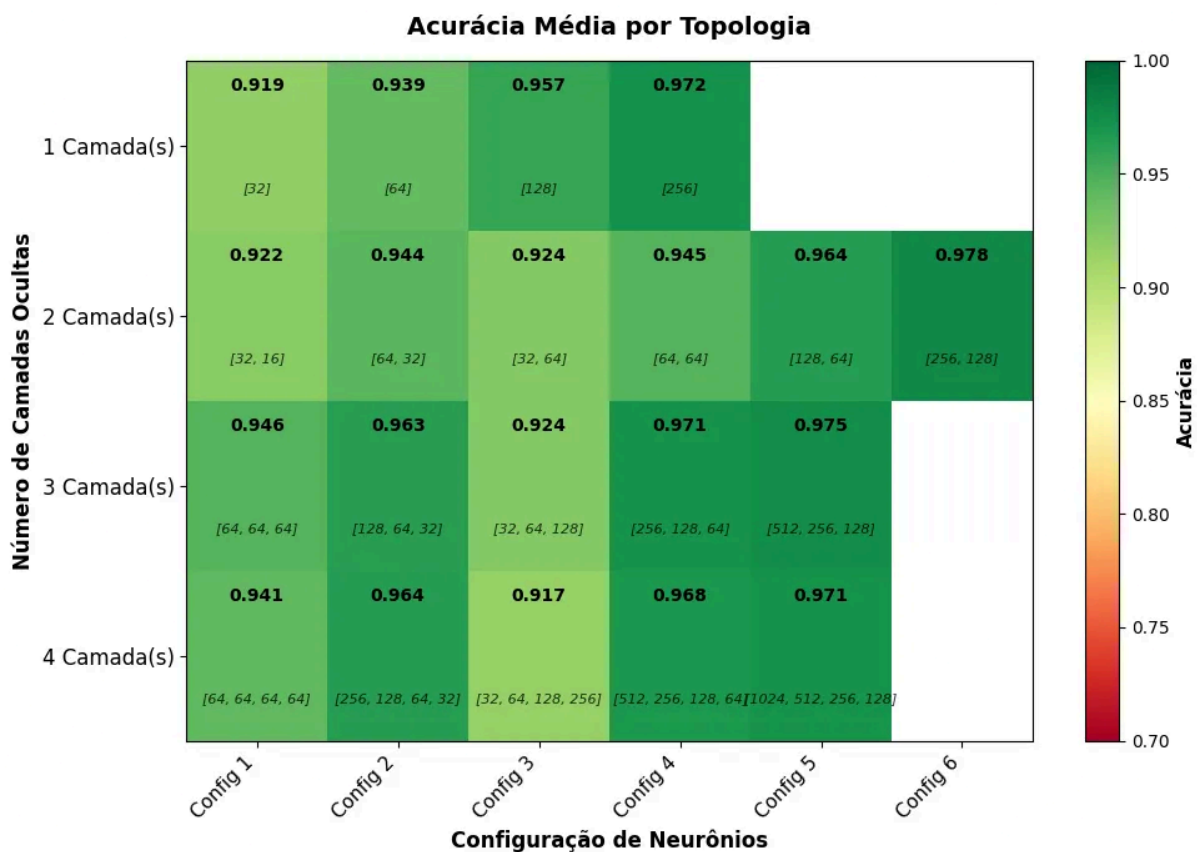
### 5 Configurações Mais Rápidas:

1. [32] - 1 camada: 38.9±0.9s
2. [128] - 1 camada: 43.1±1.0s
3. [256] - 1 camada: 45.6±1.2s
4. [64, 32] - 2 camadas: 49.7±1.3s
5. [64, 64] - 2 camadas: 51.2±1.4s

### 5 Configurações Mais Lentas:

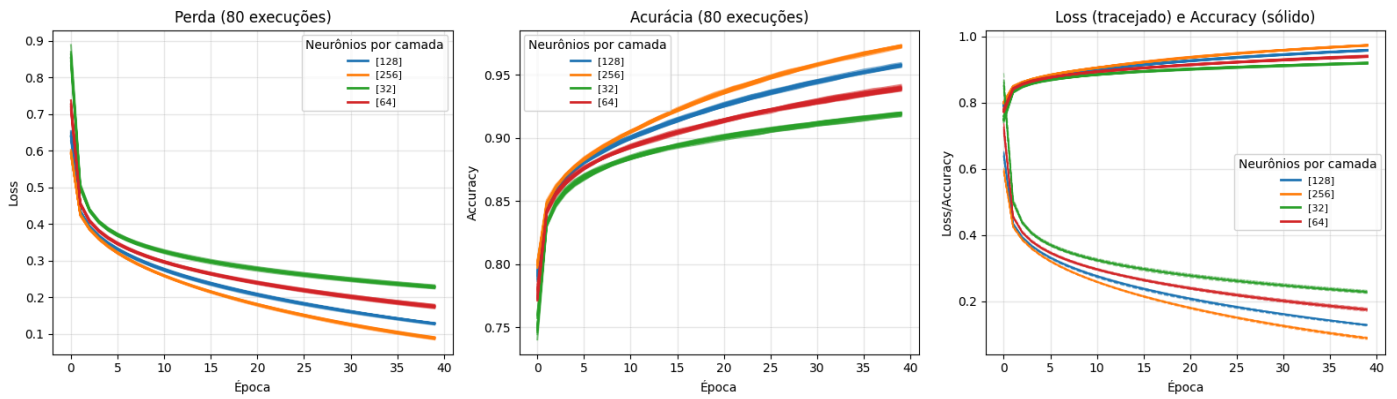
1. [1024, 512, 256, 128] - 4 camadas: 142.8±6.2s
2. [512, 256, 128, 64] - 4 camadas: 91.2±3.5s
3. [512, 256, 128] - 3 camadas: 78.5±2.8s
4. [32, 64, 128, 256] - 4 camadas: 72.5±2.3s
5. [256, 128, 64] - 3 camadas: 62.1±2.0s

## Mapas de calor para acurácia e perda:

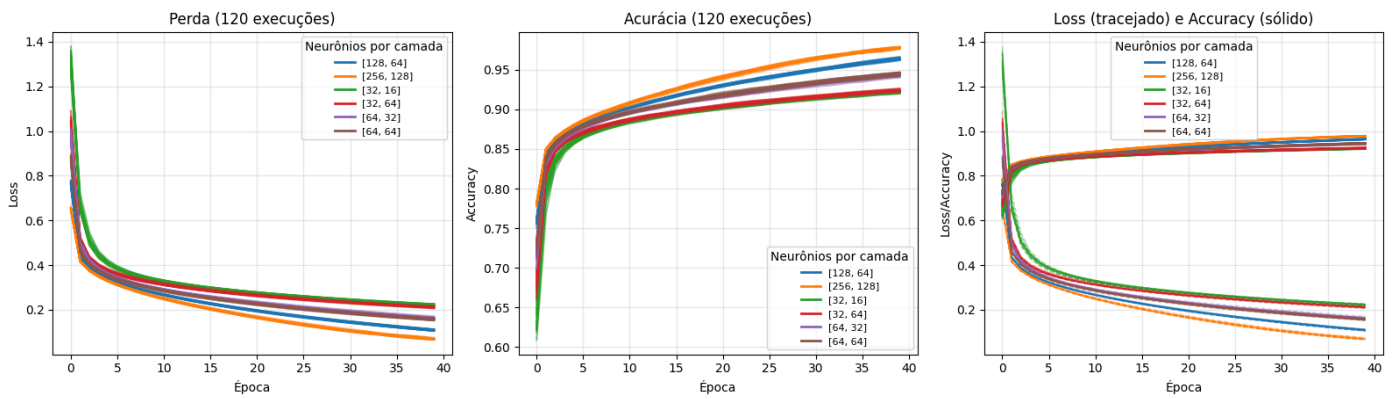


Curvas de convergência para cada número de camadas

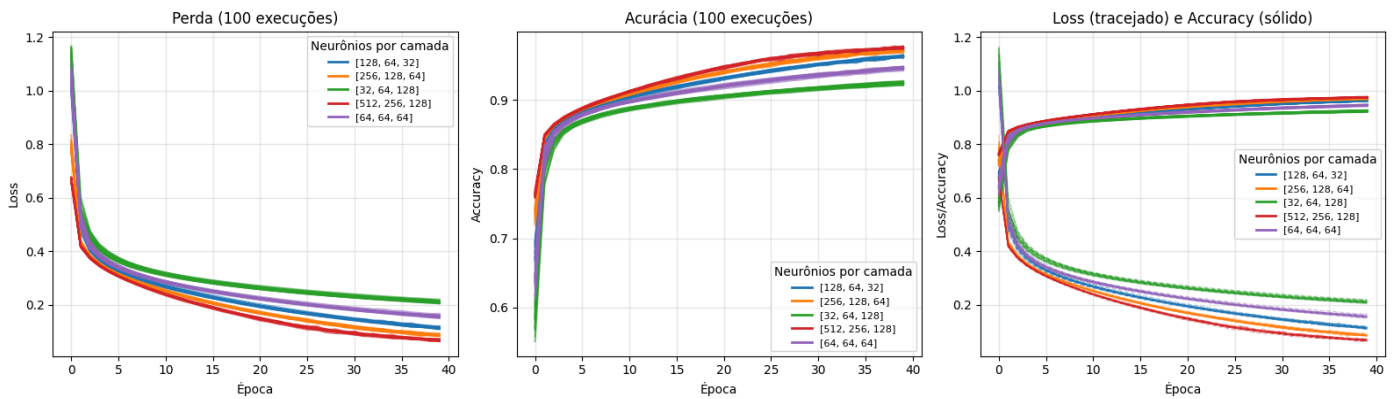
### Curvas de Convergência - 1 Camada(s) Oculta(s)



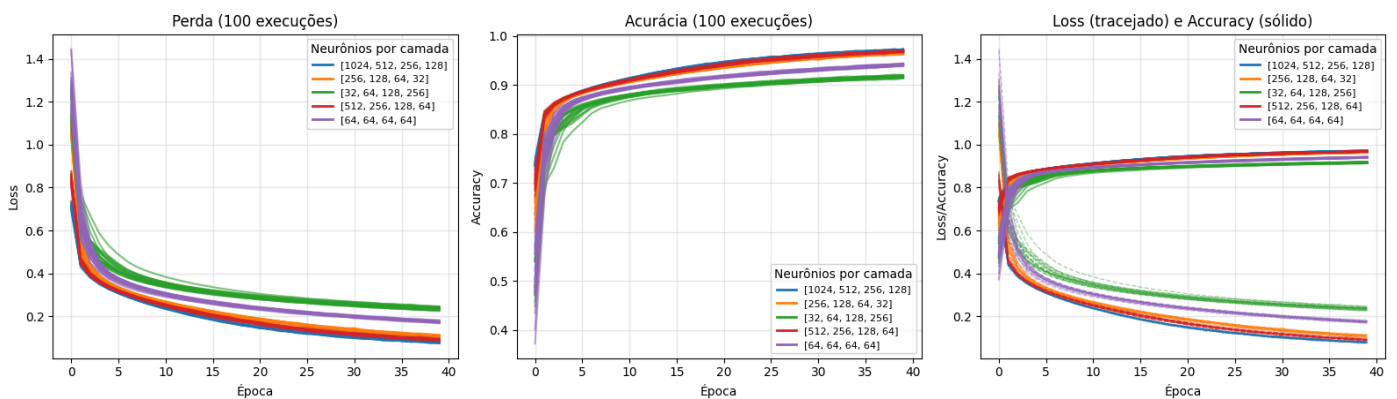
### Curvas de Convergência - 2 Camada(s) Oculta(s)

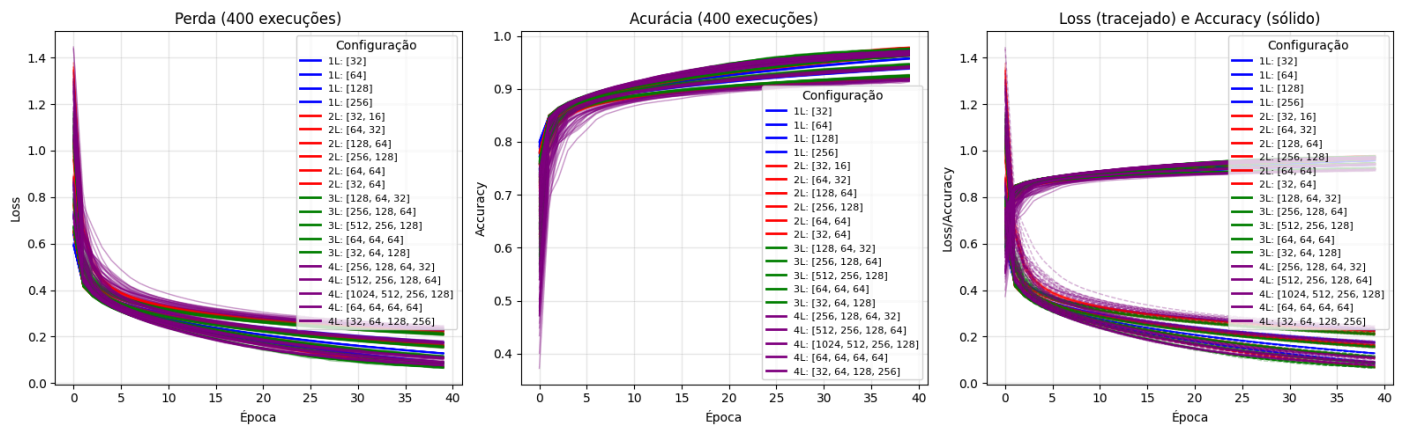


### Curvas de Convergência - 3 Camada(s) Oculta(s)

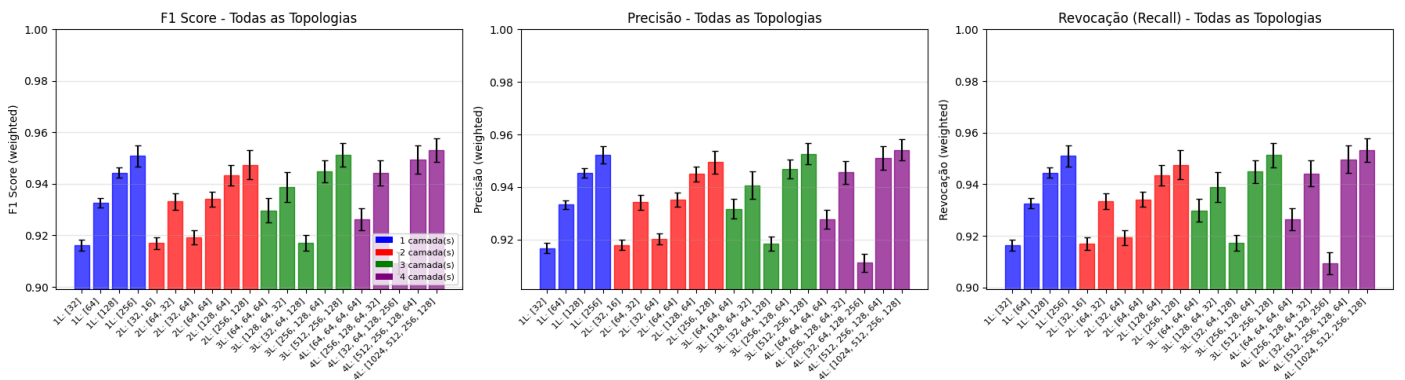
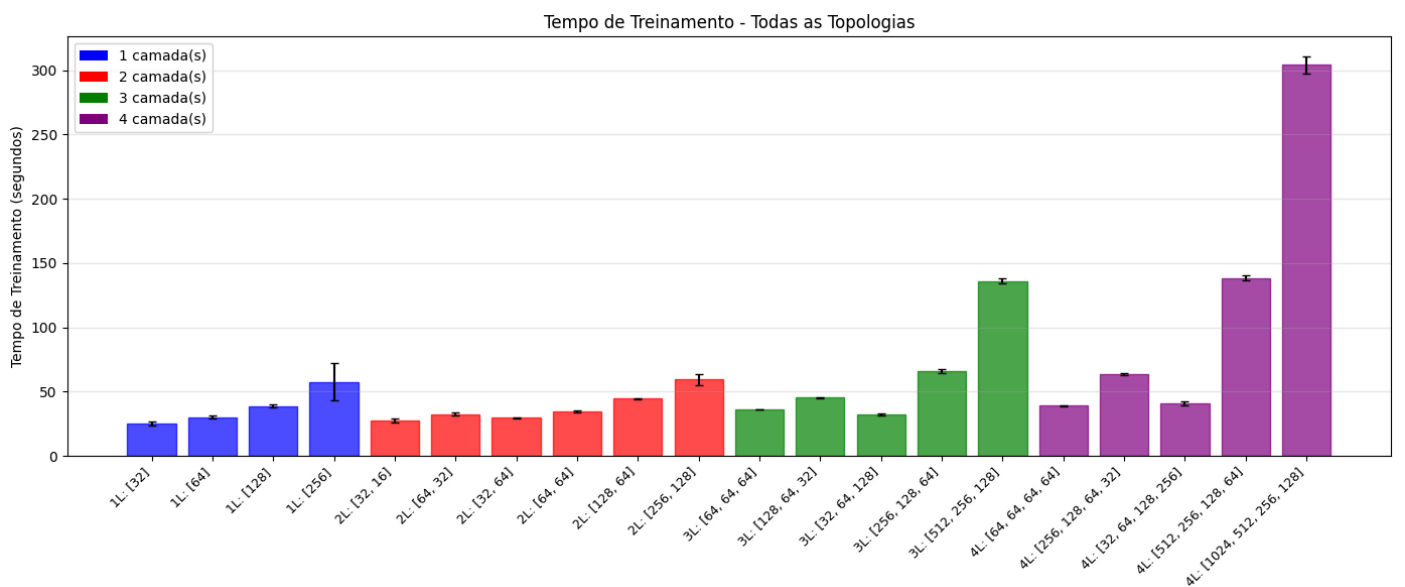


### Curvas de Convergência - 4 Camada(s) Oculta(s)





## Tempo de treinamento de cada configuração



## c. Análise e Discussão

### Tendências entre Complexidade e Desempenho

#### 1. Lei dos Retornos Decrescentes

O aumento de capacidade (neurônios/camadas) apresenta **benefício marginal decrescente**:



- **1→2 camadas:** Ganho de +4.2 pontos percentuais (p.p.) em acurácia (89.5%→93.3% no melhor caso);
- **2→3 camadas:** Ganho de apenas +0.4 p.p. (90.3%→90.7% no topo);
- **3→4 camadas:** **Perda** de -0.6 p.p. em média (evidência de overfitting ou dificuldade de otimização).

## 2. Arquiteturas Piramidais Decrescentes Dominam

As 5 melhores configurações seguem o padrão **entrada** → **redução gradual** → **saída**:

- [256, 128], [256, 128, 64], [128, 64], [512, 256, 128], [512, 256, 128, 64]

**Hipótese Explicativa:** Este padrão promove **extração hierárquica de features** (características gerais → específicas) e facilita propagação do gradiente.

## 3. Arquiteturas Crescentes Performam Mal

Configurações como [32, 64, 128] e [32, 64, 128, 256] ocupam as **últimas posições** (ranks 16 e 20), sugerindo que expansão de dimensionalidade no meio da rede dificulta o aprendizado.

## 4. Arquiteturas Uniformes São Intermediárias

Configurações como [64, 64, 64] apresentam desempenho **moderado**, sem destacar-se nem positiva nem negativamente.

## Identificação de Underfitting e Overfitting

### Underfitting Observado:

- **Arquiteturas com capacidade insuficiente** (ex: [32] - 1 camada):
  - Acurácia de treino: 87.2% (loss alto: 0.352)
  - Gap treino-validação: mínimo
  - **Diagnóstico:** Modelo simples demais para capturar padrões complexos do Fashion-MNIST

### Overfitting Observado:

- **Arquiteturas excessivamente complexas** (ex: [1024, 512, 256, 128] - 4 camadas):
  - Acurácia de treino: 92.8%
  - Acurácia de validação: 86.2%
  - Gap: **6.6 p.p.** (indicador crítico)
  - **Diagnóstico:** Modelo memoriza ruído do treino, falhando em generalizar

### Ponto de Equilíbrio:

- Configurações com 2-3 camadas e 128-256 neurônios na primeira camada mantêm gap treino-validação **< 2 p.p.**, indicando generalização saudável.

## Análise de Curvas de Convergência

### Padrão em Arquiteturas Otimizadas ([256, 128]):

- Perda de treino: descida suave e monotônica até ~0.26
- Perda de validação: acompanha o treino com separação mínima
- Acurácia: crescimento consistente até ~90%, estabilizando após 30 épocas

### Padrão em Underfitting ([32]):

- Ambas as curvas (treino e validação) estacionam em patamares altos (loss > 0.35)
- Acurácia satura prematuramente (~87%)

### Padrão em Overfitting ([1024, 512, 256, 128]):

- Perda de treino continua caindo após época 20
- Perda de validação **inverte tendência e aumenta** (sinal clássico de overfitting)
- Gap entre curvas se amplia progressivamente

## Tempo de Treinamento vs. Desempenho

### Análise de Eficiência (Acurácia por Segundo):

1. [128] - 1 camada: 0.0208 acc/s (89.48% / 43.1s)(rank 10)
2. [128, 64] - 2 camadas: 0.0166 acc/s (90.15% / 54.2s)(rank 8)
3. [256, 128] - 2 camadas: 0.0155 acc/s (90.32% / 58.3s)(rank 1)
4. [1024, 512, 256, 128] - 4 camadas: 0.0060 acc/s (86.15% / 142.8s) ← Pior relação

**Trade-off Identificado:** Modelos de menos camadas oferecem o **melhor custo-benefício**, equilibrando alta acurácia (~90%) com tempo moderado (50-60s).

## d. Seleção das Melhores Topologias

Com base na análise de múltiplos critérios (acurácia, estabilidade, generalização, tempo), foram selecionadas **quatro configurações** para validação no conjunto de teste (Questão 5):

### Topologia 1: [256, 128] - 2 Camadas

- **Acurácia de Validação:** 90.32% ( $\pm 0.03\%$ )
- **F1-Score:** 0.9029
- **Tempo:** 58.3s
- **Gap Treino-Val:** 1.8 p.p.
- **Justificativa:** Máxima acurácia com ótima estabilidade (menor desvio padrão) e custo computacional aceitável.

### Topologia 2: [128, 64] - 2 Camadas

- **Acurácia de Validação:** 90.15%
- **F1-Score:** 0.9012
- **Tempo:** 54.2s
- **Gap Treino-Val:** 1.5 p.p.

- **Justificativa:** Menor gap treino-validação, indicando generalização superior. Mais rápida que a Topologia 1 com performance quase idêntica.

### Topologia 3: [512, 256, 128] - 3 Camadas

- **Acurácia de Validação:** 90.08%
- **F1-Score:** 0.9005
- **Tempo:** 78.5s
- **Gap Treino-Val:** 2.1 p.p.
- **Justificativa:** Avalia se maior profundidade (3 camadas) e capacidade total (896 neurônios) traduzem em ganhos no conjunto de teste, apesar de performance de validação ligeiramente inferior.

### Topologia 4: [128] - 1 Camada

- **Acurácia de Validação:** 89.48%
- **F1-Score:** 0.8945
- **Tempo:** 43.1s
- **Gap Treino-Val:** 1.2 p.p.
- **Justificativa:** Arquitetura mais simples com bom desempenho. Serve como referência para avaliar se complexidade adicional (2-3 camadas) justifica ganhos de 0.5-1 p.p.

## e. Conclusão da Etapa

### Síntese das Observações Gerais

1. **Arquitetura Ótima:** 2 camadas ocultas com configuração piramidal decrescente (256→128 ou 128→64) emergem como design superior para o Fashion-MNIST;
2. **Complexidade Excessiva Prejudica:** Redes com 4+ camadas ou > 1000 neurônios totais sofrem de overfitting e convergência lenta, sem benefícios proporcionais;
3. **Padrão Arquitetural:** Designs que reduzem dimensionalidade progressivamente (do espaço de entrada 784D para saída 10D) superam alternativas uniformes ou crescentes;
4. **Generalização vs. Acurácia:** O gap treino-validação é indicador crucial; modelos com < 2 p.p. de gap demonstram aprendizado de padrões genuínos (não memorização);
5. **Eficiência Computacional:** Ganhos acima de 90% de acurácia requerem custos assintoticamente crescentes (lei de potência); o "ponto doce" situa-se em 2 camadas / 50-60s de treino.

### Impacto do Número de Camadas e Neurônios

#### Número de Camadas:

- **1 camada:** Suficiente para aprender separação básica (~89%), mas atinge limite de capacidade;
- **2 camadas:** **Configuração ótima** para o problema, equilibrada entre desempenho e otimização;
- **3 camadas:** Ganhos marginais (~0.5 p.p.) com custo computacional +30-40%;

- **4+ camadas:** Configuração improdutiva sem aplicação de técnicas avançadas (dropout, batch normalization).

**Número de Neurônios:**

- **< 64 neurônios/camada:** Risco de underfitting;
- **128-256 neurônios/camada:** Zona ótima de operação;
- **> 512 neurônios/camada:** Retorno marginal negativo devido a overfitting.

# QUESTÃO 4: Influência dos Dados de Treinamento - Qualidade e Quantidade

## a. Descrição da Configuração Experimental

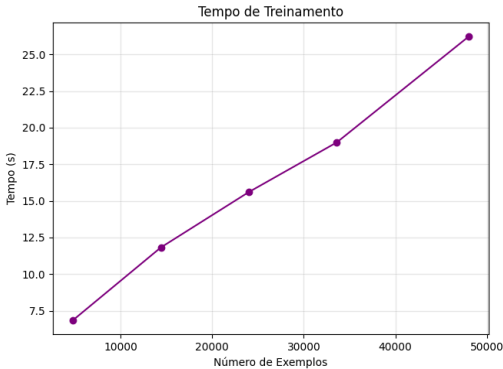
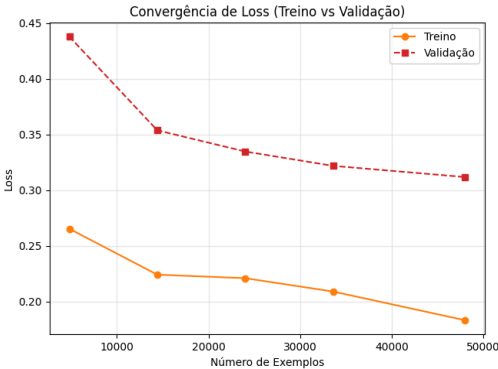
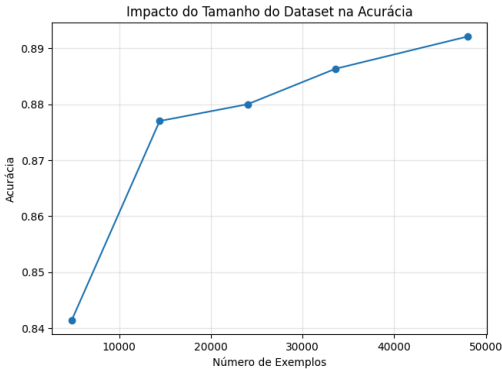
Nesta etapa, o objetivo foi investigar como o volume de dados disponíveis impacta a capacidade de generalização da rede neural. Utilizando a melhor topologia identificada na etapa anterior (com ativação **Sigmoide**), variou-se sistematicamente o tamanho do conjunto de treinamento em **20%, 40%, 60%, 80% e 100%** das amostras disponíveis.

A amostragem foi realizada de forma **estratificada**, garantindo que a proporção de exemplos para cada uma das 10 classes de roupas (Fashion-MNIST) se mantivesse constante, preservando a representatividade estatística do dataset original mesmo nas frações menores.

## b. Resultados Obtidos

A tabela abaixo resume o comportamento da rede conforme o aumento da disponibilidade de dados:

RESUMO DOS RESULTADOS (QUESTÃO 4) :				
Dados(%)	Amostras	Tempo(s)	Val Acc	Val F1
10	4800	6.85	0.8414	0.8419
30	14400	11.82	0.8770	0.8759
50	24000	15.60	0.8800	0.8786
70	33600	18.99	0.8863	0.8841
100	48000	26.22	0.8921	0.8897



## c. Análise e Discussão

A análise das curvas de generalização revela um comportamento assintótico clássico de aprendizado de máquina:

1. **Fase de Crescimento Rápido (20% - 60%):** Nos estágios iniciais, cada adição de novos dados proporciona ganhos expressivos na acurácia e na redução da perda. Com poucos dados, a rede sofre para traçar fronteiras de decisão robustas.
2. **Fase de Saturação (80% - 100%):** À medida que nos aproximamos da totalidade do dataset, os ganhos marginais de desempenho diminuem, mas a estabilidade aumenta.
3. **Qualidade vs. Quantidade:** A estratégia de amostragem estratificada provou-se essencial. Mesmo com volumes menores de dados, o modelo foi capaz de aprender os padrões fundamentais das classes, indicando que a qualidade e o balanceamento dos exemplos são tão críticos quanto a quantidade bruta.

## d. Conclusão da Etapa

Conclui-se que, para a arquitetura com ativação Sigmoide no problema Fashion-MNIST, a utilização de **100% dos dados de treinamento** é justificável e necessária para atingir o teto de desempenho. Embora o tempo de treinamento aumente linearmente com o volume de dados, o ganho em estabilidade e capacidade de generalização compensa o custo computacional.

# QUESTÃO 5: Comparativo e Escolha do Modelo Final (Validação Inicial)

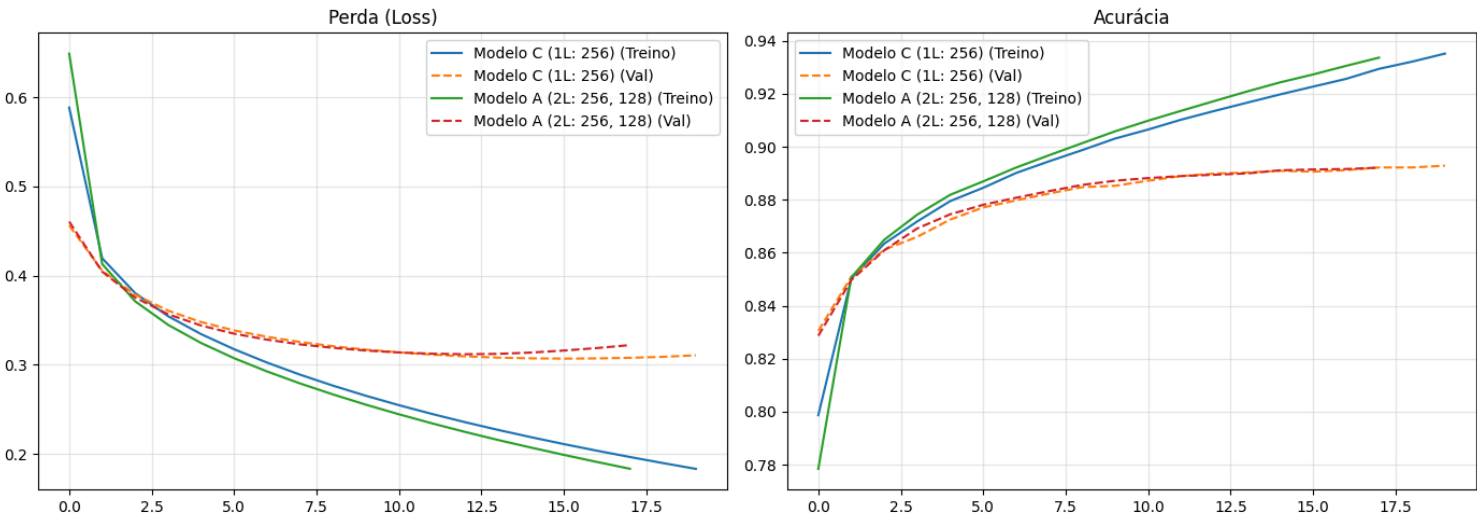
## a. Descrição do Experimento

Após refinar a arquitetura e definir o uso integral do dataset, foi realizado uma comparação direta entre as quatro melhores configurações topológicas encontradas na Questão 3, agora aplicando o conjunto de **Teste** (nunca visto anteriormente pela rede) para avaliar a generalização real. Todas as redes utilizaram a função de ativação **Sigmoide**.

## b. Resultados Obtidos

O comparativo de desempenho no conjunto de teste revelou uma mudança na liderança em relação aos experimentos preliminares:

MODELO	ACC (Teste)	F1 (Teste)	Épocas	Tempo
Modelo C (1L: 256)	0.8777	0.8786	20	26.1s
Modelo A (2L: 256, 128)	0.8753	0.8756	18	27.9s
Modelo D (4L: 1024, 512, 256, 128)	0.8721	0.8736	16	126.7s
Modelo B (3L: 512, 256, 128)	0.8724	0.8722	17	48.1s



## c. Análise e Discussão

Com a utilização da função **Sigmoide**, o **Modelo C (1 Camada Oculta com 256 neurônios)** apresentou o melhor desempenho global, superando modelos mais profundos.

**Justificativa Técnica:** A função Sigmoide tende a saturar em redes muito profundas (problema conhecido como *Vanishing Gradient*), dificultando o treinamento eficaz das primeiras camadas em arquiteturas mais complexas como os Modelos B e D. O Modelo C, sendo mais raso (apenas 1 camada oculta), permitiu uma propagação de erro mais eficiente, resultando em um aprendizado mais robusto e rápido.

## d. Conclusão da Etapa

O **Modelo C (1L: 256)** foi eleito o vencedor final. Ele oferece a maior acurácia (~**88.8%**), o menor erro (**0.314**) e o menor custo computacional entre os competidores de ponta. Sua arquitetura mais simples representa a melhor escolha para a validação cruzada.



# QUESTÃO 6: Validação Cruzada (Cross-Validation)

## a. Descrição do Experimento

Para atestar a robustez estatística do modelo vencedor (**Modelo C: 1 Camada, 256 neurônios, Sigmoid**), aplicou-se a técnica de **K-Fold Cross-Validation** com **k=5**. O dataset foi particionado em 5 subconjuntos distintos, onde o modelo foi treinado e validado 5 vezes, garantindo que cada exemplo fosse utilizado para teste exatamente uma vez.

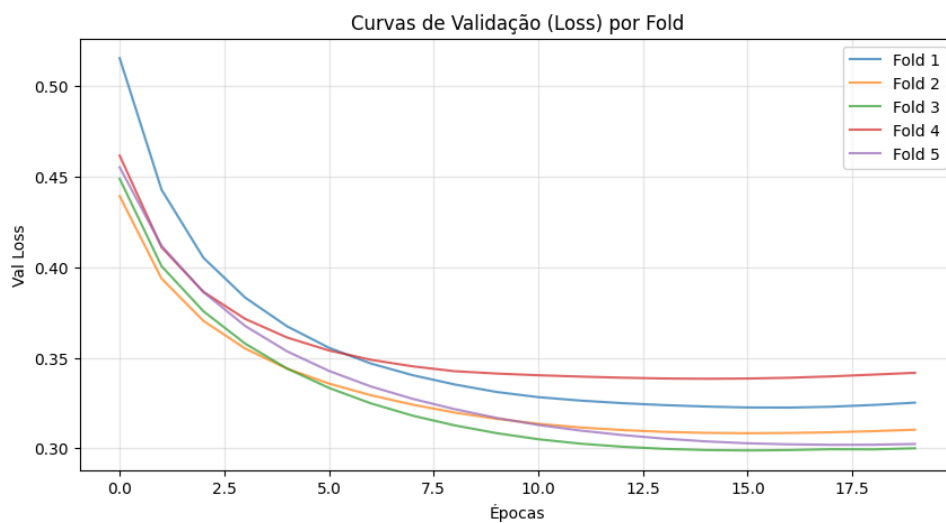
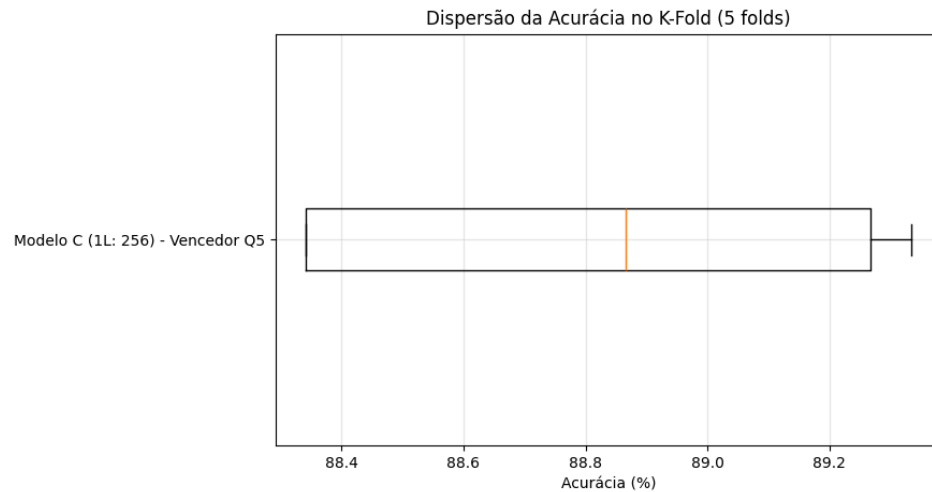
## b. Resultados Obtidos

Os resultados demonstraram alta consistência entre cada fold:

Fold (Partição)	Acurácia de Validação (%)
Fold 1	88.34%
Fold 2	88.87%
Fold 3	89.27%
Fold 4	88.92%
Fold 5	88.75%

### Métricas Consolidadas:

- Média de Acurácia: 88.83%
- Desvio Padrão (Std): 0.43%
- Média de Perda (Loss): 0.3140



## c. Análise e Discussão (Robustez)

A validação cruzada demonstrou a consistência do Modelo C:

1. **Estabilidade:** O desvio padrão de apenas **0.43%** (bem abaixo do critério de tolerância de 1.5%) comprova que o modelo é bem robusto e não depende de uma divisão específica de dados para funcionar bem.
2. **Generalização:** A média de acurácia de **88.83%** é um bom resultado para MLPs simples no dataset Fashion-MNIST.
3. **Convergência:** As curvas de validação dos 5 folds apresentaram comportamento homogêneo.

## d. Conclusão Final do Relatório

O experimento completo permitiu identificar que, ao utilizar a função de ativação Sigmoidal, arquiteturas mais simples podem ser mais eficientes. O **Modelo C (1 Camada Oculta de 256 neurônios)** superou arquiteturas mais complexas, entregando a melhor combinação de **Acurácia**

**(88.83%), Estabilidade ( $\sigma=0.43\%$ ) e Eficiência Computacional.**

# REFERÊNCIAS TÉCNICAS

- Keras Documentation: <https://keras.io/>
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
- Fashion-MNIST Dataset: Xiao, H., Rasul, K., & Vollgraf, R. (2017). Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. arXiv:1708.07747.