

Teoria da Informação: Fundamentos, Entropia e Aplicações na Era Digital

1. Introdução à Teoria da Informação

A Teoria da Informação, um campo que revolucionou a maneira como entendemos a comunicação e o processamento de dados, emergiu da necessidade premente de otimizar os sistemas de comunicação no início e meados do século XX. Antes do trabalho seminal de Claude Shannon, a análise de sistemas de comunicação era predominantemente focada em resolver problemas práticos de engenharia, muitas vezes utilizando ferramentas como a análise de Fourier para otimizar a transmissão de sinais em canais específicos.¹ Contudo, faltava um arcabouço matemático unificado que pudesse quantificar a informação e estabelecer os limites fundamentais da comunicação.

1.1. O Contexto Histórico e a Necessidade de uma Teoria

Engenheiros e cientistas enfrentam desafios crescentes na transmissão eficiente e confiável de mensagens através de diversos meios, como telégrafos, telefones e rádio. A interferência, ou ruído, era um problema onipresente, e a capacidade dos canais de comunicação parecia limitada de maneiras que não eram completamente compreendidas teoricamente. O trabalho anterior, embora valioso para aplicações específicas, carecia de uma generalidade que pudesse abranger diferentes tipos de mensagens e canais sob um mesmo conjunto de princípios. Foi nesse cenário que a necessidade de uma teoria matemática rigorosa da comunicação se tornou evidente, uma teoria que pudesse definir o que é informação, como medi-la e quais são os limites intransponíveis para sua compressão e transmissão.²

1.2. Claude Elwood Shannon: O Pai da Teoria da Informação

Claude Elwood Shannon (1916-2001), um matemático e engenheiro eletricista americano, é universalmente reconhecido como o "pai da teoria da informação".³ Sua trajetória acadêmica e profissional foi marcada por contribuições fundamentais que lançaram as bases para a era digital. Após graduar-se pela Universidade de Michigan, Shannon ingressou no Massachusetts Institute of Technology (MIT), onde, entre outras atividades, trabalhou com Vannevar Bush em seu analisador diferencial.¹ Uma experiência crucial foi seu estágio de verão nos Laboratórios Bell da American Telephone and Telegraph em 1937, que inspirou muitos de seus interesses de pesquisa subsequentes.¹

Em sua tese de mestrado de 1937 (publicada em 1938), "A Symbolic Analysis of Relay

and Switching Circuits", Shannon demonstrou de forma pioneira como a álgebra booleana poderia ser aplicada ao design e análise de circuitos de comutação elétrica, estabelecendo a base teórica para os circuitos digitais e, por extensão, para a computação moderna.¹ Este trabalho é considerado por muitos como uma das teses de mestrado mais importantes de todos os tempos. Shannon continuou sua pesquisa nos Laboratórios Bell, onde também fez contribuições significativas para a criptografia durante a Segunda Guerra Mundial, incluindo o artigo "Communication Theory of Secrecy Systems" (1949), que é uma peça fundamental da criptografia moderna.³ Sua genialidade eclética era notória, envolvendo desde malabarismos em um monociclo pelos corredores dos Bell Labs até a construção de máquinas que aprendiam, como o "Teseu", um rato mecânico capaz de encontrar a saída de um labirinto.¹

1.3. O Problema Fundamental da Comunicação

No cerne do trabalho de Shannon estava o que ele definiu como o "problema fundamental da comunicação": reproduzir em um ponto, seja exata ou aproximadamente, uma mensagem selecionada em outro ponto.⁴ Crucialmente, Shannon fez uma distinção fundamental ao separar os aspectos técnicos da transmissão da mensagem dos aspectos semânticos, ou seja, do significado da mensagem.¹ Ele argumentou que, do ponto de vista da engenharia, o significado da mensagem era, em grande parte, irrelevante para o problema de sua transmissão eficiente e confiável. O foco deveria ser nas características estatísticas da fonte da mensagem e nas propriedades do canal de comunicação. Esta abstração permitiu que ele desenvolvesse uma teoria matemática de grande aplicabilidade, capaz de modelar desde sistemas discretos (como a telegrafia) até sistemas contínuos (como o rádio e a telefonia).¹

1.4. "A Mathematical Theory of Communication" (1948): O Marco Zero

Em 1948, Claude Shannon publicou seu artigo seminal "A Mathematical Theory of Communication" no Bell System Technical Journal.² Este trabalho não apenas formalizou os conceitos de informação e entropia, mas também introduziu teoremas fundamentais sobre os limites da compressão de dados (codificação de fonte) e da taxa de transmissão de dados confiável através de canais ruidosos (codificação de canal). O impacto deste artigo foi imediato e profundo, fornecendo aos engenheiros de comunicação um modo de determinar a capacidade de um canal de comunicação e uma base teórica para o desenvolvimento de técnicas de codificação mais eficientes.¹ Referido como um "projeto para a era digital" e a "Magna Carta da Era da Informação", o trabalho de Shannon influenciou o desenvolvimento de tecnologias onipresentes como o disco compacto (CD), a internet, a telefonia móvel e até mesmo

a compreensão de fenômenos físicos como buracos negros.³

2. O Conceito de Informação

Para desenvolver uma teoria matemática da comunicação, Shannon precisou primeiro definir rigorosamente o que é "informação". Sua abordagem, contudo, divergiu radicalmente da noção intuitiva de informação ligada ao significado ou conteúdo semântico de uma mensagem.

2.1. O Que NÃO é Informação (no Sentido de Shannon): Semântica vs. Probabilidade

A teoria de Shannon deliberadamente ignora o significado das mensagens.⁵ Para ele, a questão não era o *que* a mensagem dizia, mas sim a probabilidade de uma determinada mensagem ser escolhida de um conjunto de mensagens possíveis. O foco está nos aspectos sintáticos e estatísticos da comunicação, não nos semânticos. A teoria visa otimizar a atribuição de símbolos (como bits) às mensagens e projetar algoritmos de codificação e decodificação eficientes para transmitir esses símbolos de forma confiável, independentemente do seu conteúdo interpretativo.⁵

- **Analogia 1 (O Carteiro Eficiente):** Imagine um carteiro encarregado de entregar cartas. Para o carteiro realizar seu trabalho eficientemente (entregar a carta correta no endereço correto, sem perdas ou danos), ele não precisa ler ou entender o conteúdo das cartas. Seu foco está no envelope, no endereço e na rota mais eficiente. Da mesma forma, a teoria da informação se concentra na "embalagem" e "transporte" da mensagem (os símbolos e o canal), não no seu "conteúdo" (o significado).
- **Analogia 2 (O Tradutor Universal de Símbolos):** Pense na teoria da informação como um sistema de tradução universal, mas não de idiomas, e sim de sequências de símbolos. Ela busca a maneira mais eficiente de converter qualquer conjunto de símbolos de uma fonte (letras, números, pixels) em um conjunto padronizado de símbolos (tipicamente bits) para transmissão ou armazenamento, e depois reconvertê-los de volta à forma original, sem se preocupar com o que esses símbolos representam em termos de ideias ou conceitos.
- **Exemplo Prático (Transmissão de "Fogo!" vs. "Bom dia"):** Do ponto de vista da teoria da informação de Shannon, a mensagem "FOGO!" e a mensagem "BOM DIA" são ambas sequências de símbolos (letras). A teoria se preocupa com a probabilidade de ocorrência dessas sequências específicas, o número de bits necessários para representá-las de forma única e a melhor maneira de transmiti-las através de um canal, minimizando erros. A urgência, o impacto

emocional ou o significado intrínseco de "FOGO!" em comparação com "BOM DIA" estão fora do escopo da quantificação da informação proposta por Shannon. A informação, neste contexto, está ligada à surpresa ou à redução da incerteza ao receber uma dessas sequências, não ao seu valor semântico.

2.2. A Definição de Informação: Redução da Incerteza

Se a informação não é sobre significado, então o que ela é? Para Shannon, a quantidade de informação produzida quando uma mensagem é escolhida de um conjunto de possíveis mensagens está relacionada à incerteza que é resolvida por essa escolha.³ Se um evento é altamente previsível, sua ocorrência fornece pouca informação nova. Por outro lado, se um evento é muito improvável ou surpreendente, sua ocorrência carrega uma quantidade maior de informação. Assim, informação é fundamentalmente uma medida da redução da incerteza.

- **Analogia 1 (Jogo de Adivinhação "Quente ou Frio"):** Suponha que alguém escolheu secretamente um objeto em uma sala com 100 objetos diferentes. Sua incerteza inicial é alta. Se a pessoa lhe dá uma pista como "está na metade esquerda da sala", essa mensagem reduz sua incerteza pela metade, fornecendo uma quantidade significativa de informação. Se a pista fosse "não é o elefante cor-de-rosa no canto", e só há um elefante cor-de-rosa, a redução da incerteza (e, portanto, a informação) seria menor se você já suspeitasse que não era ele.
- **Analogia 2 (Previsão do Tempo no Deserto vs. Floresta Temperada):** Se um meteorologista anuncia "hoje fará sol e calor" no Deserto do Saara durante o verão, essa mensagem carrega pouquíssima informação, pois é um evento altamente provável e esperado. Sua incerteza sobre o tempo já era baixa. No entanto, se o mesmo meteorologista anunciasse "previsão de neve para hoje" no Saara, essa mensagem seria extremamente surpreendente e, portanto, carregaria uma quantidade imensa de informação, pois reduziria drasticamente uma incerteza que considerava a neve como virtualmente impossível.
- **Exemplo Prático (Resultado de um Lançamento de Moeda Honesta):** Antes de lançar uma moeda honesta, há duas possibilidades igualmente prováveis: cara ou coroa. Você está em um estado de incerteza. Ao observar o resultado (digamos, "cara"), essa incerteza é completamente resolvida. A mensagem "cara" forneceu uma quantidade específica de informação. Se a moeda fosse viciada, com 99% de chance de dar "cara", a mensagem "cara" forneceria muito pouca informação, pois já era quase certa.

2.3. A Escolha da Medida Logarítmica

Shannon, seguindo uma sugestão anterior de Ralph Hartley, adotou uma função

logarítmica como a medida mais natural para a informação.⁴ Existem várias razões para essa escolha:

1. **Aditividade:** Se tivermos dois eventos independentes, a informação total obtida ao observar ambos os eventos deve ser a soma das informações obtidas de cada evento individualmente. As funções logarítmicas têm a propriedade de que $\log(ab) = \log(a) + \log(b)$. Se as probabilidades de eventos independentes se multiplicarem, suas informações (medidas logaritmicamente) se somam.
 2. **Crescimento com o Número de Escolhas:** Se o número de mensagens possíveis em um conjunto aumenta, a quantidade de informação associada à escolha de uma mensagem desse conjunto também deve aumentar. O logaritmo é uma função monotônica crescente.
 3. **Relação com a Codificação:** Como veremos, o logaritmo da base 2 do número de mensagens equiprováveis corresponde diretamente ao número de bits necessários para distinguir essas mensagens. Por exemplo, se temos N mensagens equiprováveis, e $N = 2^k$, então $k = \log_2(N)$ bits são necessários.
 4. **Escalabilidade:** O número de sequências de mensagens possíveis cresce exponencialmente com o comprimento da sequência. O logaritmo transforma esse crescimento exponencial em um crescimento linear, tornando a medida de informação mais gerenciável e intuitiva. Shannon observou que dobrar o tempo de transmissão (ou o comprimento da sequência) aproximadamente eleva ao quadrado o número de mensagens possíveis, o que equivale a dobrar o logaritmo desse número.⁴
- **Analogia 1 (Escala Richter para Terremotos):** A energia liberada por terremotos pode variar por muitas ordens de magnitude. A Escala Richter usa uma base logarítmica para converter essas vastas diferenças em uma escala numérica mais compacta e compreensível, onde cada aumento de unidade representa um aumento de dez vezes na amplitude medida (e cerca de 31 vezes na energia). Da mesma forma, o número de mensagens possíveis pode crescer astronomicamente com o comprimento; o logaritmo "comprime" essa escala para uma medida de informação mais linear.
 - **Analogia 2 (Decibéis para Intensidade Sonora):** A intensidade do som também é medida em uma escala logarítmica (decibéis). Isso ocorre porque nossa percepção da sonoridade é aproximadamente logarítmica em relação à intensidade real da onda sonora. Pequenas mudanças em decibéis podem corresponder a grandes mudanças na potência sonora. A medida logarítmica da informação espelha essa compressão de escala para lidar com grandes números de possibilidades.
 - **Exemplo Prático (Número de Perguntas "Sim/Não"):** Suponha que você

precise adivinhar um número inteiro escolhido aleatoriamente entre 1 e N , onde todas as escolhas são igualmente prováveis. A estratégia ótima é fazer perguntas que dividam o espaço de possibilidades pela metade a cada vez. O número de perguntas "sim/não" necessárias é $\lceil \log_2 N \rceil$.

- Se $N=8$ ($=2^3$), são necessárias $\log_2 8=3$ perguntas. (Ex: "É > 4?", "É > 2 (ou 6)?", etc.)
- Se $N=16$ ($=2^4$), são necessárias $\log_2 16=4$ perguntas.
- Se $N=1000$ (próximo de $2^{10}=1024$), são necessárias cerca de $\log_2 1000 \approx 9.96$, ou seja, 10 perguntas. A quantidade de informação (medida pelo número de perguntas) cresce linearmente (3, 4, 10), enquanto o número de possibilidades cresce exponencialmente (8, 16, 1000). A função logarítmica captura essa relação fundamental.

2.4. O Bit: A Unidade Fundamental da Informação

A escolha da base do logaritmo corresponde à escolha da unidade para medir a informação.⁴ Se a base 2 é usada, a unidade resultante é chamada de **dígito binário**, ou mais brevemente, **bit**. A palavra "bit" foi sugerida a Shannon por John W. Tukey.³ Um bit representa a quantidade de informação necessária para resolver a incerteza entre duas alternativas igualmente prováveis.

- **Analogia 1 (Interruptor de Luz):** Um bit é como um interruptor de luz. Ele pode estar em um de dois estados: ligado (que pode ser representado por 1) ou desligado (representado por 0). Observar o estado do interruptor fornece 1 bit de informação se ambos os estados fossem igualmente prováveis antes da observação.
- **Analogia 2 (Cara ou Coroa):** O resultado de um único lançamento de uma moeda honesta (não viciada) é "cara" ou "coroa". Como há duas possibilidades equiprováveis, saber o resultado fornece exatamente 1 bit de informação.
- **Exemplo Prático (Codificação de Caracteres em Teletipo):** Nos sistemas de teletipo mais antigos, frequentemente se usava um conjunto de 32 caracteres distintos (letras, números, símbolos de pontuação, caracteres de controle). Para representar cada um desses 32 caracteres de forma única usando uma sequência de bits, precisamos de um número k de bits tal que $2^k \geq 32$. O menor k que satisfaz isso é 5, pois $2^5=32$. Portanto, cada símbolo em tal sistema de teletipo representa 5 bits de informação.⁴ Se o sistema transmite n símbolos por segundo, diz-se que o canal tem uma capacidade de $5n$ bits por segundo.⁴ É importante notar que esta é a capacidade máxima; a taxa real de transmissão de informação pode ser menor, dependendo da fonte que alimenta o canal.

Se outras bases logarítmicas são usadas, as unidades de informação recebem nomes

diferentes. Por exemplo, se o logaritmo natural (base e) é usado, a unidade é chamada de "nat" (unidade natural).⁴ A conversão entre unidades é simples: para mudar da base a para a base b , basta multiplicar por $\log_b a$.⁴

3. Entropia de Shannon ($H(X)$): A Medida da Incerteza

Com a informação definida em termos de redução de incerteza e o bit estabelecido como sua unidade fundamental (usando logaritmo de base 2), Shannon introduziu o conceito de **entropia** para quantificar a incerteza média ou a quantidade média de informação associada a uma fonte de informação.¹ A entropia, denotada por $H(X)$ para uma variável aleatória X (que representa a fonte), tornou-se uma pedra angular da teoria da informação.

3.1. Definição Formal e Intuição

A entropia de Shannon de uma variável aleatória discreta X , que pode assumir valores x_1, x_2, \dots, x_n com probabilidades $p(x_1), p(x_2), \dots, p(x_n)$ respectivamente, é uma medida da sua incerteza média. Intuitivamente, ela representa a quantidade média de informação que obtemos ao observar uma realização da variável aleatória X , ou, equivalentemente, o número médio de bits necessários para descrever o resultado de X da maneira mais eficiente possível.

A entropia pode ser vista como o valor esperado da "surpresa" de um resultado.⁷ A "surpresa" de um resultado específico x_i é definida como $\log_2(1/p(x_i))$ (ou $-\log_2 p(x_i)$). Resultados menos prováveis são mais "surpreendentes" e carregam mais informação. A entropia $H(X)$ é então a média dessas surpresas, ponderada pelas probabilidades de cada resultado: $H(X) = E[\log_2(1/P(X))]$.⁷

- **Analogia 1 (O Nível de "Bagunça" numa Caixa de Brinquedos):** Considere uma caixa de brinquedos. Se todos os brinquedos estão completamente misturados e desorganizados, pegar um brinquedo aleatoriamente da caixa envolve alta incerteza – a caixa tem alta entropia. Você não tem muita ideia de qual brinquedo virá. Se, por outro lado, os brinquedos estão meticulosamente organizados em compartimentos (e.g., todos os carrinhos juntos, todas as bonecas juntas), a incerteza ao pegar um brinquedo de um compartimento específico é muito menor – a caixa (ou o compartimento) tem baixa entropia.
- **Analogia 2 (A Previsibilidade de um Texto em um Idioma):** Um texto escrito em português (ou qualquer linguagem natural) possui uma certa estrutura e previsibilidade. Por exemplo, a letra 'q' é quase invariavelmente seguida pela letra 'u'. A probabilidade de outras letras seguirem 'q' é muito baixa. Portanto, a incerteza sobre a próxima letra após um 'q' é pequena. Em contraste, uma

sequência de letras gerada completamente ao acaso, onde cada letra do alfabeto tem a mesma chance de aparecer em qualquer posição, teria uma entropia muito maior, pois não haveria padrões ou previsibilidade.

- **Exemplo Prático (Lançamento de Dados Honestos vs. Viciados):**

- Um dado honesto de seis faces tem seis resultados possíveis, cada um com probabilidade $1/6$. Há uma quantidade significativa de incerteza sobre qual face aparecerá.
- Agora, considere um dado viciado onde a face '6' aparece com probabilidade 0.9 (90%) e as outras cinco faces aparecem com probabilidade 0.02 (2%) cada. A incerteza sobre o resultado deste dado é muito menor. Na maioria das vezes, você esperaria ver um '6'. Portanto, a entropia do dado viciado é consideravelmente menor que a do dado honesto.⁹

A entropia na teoria da informação é análoga ao conceito de entropia na termodinâmica, que mede a quantidade de desordem em sistemas físicos.¹ Embora haja debates sobre a profundidade e exatidão dessa analogia, a intuição de "desordem" ou "incerteza" é útil em ambos os contextos.

3.2. A Fórmula da Entropia: $H(X) = -\sum p(x_i) \log_b p(x_i)$

A fórmula matemática para a entropia de Shannon de uma variável aleatória discreta X que pode assumir os estados x_1, x_2, \dots, x_k com probabilidades $p(x_1), p(x_2), \dots, p(x_k)$ (onde $\sum_{i=1}^k p(x_i) = 1$) é:

$$H(X) = -\sum_{i=1}^k p(x_i) \log_b p(x_i)$$

Vamos detalhar os componentes desta fórmula:

- X : A variável aleatória que representa a fonte de informação ou o conjunto de eventos possíveis.
- x_i : Um dos k resultados ou símbolos possíveis que X pode assumir.
- $p(x_i)$: A probabilidade de ocorrência do resultado x_i .
- $\log_b p(x_i)$: O logaritmo da probabilidade $p(x_i)$. A base b do logaritmo determina a unidade da entropia. Comumente, $b=2$, e a entropia é medida em bits. Se $b=e$ (número de Euler), a unidade é "nats". Se $b=10$, a unidade é "hartleys" ou "dits".
- $-\log_b p(x_i)$: Esta quantidade é chamada de **auto-informação** ou "surpresa" do evento x_i . Como $0 \leq p(x_i) \leq 1$, $\log_b p(x_i)$ será ≤ 0 (para $b > 1$). O sinal negativo torna a auto-informação não negativa. Eventos menos prováveis ($p(x_i)$ pequeno) têm maior auto-informação (maior surpresa). Eventos mais prováveis ($p(x_i)$ grande, próximo de 1) têm menor auto-informação (menor surpresa). Se um evento é certo ($p(x_i)=1$), sua auto-informação é 0.

- $p(x_i)[- \log_b p(x_i)]$: Cada termo da soma é a auto-informação de um evento x_i ponderada pela sua probabilidade de ocorrência $p(x_i)$.
- $\sum_{i=1}^k$: A soma é feita sobre todos os k resultados possíveis da variável aleatória X .

Portanto, a entropia $H(X)$ é o valor esperado (ou média ponderada) da auto-informação dos resultados da variável aleatória X . Por convenção, $0 \log_b 0 = 0$, o que é justificado por continuidade, pois $x \log x \rightarrow 0$ quando $x \rightarrow 0$.¹⁰

Shannon mostrou que qualquer medida de informação (ou incerteza) deveria satisfazer certas propriedades intuitivas, e a fórmula da entropia satisfaz essas propriedades⁵:

1. **Continuidade:** A entropia deve ser uma função contínua das probabilidades $p(x_i)$. Pequenas mudanças nas probabilidades devem levar a pequenas mudanças na entropia.
 2. **Não-negatividade e Certeza:** $H(X) \geq 0$. A entropia é zero se e somente se um dos eventos x_i tem probabilidade 1 (certeza) e todos os outros têm probabilidade 0 (nenhuma incerteza).
 3. **Máxima Incerteza:** Para um número fixo de k resultados possíveis, a entropia é máxima quando todos os resultados são igualmente prováveis, ou seja, $p(x_i) = 1/k$ para todo i . Neste caso, $H(X) = \log_b k$. Esta é a situação de maior incerteza.
 4. **Aditividade para Eventos Independentes (Composição):** Se uma escolha é decomposta em duas escolhas sucessivas e independentes, a entropia da escolha original deve ser a soma das entropias das escolhas componentes. Mais geralmente, para variáveis aleatórias independentes X e Y , $H(X, Y) = H(X) + H(Y)$, onde $H(X, Y)$ é a entropia conjunta (discutida mais adiante).
- **Analogia 1 (Cálculo da Média Ponderada de "Surpresa" em um Sorteio):**
Imagine um sorteio com diferentes prêmios, cada um com uma probabilidade diferente de ser ganho. Cada prêmio tem um "fator de surpresa" associado (inversamente proporcional à sua probabilidade – ganhar um prêmio raro é muito surpreendente). A entropia do sorteio seria a média desses fatores de surpresa, onde cada fator é ponderado pela chance de ganhar aquele prêmio específico. Prêmios muito raros (alta surpresa) contribuem muito para a surpresa individual, mas como são raros, seu peso na média é pequeno. Prêmios comuns (baixa surpresa) contribuem pouco individualmente, mas seu peso na média é maior.
 - **Analogia 2 (Planejando um Código de Comprimento Variável Eficiente):**
Suponha que você queira criar um código (como o código Morse, mas otimizado) para transmitir mensagens compostas por diferentes símbolos. A fórmula da entropia informa o comprimento médio mínimo teórico, em bits por símbolo (se $b=2$), que seu código pode alcançar. Para fazer isso, você atribuiria sequências de

bits mais curtas aos símbolos que aparecem com mais frequência (baixo $p(x_i)$, mas alto $-\log p(x_i)$ pequeno) e sequências de bits mais longas aos símbolos que aparecem raramente (alto $p(x_i)$, mas alto $-\log p(x_i)$ grande). A entropia é a média ponderada desses comprimentos de código ideais.

- **Exemplo Prático (Cálculo da Entropia de uma Moeda Viciada):**

Considere uma moeda viciada onde a probabilidade de "Cara" (C) é $P(C)=0.9$ e a probabilidade de "Coroa" (K) é $P(K)=0.1$. Usando logaritmo de base 2 (para bits):

Auto-informação de Cara: $I(C)=-\log_2(0.9)\approx-(-0.152)=0.152$ bits.

Auto-informação de Coroa: $I(K)=-\log_2(0.1)\approx-(-3.322)=3.322$ bits.

A entropia $H(X)$ é:

$$H(X)=P(C)\cdot I(C)+P(K)\cdot I(K)$$

$$H(X)=(0.9\cdot 0.152)+(0.1\cdot 3.322)$$

$$H(X)=0.1368+0.3322=0.469 \text{ bits.}$$

Compare isso com uma moeda honesta, onde $P(C)=0.5$ e $P(K)=0.5$:

$$I(C)=-\log_2(0.5)=1 \text{ bit.}$$

$$I(K)=-\log_2(0.5)=1 \text{ bit.}$$

$$H(X)=(0.5\cdot 1)+(0.5\cdot 1)=0.5+0.5=1 \text{ bit.}$$

Como esperado, a moeda viciada (resultado mais previsível) tem uma entropia menor (0.469 bits) do que a moeda honesta (resultado menos previsível, entropia máxima para dois resultados, 1 bit). Isso significa que, em média, menos informação é transmitida pelo resultado de um lançamento da moeda viciada.

3.3. Entropia Alta vs. Baixa: Exemplos e Intuição

A magnitude da entropia fornece uma indicação direta do grau de incerteza ou aleatoriedade de uma fonte de informação.

- **Baixa Entropia:** Indica alta previsibilidade, ordem e redundância na fonte. Os resultados são, em certa medida, esperados ou seguem padrões.
 - *Analogia 1 (Um Baralho de Cartas Novo e Ordenado):* Um baralho de cartas recém-aberto, com todas as cartas em sequência (e.g., Ás de Espadas, Dois de Espadas,..., Rei de Copas), tem entropia baixíssima (quase zero, se a ordem específica é conhecida). Se você tirar uma carta do topo, sua identidade é altamente previsível.
 - *Analogia 2 (Uma Melodia Simples e Repetitiva):* Uma canção de ninar com uma melodia muito simples e frases repetitivas tem uma entropia relativamente baixa. As próximas notas são, em grande parte, previsíveis com base nas anteriores.
 - *Exemplo Prático (Texto em Português com Muitas Repetições):* A frase "o rato roeu a roupa do rei de roma o rato roeu a roupa do rei de roma" tem uma

entropia menor do que uma frase mais variada de mesmo comprimento, devido à alta redundância e previsibilidade das palavras.

- **Alta Entropia:** Indica baixa previsibilidade, alta desordem e pouca ou nenhuma redundância. Os resultados são próximos do aleatório.
 - *Analogia 1 (Um Baralho de Cartas Bem Embaralhado):* Após embaralhar completamente um baralho de cartas, a ordem das cartas é altamente imprevisível. Cada carta tem aproximadamente a mesma chance de aparecer em qualquer posição. A entropia do baralho embaralhado é alta.
 - *Analogia 2 (Ruído Branco Estático):* O som de ruído branco (como o chiado de uma rádio fora de sintonia) é caracterizado por ter todas as frequências presentes com igual intensidade, resultando em um sinal altamente aleatório e imprevisível. Sua entropia é alta. Mudanças abruptas no registro, textura e tonalidade em uma peça musical também são exemplos de contextos de alta entropia local.¹¹
 - *Exemplo Prático (Uma Senha Aleatória Forte):* Uma senha como "TrBw3!z@9*qX" gerada aleatoriamente com caracteres maiúsculos, minúsculos, números e símbolos tem alta entropia. Cada caractere é difícil de prever com base nos anteriores. Em contraste, uma senha como "123456" tem baixíssima entropia.
 - *Exemplo Prático (A "Neve" na Tela de uma TV Analógica sem Sinal):* A imagem de "neve" (pontos pretos e brancos piscando aleatoriamente) em uma televisão antiga que não está sintonizada em nenhum canal é um excelente exemplo visual de um processo de alta entropia. Cada pixel na tela parece ser determinado aleatoriamente e independentemente dos pixels vizinhos, levando a uma imagem sem padrões discerníveis e, portanto, altamente imprevisível.¹²

A entropia de uma fonte de informação tem implicações diretas na sua compressibilidade. Fontes de baixa entropia são inerentemente mais redundantes.¹⁴ Essa redundância pode ser explorada por algoritmos de compressão para representar a mesma informação usando menos bits. Por exemplo, em um texto com muitas palavras repetidas, podemos substituir as repetições por referências mais curtas. Fontes de alta entropia, por outro lado, já são "densas" em informação; cada símbolo contribui significativamente para a novidade, oferecendo poucas oportunidades para compressão sem perdas. Tentar comprimir dados de alta entropia (como um arquivo já bem comprimido ou dados verdadeiramente aleatórios) geralmente resulta em pouco ou nenhum ganho de tamanho, e às vezes pode até aumentar ligeiramente o tamanho devido ao overhead do algoritmo de compressão.

3.4. Unidades de Entropia: Bits e Nats

Como mencionado anteriormente, a unidade da entropia depende da base do logaritmo (b) usada na fórmula $H(X) = -\sum p(x_i) \log_b p(x_i)$.

- Se o logaritmo de **base 2** é usado, a entropia é medida em **bits**.⁴ Esta é a unidade mais comum na ciência da computação e na teoria da informação digital, pois os computadores operam fundamentalmente com dígitos binários.
- Se o logaritmo natural (logaritmo neperiano, **base e**) é usado, a entropia é medida em **nats** (unidades naturais).⁴ Nats são frequentemente usados em contextos matemáticos e estatísticos mais teóricos, e em algumas áreas de aprendizado de máquina.
- Se o logaritmo de **base 10** é usado, a unidade é às vezes chamada de **hartley** (em homenagem a Ralph Hartley) ou **dit** (decimal digit).

A escolha da base é, em essência, uma escolha de unidade de medida. A quantidade de incerteza subjacente é a mesma, apenas a escala da medição muda. A conversão entre entropias calculadas com bases diferentes é direta, usando a propriedade de mudança de base dos logaritmos: $H_b(X) = (\log_b a) \cdot H_a(X)$. Por exemplo, para converter de nats para bits, multiplicamos por $\log_2 e \approx 1.4427$. Para converter de bits para nats, multiplicamos por $\ln 2 \approx 0.6931$.

- **Analogia 1 (Medindo Distância em Diferentes Unidades):** Pense em medir a distância entre duas cidades. Você pode expressar essa distância em quilômetros ou em milhas. O valor numérico será diferente (e.g., 100 km vs. 62.14 milhas), mas a distância física real é a mesma. Bits e nats são como quilômetros e milhas para a informação: unidades diferentes para medir a mesma quantidade fundamental de incerteza.
- **Analogia 2 (Moedas de Diferentes Países com Taxas de Câmbio):** Bits e nats podem ser vistos como duas moedas diferentes, digamos, o Dólar Americano e o Euro. Ambas representam valor (informação/incerteza), mas têm "taxas de câmbio" diferentes entre si, dadas pelo fator de conversão logarítmico. Você pode converter uma quantia de uma moeda para outra, mas o poder de compra subjacente (a quantidade de informação) permanece o mesmo.
- **Exemplo Prático (Cálculo da Entropia de um Dado Honesto de 6 Faces):**
Um dado honesto tem 6 resultados equiprováveis, então $p(x_i) = 1/6$ para $i = 1, \dots, 6$.
A entropia em bits (usando \log_2) é:
 $H_2(X) = -\sum_{i=1}^6 \frac{1}{6} \log_2 \left(\frac{1}{6}\right) = -6 \cdot \left(\frac{1}{6} \log_2 \left(\frac{1}{6}\right)\right) = -\log_2 \left(\frac{1}{6}\right) = \log_2(6)$
 $H_2(X) \approx 2.585$ bits.¹⁵
Isso significa que, em média, são necessários cerca de 2.585 bits para especificar o resultado de um lançamento de um dado honesto de seis faces.

A entropia em nats (usando \ln , logaritmo natural) é:

$$H_e(X) = -\sum_i p_i \ln(p_i) = -\ln(61) = \ln(6)$$

$$H_e(X) \approx 1.792 \text{ nats.}$$

Podemos verificar a conversão: $H_e(X) \cdot \log_2 e \approx 1.792 \cdot 1.4427 \approx 2.585$, que é $H_2(X)$.

4. Entropia e Compressão de Dados: O Limite Teórico

Um dos resultados mais profundos e práticos da teoria da informação de Shannon é a conexão direta entre a entropia de uma fonte de informação e o limite fundamental para a compressão de dados sem perdas (lossless data compression). A compressão de dados visa representar a informação de forma mais compacta, usando menos bits, sem perder nenhum detalhe original.

4.1. O Teorema da Codificação de Fonte de Shannon (Teorema Silencioso / Noiseless Coding Theorem)

O Teorema da Codificação de Fonte de Shannon, também conhecido como teorema da codificação sem ruído (noiseless coding theorem), estabelece os limites estatísticos para a compressão de dados.¹⁶ Ele lida com a eficiência com que os dados gerados por uma fonte podem ser representados.

O teorema afirma fundamentalmente o seguinte para uma fonte que gera símbolos independentes e identicamente distribuídos (i.i.d.):

1. É impossível comprimir os dados de tal forma que a taxa de código média (o número médio de bits usados por símbolo da fonte original) seja menor que a entropia da fonte $H(X)$ sem que haja perda de informação. Ou seja, $H(X)$ é um limite inferior para o comprimento médio de qualquer código de prefixo que represente a fonte X .⁸
2. Por outro lado, é possível construir um esquema de codificação tal que a taxa de código média seja arbitrariamente próxima de $H(X)$ (especificamente, para qualquer $\epsilon > 0$, pode-se atingir uma taxa de $H(X) + \epsilon$ com uma probabilidade de erro de decodificação que pode ser tornada arbitrariamente pequena, à medida que o comprimento dos blocos de símbolos codificados (n) tende ao infinito).⁸

Em essência, a entropia $H(X)$ de uma fonte não é apenas uma medida de incerteza, mas também define o limite teórico fundamental para a compressão de dados sem perdas dessa fonte.⁸ Não se pode comprimir os dados, em média, para usar menos bits por símbolo do que a entropia da fonte, se se deseja ser capaz de reconstruir os dados originais perfeitamente.

- **Analogia 1 (Mala de Viagem Eficiente):** Imagine que você está fazendo uma

mala para uma viagem. A "entropia" de suas roupas e pertences (determinada pela variedade, necessidade de cada item, quão compactáveis são, etc.) dita o tamanho mínimo da mala que você precisa para levar tudo essencial sem amassar ou danificar nada (perder informação). O Teorema da Codificação de Fonte diz que existe um "tamanho de mala ideal" (a entropia). Você pode tentar usar uma mala menor, mas inevitavelmente terá que deixar coisas importantes para trás ou amassar tudo de forma irreversível. Você pode usar uma mala um pouco maior que o ideal e organizar tudo perfeitamente para caber.

- **Analogia 2 (Resumindo um Livro Detalhado):** Considere um livro longo e detalhado. A entropia do livro representa sua "essência informativa" – a quantidade mínima de informação que captura todas as suas ideias cruciais e enredo. Você pode resumir o livro (comprimi-lo), tornando-o mais curto. No entanto, há um ponto a partir do qual qualquer resumo adicional levaria à perda de informações vitais, alterando o significado ou omitindo partes importantes da história. A entropia define o limite desse "resumo perfeito" sem perdas.
- **Exemplo Prático (Compressão de Arquivos de Texto com ZIP):**
 - Um arquivo de texto que contém muitas palavras repetidas, frases comuns ou longas sequências do mesmo caractere (por exemplo, um arquivo de log com muitas entradas "SUCESSO SUCESSO SUCESSO") tem baixa entropia. Algoritmos de compressão como o ZIP podem explorar essa redundância e reduzir significativamente o tamanho do arquivo.
 - Por outro lado, um arquivo de texto que já foi comprimido (e, portanto, tem a maior parte de sua redundância removida) ou um arquivo contendo dados verdadeiramente aleatórios (como a saída de um gerador de números aleatórios criptograficamente seguro) tem alta entropia. Tentar comprimir tal arquivo com ZIP resultará em pouca ou nenhuma redução de tamanho, e às vezes o arquivo pode até ficar ligeiramente maior devido ao cabeçalho e metadados adicionados pelo formato ZIP.¹⁴ O Teorema da Codificação de Fonte de Shannon implica que a taxa de compressão alcançada pelo arquivo.zip (número de bits no arquivo comprimido dividido pelo número de símbolos no arquivo original) não pode ser, em média, menor que a entropia do conteúdo original do arquivo.

Este teorema não apenas estabelece um limite, mas também destaca que as características estatísticas da fonte – as probabilidades de seus símbolos – são cruciais para a compressão. Fontes cujos símbolos não são uniformemente distribuídos (ou seja, alguns símbolos são muito mais prováveis que outros) geralmente têm entropia mais baixa e, portanto, são mais compressíveis.⁸ Se todos os símbolos fossem equiprováveis, a entropia seria máxima para aquele número de

símbolos, indicando menos redundância inerente para ser explorada pela compressão. A genialidade do teorema é formalizar essa intuição: códigos mais curtos podem ser atribuídos a símbolos mais frequentes e códigos mais longos a símbolos mais raros para alcançar uma representação média mais compacta.

4.2. Codificação de Huffman: Um Exemplo Prático de Compressão Eficiente

Enquanto o Teorema da Codificação de Fonte de Shannon prova a existência de códigos ótimos, ele não especifica como construí-los. A **Codificação de Huffman** é um exemplo prático e amplamente utilizado de um algoritmo que constrói códigos de prefixo de comprimento variável quase ótimos para compressão de dados sem perdas.¹⁸

O princípio fundamental da Codificação de Huffman é atribuir códigos binários mais curtos aos símbolos que ocorrem com mais frequência na fonte de dados e códigos binários mais longos aos símbolos que ocorrem com menos frequência.¹⁹ Uma propriedade crucial dos códigos de Huffman é que eles são **códigos de prefixo** (também conhecidos como códigos instantâneos). Isso significa que nenhum código atribuído a um símbolo é o prefixo do código atribuído a qualquer outro símbolo.¹⁸ Essa propriedade é essencial porque garante que o fluxo de bits codificado possa ser decodificado de forma única e não ambígua. Por exemplo, se 'A' fosse codificado como "0" e 'B' como "01", ao receber "01", o decodificador não saberia se é 'A' seguido por algo ou se é 'B'. Com códigos de prefixo, assim que uma sequência de bits corresponde a um código válido, ela pode ser decodificada imediatamente.

O algoritmo de Huffman constrói uma árvore binária (a árvore de Huffman) de forma bottom-up (de baixo para cima), geralmente utilizando uma fila de prioridade (min-heap) para gerenciar os nós da árvore com base em suas frequências ¹⁸:

1. **Inicialização:** Crie um nó folha para cada símbolo único da fonte. Cada nó folha armazena o símbolo e sua frequência (ou probabilidade) de ocorrência. Insira todos os nós folha na fila de prioridade, onde a prioridade é determinada pela frequência (menor frequência tem maior prioridade para ser extraída).
2. **Construção da Árvore:** Enquanto houver mais de um nó na fila de prioridade: a. Extraia os dois nós com as menores frequências da fila. Estes serão os nós com menor probabilidade de ocorrência até o momento. b. Crie um novo nó interno. A frequência deste novo nó será a soma das frequências dos dois nós extraídos. c. Faça os dois nós extraídos serem os filhos (esquerdo e direito) do novo nó interno. Não importa qual é o esquerdo ou o direito inicialmente, mas uma convenção consistente deve ser usada (e.g., o de menor frequência à esquerda). d. Insira o novo nó interno de volta na fila de prioridade.

3. **Árvore Completa:** Quando restar apenas um nó na fila de prioridade, este nó é a raiz da árvore de Huffman, e a construção está completa.
 4. **Atribuição de Códigos:** Para obter o código binário para cada símbolo, percorra a árvore da raiz até o nó folha correspondente ao símbolo. Atribua '0' para cada vez que você seguir para um filho esquerdo e '1' para cada vez que seguir para um filho direito (ou vice-versa, desde que a convenção seja consistente). A sequência de 0s e 1s acumulada ao longo do caminho é o código de Huffman para aquele símbolo.
- Exemplo Prático (Codificação de uma pequena string de caracteres):
Vamos usar o exemplo fornecido em 18 e 19, com os seguintes caracteres e suas frequências:
 - a: 5
 - b: 9
 - c: 12
 - d: 13
 - e: 16
 - f: 45

Construção da Árvore de Huffman (passo a passo, como em ¹⁸):

1. **Inicial:** Nós folha na fila de prioridade (símbolo: frequência): (a:5), (b:9), (c:12), (d:13), (e:16), (f:45).
2. **Extraí (a:5) e (b:9).** Cria nó interno N1 com frequência 5+9=14. Filhos: a, b. Fila: (c:12), (d:13), (N1:14), (e:16), (f:45).
3. **Extraí (c:12) e (d:13).** Cria nó interno N2 com frequência 12+13=25. Filhos: c, d. Fila: (N1:14), (e:16), (N2:25), (f:45).
4. **Extraí (N1:14) e (e:16).** Cria nó interno N3 com frequência 14+16=30. Filhos: N1, e. Fila: (N2:25), (N3:30), (f:45).
5. **Extraí (N2:25) e (N3:30).** Cria nó interno N4 com frequência 25+30=55. Filhos: N2, N3. Fila: (f:45), (N4:55).
6. **Extraí (f:45) e (N4:55).** Cria nó raiz N5 com frequência 45+55=100. Filhos: f, N4. Fila: (N5:100). Fim.

Atribuição de Códigos (convenção: 0 para esquerda, 1 para direita, e assumindo que o filho de menor frequência foi para a esquerda nas combinações): A árvore resultante (simplificada) seria algo como: N5(100)

```

      /   \
    f(45)  N4(55)
   (0)    /   \ (1)
        N2(25) N3(30)
       / \ (0) / \ (1)
      c(12) d(13) N1(14) e(16)
     (0) (1) / \ (0) (1)

```

a(5) b(9)
(0) (1)

Percorrendo da raiz aos folhas:

- f: 0
- c: 100
- d: 101
- a: 1100
- b: 1101
- e: 111

Comprimento Médio do Código: Soma das (frequência do símbolo * comprimento do seu código) / soma total das frequências. Total de frequências = 5+9+12+13+16+45 = 100. Comprimento médio =

$(5 \cdot 4 + 9 \cdot 4 + 12 \cdot 3 + 13 \cdot 3 + 16 \cdot 3 + 45 \cdot 1) / 100 = (20 + 36 + 36 + 39 + 48 + 45) / 100 = 224 / 100 = 2.24$

bits/símbolo. Se usássemos um código de comprimento fixo para 6 símbolos, precisaríamos de $\lceil \log_2 6 \rceil = 3$ bits por símbolo. A codificação de Huffman economiza, em média, $3 - 2.24 = 0.76$ bits por símbolo para esta distribuição de frequências. A entropia desta fonte

é: $H(X) = -[(5/100)\log_2(5/100) + (9/100)\log_2(9/100) + (12/100)\log_2(12/100) + (13/100)\log_2(13/100) + (16/100)\log_2(16/100) + (45/100)\log_2(45/100)]$
 $H(X) \approx -[0.05(-4.32) + 0.09(-3.47) + 0.12(-3.06) + 0.13(-2.94) + 0.16(-2.64) + 0.45(-1.15)]$
 $H(X) \approx -[-0.216 - 0.312 - 0.367 - 0.382 - 0.422 - 0.518] = -[-2.217] = 2.217$ bits/símbolo. O código de Huffman (2.24 bits/símbolo) chega muito perto do limite da entropia (2.217 bits/símbolo).

A codificação de Huffman é um algoritmo *guloso* (greedy), pois a cada passo toma a decisão localmente ótima de combinar os dois nós de menor frequência.¹⁹ Para o problema de encontrar um código de prefixo ótimo para um conjunto conhecido de probabilidades de símbolos (codificação símbolo a símbolo), a codificação de Huffman de fato atinge o comprimento médio de código esperado mais curto possível, aproximando-se do limite da entropia de Shannon. A estratégia gulosa funciona porque, intuitivamente, "empurra" os símbolos menos prováveis para níveis mais profundos na árvore (resultando em códigos mais longos para eles) e, consequentemente, permite que os símbolos mais prováveis fiquem mais próximos da raiz (resultando em códigos mais curtos). Embora seja uma heurística, ela leva a uma solução ótima para este tipo específico de problema de codificação.

5. Explorando Múltiplas Variáveis: Entropias Avançadas

A entropia de Shannon, $H(X)$, quantifica a incerteza de uma única variável aleatória. No entanto, muitos sistemas e problemas do mundo real envolvem múltiplas variáveis que podem interagir e depender umas das outras. Para analisar tais sistemas, a teoria da informação estende o conceito de entropia para lidar com múltiplas variáveis. Introduzimos aqui a entropia conjunta, a entropia condicional e a regra da cadeia que

as relaciona.

Antes de detalhar cada uma, a tabela a seguir resume essas medidas, que serão cruciais para entender as seções subsequentes, incluindo a Informação Mútua.

Medida	O Que Mede	Fórmula Principal (base 2)	Exemplo Intuitivo
Entropia (Shannon) $H(X)$	Incerteza de uma variável aleatória X .	$-\sum p(x)\log p(x)$	Incerteza no resultado do lançamento de um dado.
Entropia Conjunta $H(X,Y)$	Incerteza combinada de duas variáveis aleatórias X e Y ocorrendo juntas.	$-\sum \sum p(x,y)\log p(x,y)$	Incerteza sobre o tempo E o trânsito amanhã.
Entropia Condicional $H(Y X)$	X)	Incerteza da variável Y , dado que o valor da variável X já é conhecido.	$-\sum \sum p(x,y)\log p(y x)$
Informação Mútua $I(X;Y)$	Quantidade de informação que X contém sobre Y (ou Y sobre X). Redução da incerteza.	$H(X) - H(X Y)$	Y) ou $H(Y) - H(Y X)$

Esta tabela serve como um guia rápido. À medida que introduzimos múltiplas variáveis, os conceitos de entropia se ramificam. Uma visão comparativa é vital para distinguir claramente entre $H(X)$, $H(X,Y)$, $H(Y|X)$ e, posteriormente, $I(X;Y)$. Ela fornece um resumo conciso de suas definições, o que medem, suas fórmulas e um exemplo simples para ancorar a intuição, prevenindo confusão e facilitando a absorção do material subsequente.

5.1. Entropia Conjunta ($H(X,Y)$): Incerteza de um Sistema de Variáveis

A **entropia conjunta** $H(X,Y)$ de um par de variáveis aleatórias discretas (X,Y) mede a incerteza total associada à ocorrência conjunta de valores específicos para X e Y .¹⁰ Se $p(x,y)$ é a probabilidade conjunta de X assumir o valor x e Y assumir o valor y , a fórmula para a entropia conjunta (usando logaritmo de base 2) é:

$$H(X,Y) = -\sum_{x \in X} \sum_{y \in Y} p(x,y) \log_2 p(x,y)$$

Intuitivamente, $H(X,Y)$ quantifica a quantidade média de informação necessária para especificar os valores de *ambas* as variáveis simultaneamente.

A entropia conjunta possui algumas propriedades importantes ²⁰:

1. **Maior ou igual às entropias marginais:** $H(X,Y) \geq H(X)$ e $H(X,Y) \geq H(Y)$. A incerteza sobre o par de variáveis é sempre pelo menos tão grande quanto a incerteza sobre qualquer uma das variáveis individualmente. Adicionar um novo sistema (variável) nunca pode reduzir a incerteza disponível.
 2. **Subaditividade:** $H(X,Y) \leq H(X) + H(Y)$. A incerteza conjunta de duas variáveis nunca é maior que a soma de suas incertezas individuais. A igualdade $H(X,Y) = H(X) + H(Y)$ ocorre se, e somente se, as variáveis X e Y são estatisticamente independentes. Se elas são dependentes, conhecer uma fornece alguma informação sobre a outra, então a incerteza conjunta é menor do que a soma das incertezas individuais.
- **Analogia 1 (Prever o Resultado de Duas Moedas Lançadas Simultaneamente):** Se você lança duas moedas, X e Y . A entropia conjunta $H(X,Y)$ refere-se à sua incerteza sobre qual dos quatro resultados possíveis ocorrerá: (Cara, Cara), (Cara, Coroa), (Coroa, Cara), (Coroa, Coroa).
 - Se as moedas são independentes e honestas, $p(x,y) = 1/4$ para cada par. $H(X,Y) = \log_2 4 = 2$ bits. Também, $H(X) = 1$ bit e $H(Y) = 1$ bit, então $H(X,Y) = H(X) + H(Y)$.
 - Se as moedas são perfeitamente correlacionadas (e.g., uma é a cópia da outra, sempre caem iguais), então saber X determina Y . Os únicos resultados possíveis são (Cara, Cara) e (Coroa, Coroa), cada um com $p = 1/2$. Neste caso, $H(X,Y) = \log_2 2 = 1$ bit, que é igual a $H(X)$ (e $H(Y)$). Aqui, $H(X,Y) < H(X) + H(Y)$.
 - **Analogia 2 (Diagnóstico Médico Baseado em Múltiplos Sintomas Iniciais):** Um médico está avaliando um paciente que apresenta dois sintomas iniciais: febre (variável X : tem/não tem febre) e dor de garganta (variável Y : tem/não tem dor de garganta). A entropia conjunta $H(\text{Febre}, \text{DorGarganta})$ representa a incerteza total do médico sobre a combinação específica desses dois sintomas que o paciente apresenta, com base na prevalência dessas combinações na população de pacientes que ele vê.
 - **Exemplo Prático (Relação entre Condições Climáticas e Uso de Guarda-Chuva):** Seja X uma variável aleatória que representa se está chovendo ou não (eventos $x_1 = \text{chove}$, $x_2 = \text{não chove}$) e Y uma variável aleatória que representa se uma pessoa está usando um guarda-chuva ou não (eventos $y_1 = \text{usa}$, $y_2 = \text{não usa}$). A entropia conjunta $H(X,Y)$ mede a incerteza sobre a ocorrência conjunta desses eventos. Por exemplo, a probabilidade $p(x_1, y_1)$ (chove

E usa guarda-chuva) é provavelmente alta, enquanto $p(x_2, y_1)$ (não chove E usa guarda-chuva) é provavelmente baixa. A entropia conjunta considera as probabilidades de todas as quatro combinações possíveis para quantificar a incerteza total sobre o "estado do sistema" (clima e guarda-chuva).²³

A diferença $H(X) + H(Y) - H(X, Y)$ é uma medida de quanta dependência existe entre as variáveis X e Y . Se elas são independentes, $H(X, Y) = H(X) + H(Y)$, e essa diferença é zero. Se elas são dependentes, $H(X, Y) < H(X) + H(Y)$, e a diferença é positiva. Essa diferença, como veremos na Seção 6, é precisamente a **Informação Mútua** $I(X; Y)$. A "lacuna" entre a incerteza conjunta e a soma das incertezas individuais surge porque, se as variáveis são dependentes, conhecer uma delas pode fornecer alguma informação sobre a outra, reduzindo assim a incerteza total necessária para descrever ambas em comparação com descrevê-las separadamente e depois somar as incertezas. Essa "redução de incerteza devido à dependência" é a essência da informação compartilhada entre elas.

5.2. Entropia Condicional ($H(Y|X)$): Incerteza Restante Após Conhecimento

A **entropia condicional** $H(Y|X)$ mede a incerteza média que resta sobre a variável aleatória Y quando o valor da variável aleatória X já é conhecido.¹⁰ Ela responde à pergunta: "Quanto ainda não sei sobre Y , mesmo depois que você me contou o valor de X ?"

A fórmula para a entropia condicional é derivada como a esperança da entropia de Y condicionada a valores específicos de X :

$$H(Y|X) = \sum_{x \in X} p(x) H(Y|X=x)$$

onde $H(Y|X=x) = -\sum_{y \in Y} p(y|x) \log_2 p(y|x)$ é a entropia de Y quando X tem o valor específico x , e $p(y|x)$ é a probabilidade condicional de $Y=y$ dado $X=x$.

Substituindo, obtemos a forma mais comum:

$$H(Y|X) = -\sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(y|x)$$

A entropia condicional $H(Y|X)$ é, por vezes, chamada de "entropia de ruído" de Y com respeito a X , porque representa a variabilidade ou incerteza em Y que não pode ser explicada ou prevista pelo conhecimento de X .²¹ É a parte da entropia de Y que não é informativa sobre X .

Propriedades da entropia condicional⁷:

1. **Não-negatividade:** $H(Y|X) \geq 0$.
2. **Conhecimento não aumenta incerteza:** $H(Y|X) \leq H(Y)$. Saber o valor de X nunca aumenta a incerteza média sobre Y . A incerteza pode diminuir ou permanecer a mesma.

3. **Independência:** $H(Y|X)=H(Y)$ se, e somente se, X e Y são variáveis aleatórias independentes. Se são independentes, conhecer X não fornece nenhuma informação sobre Y, então a incerteza sobre Y permanece inalterada.
 4. **Determinação:** Se Y é completamente determinada por X (ou seja, Y é uma função de X), então $H(Y|X)=0$. Uma vez que X é conhecido, não há mais incerteza sobre Y.
- **Analogia 1 (Adivinhar a Segunda Carta de um Baralho Após Ver a Primeira):** Suponha que você retire duas cartas de um baralho padrão de 52 cartas, sem reposição. Seja X a primeira carta retirada e Y a segunda. A entropia $H(Y)$ representa sua incerteza inicial sobre a segunda carta (antes de qualquer carta ser retirada). A entropia condicional $H(Y|X)$ representa sua incerteza sobre a segunda carta DEPOIS que você já viu qual foi a primeira carta. Como a primeira carta não é repostada, conhecer X muda as probabilidades para Y. Por exemplo, se X foi um Ás de Espadas, então Y não pode ser um Ás de Espadas, e as probabilidades das outras 51 cartas aumentam ligeiramente. Assim, $H(Y|X)$ será menor que $H(Y)$, pois conhecer X reduziu sua incerteza sobre Y.
 - **Analogia 2 (Incerteza do Diagnóstico Após um Resultado de Exame):** Seja D uma variável que indica se um paciente tem uma determinada doença (presente/ausente) e T uma variável que indica o resultado de um teste diagnóstico para essa doença (positivo/negativo). $H(D)$ é a incerteza inicial sobre o paciente ter a doença (baseada, por exemplo, na prevalência da doença). $H(D|T=\text{positivo})$ é a incerteza que resta sobre o paciente ter a doença APÓS a observação de um resultado positivo no teste. Um bom teste diagnóstico deve reduzir significativamente essa incerteza, ou seja, $H(D|T=\text{positivo})$ deve ser muito menor que $H(D)$.
 - **Exemplo Prático (Previsão do Próximo Caractere em um Texto Dado o Caractere Anterior):** Em um texto em português, seja X o caractere atual e Y o próximo caractere.
 - Se o caractere atual X é 'q', a incerteza sobre qual será o próximo caractere Y, $H(Y|X='q')$, é muito baixa, pois Y é quase certamente 'u'.
 - Se o caractere atual X é 'a', a incerteza $H(Y|X='a')$ é consideravelmente maior, pois muitas letras diferentes podem seguir um 'a' (e.g., 'b', 'c', 'd', 'l', 'm', 'n', 'r', 's', 't', 'v', vogais, espaço). A entropia condicional $H(Y|X)$ seria a média dessas incertezas $H(Y|X=x)$ sobre todas as possíveis letras x do alfabeto, ponderadas por suas probabilidades de ocorrência $p(x)$.

A entropia condicional é um conceito chave para entender a dependência entre variáveis. Se $H(Y|X) < H(Y)$, isso implica que X fornece alguma informação sobre Y, pois conhecer X reduziu a incerteza sobre Y. A quantidade exata dessa redução,

$H(Y) - H(Y|X)$, é, como veremos, a Informação Mútua. Se conhecer X não altera nossa incerteza sobre Y (ou seja, $H(Y|X) = H(Y)$), isso significa que X e Y são independentes. A entropia condicional quantifica a incerteza que *permanece*; a diferença entre a entropia original e a condicional representa a informação que foi *compartilhada* ou *transmitida* por X sobre Y .

5.3. A Regra da Cadeia para Entropia: Relacionando as Entropias

A **Regra da Cadeia para Entropia** (Chain Rule for Entropy) estabelece uma relação fundamental e elegante entre a entropia conjunta $H(X,Y)$, a entropia marginal $H(X)$ (ou $H(Y)$), e a entropia condicional $H(Y|X)$ (ou $H(X|Y)$).¹⁰ A regra afirma:

$$H(X,Y) = H(X) + H(Y|X)$$

E, por simetria (já que $H(X,Y) = H(Y,X)$):

$$H(X,Y) = H(Y) + H(X|Y)$$

Interpretação: A incerteza total sobre o par de variáveis (X,Y) é igual à incerteza sobre X mais a incerteza que resta sobre Y uma vez que X é conhecido. Alternativamente, é a incerteza sobre Y mais a incerteza que resta sobre X uma vez que Y é conhecido.

A prova desta regra decorre diretamente das definições de entropia conjunta e condicional e da regra de probabilidade $p(x,y) = p(x)p(y|x)$:

$$H(X,Y) = -\sum_x \sum_y p(x,y) \log p(x,y)$$

$$= -\sum_x \sum_y p(x,y) \log [p(x)p(y|x)]$$

$$= -\sum_x \sum_y p(x,y) [\log p(x) + \log p(y|x)]$$

$$= -\sum_x \sum_y p(x,y) \log p(x) - \sum_x \sum_y p(x,y) \log p(y|x)$$

O primeiro termo é $-\sum_x p(x) \log p(x) = H(X)$ (já que $\sum_y p(x,y) = p(x)$).

O segundo termo é, por definição, $H(Y|X)$.

Portanto, $H(X,Y) = H(X) + H(Y|X)$.

A Regra da Cadeia pode ser generalizada para um conjunto de n variáveis aleatórias

X_1, X_2, \dots, X_n 22:

$$H(X_1, X_2, \dots, X_n) = H(X_1) + H(X_2|X_1) + H(X_3|X_1, X_2) + \dots + H(X_n|X_1, \dots, X_{n-1})$$

Ou, de forma mais compacta:

$$H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i|X_1, \dots, X_{i-1})$$

(onde $H(X_1|X_0)$ é definido como $H(X_1)$).

- **Analogia 1 (Construindo um Quebra-Cabeça de Duas Peças):** A incerteza total sobre a imagem final formada por duas peças de um quebra-cabeça ($H(X,Y)$, onde X é a primeira peça e Y é a segunda) é igual à sua incerteza sobre qual é a primeira peça que você pega ($H(X)$), mais a sua incerteza sobre qual será a

segunda peça, dado que você já viu e posicionou a primeira ($H(Y|X)$).

- **Analogia 2 (Desvendando um Mistério de Crime):** Em um jogo de detetive, a incerteza total sobre "quem cometeu o crime E com qual arma" ($H(\text{Criminoso}, \text{Arma})$) pode ser decomposta. É igual à sua incerteza inicial sobre quem é o criminoso ($H(\text{Criminoso})$), mais a incerteza que resta sobre qual arma foi usada, uma vez que você já descobriu a identidade do criminoso ($H(\text{Arma}|\text{Criminoso})$).
- **Exemplo Prático (Previsão do Tempo para Hoje e Amanhã):** Seja X o estado do tempo hoje (e.g., ensolarado, chuvoso, nublado) e Y o estado do tempo amanhã. A incerteza conjunta sobre o tempo nos dois dias, $H(X,Y)$, é igual à incerteza sobre o tempo hoje, $H(X)$, mais a incerteza sobre o tempo de amanhã, dado que já sabemos como está o tempo hoje, $H(Y|X)$. Como o tempo de amanhã geralmente tem alguma dependência do tempo de hoje (e.g., uma frente fria hoje aumenta a chance de chuva amanhã), então $H(Y|X)$ será geralmente menor que $H(Y)$ (a incerteza sobre o tempo de amanhã sem saber o de hoje).

A Regra da Cadeia é mais do que uma simples identidade matemática; ela reflete a natureza sequencial pela qual a informação pode ser revelada ou acumulada em muitos sistemas. Ela decompõe a incerteza conjunta em componentes condicionais, o que é fundamental para modelar processos estocásticos (como cadeias de Markov, onde o próximo estado depende do estado atual) e sistemas onde a informação é adquirida passo a passo (como em um processo diagnóstico ou na compressão de dados sequenciais onde o próximo símbolo é codificado com base nos anteriores). A regra da cadeia fornece a base para quantificar a informação em tais processos dinâmicos e é uma ferramenta essencial na teoria da informação e suas aplicações.

6. Informação Mútua ($I(X;Y)$): Medindo a Dependência e o Compartilhamento de Informação

Enquanto a entropia conjunta $H(X,Y)$ mede a incerteza total de um par de variáveis e a entropia condicional $H(Y|X)$ mede a incerteza de Y que permanece após X ser conhecido, a **Informação Mútua** $I(X;Y)$ quantifica a quantidade de informação que uma variável aleatória X contém sobre outra variável aleatória Y, e vice-versa. É uma medida da dependência estatística entre as duas variáveis.

6.1. Definição e Fórmulas

A Informação Mútua $I(X;Y)$ entre duas variáveis aleatórias discretas X e Y pode ser definida de várias maneiras equivalentes, cada uma oferecendo uma perspectiva ligeiramente diferente sobre seu significado ²¹:

1. **Em termos de redução da entropia (incerteza):**

- $I(X;Y)=H(X)-H(X|Y)$ ²¹ Interpretação: A informação mútua é a redução na incerteza sobre X que ocorre quando Y se torna conhecido.
 - $I(X;Y)=H(Y)-H(Y|X)$ ²¹ Interpretação: A informação mútua é a redução na incerteza sobre Y que ocorre quando X se torna conhecido. (Como $I(X;Y)$ é simétrica, ambas as formas são iguais).
2. **Em termos de entropias marginais e conjunta:**
- $I(X;Y)=H(X)+H(Y)-H(X,Y)$ ²¹ Interpretação: A informação mútua é a soma das entropias individuais menos a entropia conjunta. Se X e Y fossem independentes, $H(X,Y)=H(X)+H(Y)$, então $I(X;Y)=0$. Qualquer valor de $I(X;Y)>0$ indica que $H(X,Y)<H(X)+H(Y)$, ou seja, há alguma redundância ou informação compartilhada.
3. **Em termos da Divergência de Kullback-Leibler (discutida na Seção 7):**
- $I(X;Y)=DKL(p(x,y) || p(x)p(y))$
 - $I(X;Y)=\sum_{x \in X} \sum_{y \in Y} p(x,y) \log_2 \frac{p(x,y)}{p(x)p(y)}$ ²⁷ Interpretação: A informação mútua mede quão diferente a distribuição conjunta real $p(x,y)$ é do produto das distribuições marginais $p(x)p(y)$ (que seria a distribuição conjunta se X e Y fossem independentes). É a "distância" informacional entre a verdadeira relação das variáveis e a hipótese de independência.

Propriedades da Informação Mútua:

- **Não-negatividade:** $I(X;Y) \geq 0$. A informação mútua nunca é negativa.
- **Simetria:** $I(X;Y)=I(Y;X)$. A informação que X fornece sobre Y é a mesma que Y fornece sobre X.
- **Independência:** $I(X;Y)=0$ se, e somente se, X e Y são variáveis aleatórias estatisticamente independentes. ²⁷ Se são independentes, conhecer uma não fornece informação alguma sobre a outra.
- **Auto-informação:** $I(X;X)=H(X)$. A quantidade de informação que uma variável contém sobre si mesma é sua própria entropia. Conhecer X remove toda a incerteza sobre X.

6.2. Intuição: O que X me diz sobre Y? Redução da Incerteza

A intuição central da informação mútua é que ela quantifica o quanto o conhecimento de uma variável reduz a incerteza sobre outra. ²¹ É a medida da informação que é "compartilhada" ou "mútua" entre as duas variáveis. Se $I(X;Y)$ é alta, significa que X e Y são fortemente dependentes; observar X nos diz muito sobre Y (e vice-versa). Se $I(X;Y)$ é baixa (próxima de zero), elas são fracamente dependentes ou independentes.

- **Analogia 1 (Dois Detetives Compartilhando Pistas sobre um Caso):** Imagine dois detetives, Alice e Bob, investigando o mesmo caso. Alice coleta um conjunto

de pistas (que têm uma certa incerteza/entropia $H(A)$) e Bob coleta outro conjunto de pistas ($H(B)$). A informação mútua $I(A;B)$ representa a sobreposição de informações entre eles: as pistas que ambos descobriram independentemente, ou as pistas de Alice que, quando reveladas a Bob, o ajudam a entender melhor suas próprias pistas (ou seja, reduzem a incerteza de Bob sobre o caso, $H(\text{Caso} | B)$, mais do que se ele não tivesse as pistas de Alice).

- **Analogia 2 (Sincronia entre Dois Dançarinos):** Considere dois dançarinos, X e Y, realizando uma coreografia.
 - Se eles estão perfeitamente sincronizados e executam movimentos idênticos ou complementares de forma coordenada, observar os movimentos do dançarino X lhe dará informação quase total sobre os movimentos do dançarino Y. Neste caso, $I(X;Y)$ é muito alta (próxima de $H(X)$ ou $H(Y)$).
 - Se os dois dançarinos estão improvisando de forma completamente aleatória e independente um do outro, observar X não lhe dirá nada sobre os movimentos de Y. Neste caso, $I(X;Y)=0$.
- **Exemplo Prático (Correlação (ou Dependência) entre Horas de Estudo e Notas em uma Prova):** Seja X o número de horas que um estudante dedica ao estudo para uma prova, e Y a nota obtida nessa prova. A informação mútua $I(X;Y)$ mediria o quanto saber o número de horas de estudo de um aluno reduz a nossa incerteza sobre a nota que ele provavelmente obterá. Se houver uma forte relação (e.g., mais estudo geralmente leva a notas mais altas, ou muito pouco estudo a notas baixas), então $I(X;Y)$ será alta. É importante notar que a informação mútua pode capturar dependências não-lineares complexas entre X e Y, não apenas correlações lineares simples (como o coeficiente de correlação de Pearson).³⁰ Por exemplo, se a relação fosse tal que tanto poucas horas quanto muitas horas de estudo levassem a notas baixas, enquanto um número moderado de horas levasse a notas altas (uma relação em forma de U invertido), a correlação linear poderia ser fraca, mas a informação mútua ainda poderia ser alta.

6.3. Visualizando Relações: O Diagrama de Venn da Entropia e Informação Mútua

Um diagrama de Venn é uma ferramenta pedagógica extremamente útil para visualizar as relações entre as diferentes medidas de entropia (marginal, conjunta, condicional) e a informação mútua para duas (ou até três) variáveis aleatórias.²⁹

Para duas variáveis X e Y:

- Desenhe dois círculos que se sobrepõem.
- O **círculo esquerdo inteiro** representa a entropia de X, $H(X)$.
- O **círculo direito inteiro** representa a entropia de Y, $H(Y)$.

- A **área de intersecção** dos dois círculos (a região onde eles se sobrepõem) representa a **Informação Mútua $I(X;Y)$** . Esta é a "informação compartilhada".
- A parte do círculo de X que **não** se sobrepõe com Y (a "lua" à esquerda) representa a **Entropia Condicional $H(X|Y)$** . Esta é a incerteza de X que resta quando Y é conhecido; é a informação única de X.
- A parte do círculo de Y que **não** se sobrepõe com X (a "lua" à direita) representa a **Entropia Condicional $H(Y|X)$** . Esta é a incerteza de Y que resta quando X é conhecido; é a informação única de Y.
- A **união total** dos dois círculos (toda a área coberta por pelo menos um dos círculos) representa a **Entropia Conjunta $H(X,Y)$** .

A partir deste diagrama, as seguintes identidades se tornam visualmente intuitivas:

- $H(X) = H(X|Y) + I(X;Y)$ (O círculo de X é a soma da sua parte única e da parte compartilhada)
- $H(Y) = H(Y|X) + I(X;Y)$ (O círculo de Y é a soma da sua parte única e da parte compartilhada)
- $H(X,Y) = H(X|Y) + H(Y|X) + I(X;Y)$ (A união é a soma das três regiões distintas)
- $H(X,Y) = H(X) + H(Y|X)$ (A união é o círculo de X mais a parte única de Y)
- $H(X,Y) = H(Y) + H(X|Y)$ (A união é o círculo de Y mais a parte única de X)
- $I(X;Y) = H(X) + H(Y) - H(X,Y)$ (A intersecção é a soma dos círculos menos a união, pois a intersecção foi contada duas vezes na soma $H(X) + H(Y)$ e a união a subtrai uma vez).
- **Analogia 1 (Diagrama de Conjuntos de Habilidades para Cientistas de Dados):** Se o círculo X representa o conjunto de "habilidades de programação" e o círculo Y representa o conjunto de "habilidades de estatística" possuídas por um grupo de profissionais. $H(X)$ e $H(Y)$ representam a diversidade ou incerteza dentro de cada área de habilidade. A intersecção $I(X;Y)$ representa as habilidades que são comuns a ambas, como "modelagem estatística com Python/R" – as habilidades de "ciência de dados" que se sobrepõem. $H(X|Y)$ seriam as habilidades de programação que não têm componente estatístico (e.g., desenvolvimento web front-end puro), e $H(Y|X)$ seriam as habilidades estatísticas que não requerem programação (e.g., teoria estatística pura). $H(X,Y)$ seria a diversidade total de habilidades combinadas de programação e estatística.
- **Analogia 2 (Cobertura de Notícias por Dois Jornais Concorrentes):** Seja $H(X)$ a totalidade das notícias (ou a incerteza sobre qual notícia um leitor encontrará) cobertas pelo Jornal X em um dia, e $H(Y)$ a mesma para o Jornal Y. A informação mútua $I(X;Y)$ representa as notícias que ambos os jornais cobriram (a sobreposição). $H(X|Y)$ são as notícias que foram exclusivas do Jornal X (não cobertas pelo Y), e $H(Y|X)$ são as notícias exclusivas do Jornal Y. $H(X,Y)$ é o

conjunto total de notícias distintas cobertas por pelo menos um dos jornais.

- Exemplo Prático (Sintomas e Doenças – Febre e Gripe):

Seja X a variável "paciente tem febre" (sim/não) e Y a variável "paciente tem gripe" (sim/não).

$H(Y)$ é a incerteza geral sobre um paciente ter gripe.

$I(X;Y)$ é a quantidade de informação que saber se o paciente tem febre (X) fornece sobre a probabilidade de ele ter gripe (Y).

$H(Y|X)$ é a incerteza restante sobre o paciente ter gripe, mesmo depois de sabermos se ele tem febre.

O diagrama de Venn pode ilustrar como a entropia total da "gripe" $H(Y)$ é composta pela informação que a "febre" fornece sobre ela ($I(X;Y)$) e a incerteza que ainda resta ($H(Y|X)$).

O diagrama de Venn é uma ferramenta pedagógica poderosa porque transforma relações algébricas abstratas em representações espaciais intuitivas, tornando os conceitos de entropia e informação mútua mais acessíveis e fáceis de memorizar.³³

Muitas pessoas pensam visualmente, e ver as entropias como "áreas" e a informação mútua como uma "intersecção" permite uma compreensão mais imediata das relações de soma e subtração entre essas quantidades do que apenas olhar para as equações. Por exemplo, a identidade $I(X;Y)=H(X)+H(Y)-H(X,Y)$ torna-se bastante óbvia ao observar que a união $H(X,Y)$ é a soma das áreas dos círculos $H(X)$ e $H(Y)$ menos a área da intersecção $I(X;Y)$ (que foi contada duas vezes na soma).

6.4. Aplicações da Informação Mútua

A Informação Mútua é um conceito versátil com inúmeras aplicações em diversos campos, devido à sua capacidade de medir a dependência geral entre variáveis.

6.4.1. Seleção de Características em Aprendizado de Máquina (Machine Learning)

Em aprendizado de máquina, um dos desafios é construir modelos preditivos eficazes usando apenas as características (features) mais relevantes de um conjunto de dados, que muitas vezes pode conter centenas ou milhares de características. A **seleção de características** visa reduzir a dimensionalidade, melhorar o desempenho do modelo (evitando overfitting), reduzir o tempo de treinamento e aumentar a interpretabilidade.

A Informação Mútua (IM) é frequentemente usada como um critério para medir a relevância de uma característica individual em relação à variável alvo (target) que se deseja prever.²⁸

- Calcula-se a IM entre cada característica e a variável alvo.

- Características com alta IM em relação ao alvo são consideradas mais informativas e, portanto, mais preditivas. Elas são selecionadas para treinar o modelo, enquanto aquelas com baixa IM podem ser descartadas.

Vantagens de usar IM para seleção de características ³¹:

- **Captura Relações Não-Lineares:** Diferentemente de medidas como o coeficiente de correlação de Pearson (que mede apenas dependência linear), a IM pode detectar qualquer tipo de relação estatística, incluindo as não-lineares.
- **Funciona para Dados Categóricos e Contínuos:** Pode ser aplicada tanto a características discretas/categóricas quanto a características contínuas.
- **Redução de Redundância (indireta):** Embora a IM entre uma feature e o alvo não meça diretamente a redundância entre features, ela ajuda a focar nas features que individualmente compartilham mais informação com o alvo. Técnicas mais avançadas podem usar IM para avaliar a redundância entre as próprias features.

Desvantagens ³¹:

- **Custo Computacional:** O cálculo da IM, especialmente para características contínuas (que requerem estimativas de densidade ou métodos baseados em k-vizinhos mais próximos), pode ser computacionalmente intensivo para conjuntos de dados muito grandes ou com um número muito elevado de características.
- **Não Indica Direcionalidade:** A IM mede a força da dependência, mas não indica a direção da relação (e.g., se uma aumenta quando a outra aumenta ou diminui).
- **Sensibilidade à Estimativa de Probabilidades:** A precisão da estimativa da IM depende de quão bem as distribuições de probabilidade subjacentes (ou suas estimativas a partir dos dados) são calculadas. Isso pode ser um desafio com amostras pequenas ou dados esparsos.

Bibliotecas de aprendizado de máquina como o scikit-learn em Python fornecem funções como `mutual_info_classif` (para alvos de classificação) e `mutual_info_regression` (para alvos de regressão) que estimam a informação mútua entre características e o alvo.³¹

- **Analogia 1 (Detetive Escolhendo as Pistas Mais Promissoras):** Um detetive chega à cena de um crime e encontra uma vasta quantidade de evidências e pistas (as características do conjunto de dados). Para resolver o caso (prever o culpado, que é a variável alvo), ele não pode investigar todas as pistas com o mesmo rigor devido a limitações de tempo e recursos. A IM ajuda o detetive a classificar as pistas pela quantidade de "informação" que cada uma parece

fornecer sobre a identidade do culpado. Ele focará nas pistas com maior IM.

- **Analogia 2 (Peneirando Ouro em um Rio):** Os dados brutos de um problema de aprendizado de máquina são como o cascalho e a areia retirados de um leito de rio. A seleção de características usando IM é como usar uma peneira especial. A peneira (o cálculo da IM) ajuda a separar os "grãos de ouro" (as características altamente informativas e preditivas) do material menos valioso ou irrelevante (características com baixa IM), que pode ser descartado.
- **Exemplo Prático (Dataset de Aprovação de Empréstimo):** Considere um conjunto de dados para prever se um pedido de empréstimo será aprovado (variável alvo: 'StatusDoEmpréstimo' - Sim/Não). As características podem incluir 'Gênero', 'EstadoCivil', 'NúmeroDeDependentes', 'NívelDeEducação', 'Autônomo', 'RendaDoAplicante', 'RendaDoCoaplicante', 'ValorDoEmpréstimo', 'PrazoDoEmpréstimo', 'HistóricoDeCrédito', 'ÁreaDaPropriedade'.²⁸ Usando `mutual_info_classif` (como no exemplo de código em 34), podemos calcular a IM entre cada uma dessas características e a variável 'StatusDoEmpréstimo'.

Python

```
from sklearn.datasets import make_classification # Exemplo genérico
from sklearn.feature_selection import mutual_info_classif
# Suponha que X_loan e y_loan são seus dados do dataset de empréstimo pré-processados
# X, y = make_classification(n_samples=100, n_features=10, n_informative=2, random_state=42) #
Exemplo didático
# mi_scores = mutual_info_classif(X, y, random_state=42)
# print(mi_scores)
```

As características que apresentarem os maiores scores de IM (e.g., 'HistóricoDeCrédito', 'RendaDoAplicante') seriam consideradas as mais relevantes para prever a aprovação do empréstimo e poderiam ser selecionadas para treinar um modelo classificador. Características com IM próxima de zero (e.g., talvez 'Gênero', se não houver viés ou relação estatística) poderiam ser descartadas.

6.4.2. Diagnóstico Médico: Relação entre Sintomas e Doenças

A Informação Mútua também encontra aplicações valiosas na área médica, particularmente no processo de diagnóstico. Ela pode ser usada para quantificar a força da associação entre sintomas (ou resultados de testes) e a presença de doenças.³⁵

- **Priorização de Testes:** A IM pode ajudar a determinar quais testes diagnósticos são mais informativos para uma determinada condição suspeita. Um teste que tem alta IM com o estado da doença (presente/ausente) é um teste que, quando realizado, reduz significativamente a incerteza do médico.

- **Identificação de Redundâncias:** Se dois testes diferentes, T1 e T2, ambos fornecem informação sobre uma doença D, mas T1 e T2 também têm alta IM entre si ($I(T1;T2)$ é alta), isso sugere que os testes são redundantes. Realizar ambos pode não adicionar muito mais informação do que realizar apenas um deles, o mais informativo ou o menos invasivo/custoso.³⁷
- **Analogia 1 (Peças de um Quebra-Cabeça Médico Complexo):** O diagnóstico de uma doença complexa é como montar um quebra-cabeça. Cada sintoma que o paciente relata e cada resultado de teste diagnóstico que chega é uma peça do quebra-cabeça. A IM ajuda o médico a entender quais peças (sintomas/testes) se "encaixam" melhor com a "imagem" de uma doença específica, ou seja, quais peças fornecem mais informação para confirmar ou descartar um diagnóstico. Também ajuda a ver se duas peças diferentes estão basicamente mostrando a mesma parte da imagem (redundância).
- **Analogia 2 (Construindo um Caso Legal para um Diagnóstico):** Um médico, como um promotor, usa evidências (sintomas, histórico do paciente, resultados de exames) para construir um "caso" a favor de um diagnóstico específico e contra outros. A IM ajuda a pesar o valor de cada peça de evidência. Uma evidência com alta IM em relação a uma doença fortalece o caso para essa doença. Se duas peças de evidência têm alta IM entre si, elas podem estar apenas reiterando o mesmo ponto.
- **Exemplo Prático (Redundância entre Testes de BUN e Creatinina em UTI):** Um estudo de Lee e Maslove (2015), citado em ³⁷, aplicou conceitos da teoria da informação, incluindo IM, para analisar a redundância em testes laboratoriais comumente realizados em Unidades de Terapia Intensiva (UTI). Eles descobriram, por exemplo, um alto nível de informação mútua entre os resultados dos testes de Nitrogênio Ureico Sanguíneo (BUN) e os de creatinina. Isso sugere que, uma vez que o resultado de um desses testes é conhecido, o resultado do outro adiciona relativamente pouca informação nova sobre o estado do paciente (pelo menos no que diz respeito à informação compartilhada por ambos). Tal análise pode levar a protocolos de teste mais otimizados, reduzindo o número de coletas de sangue desnecessárias, custos e desconforto para o paciente, sem comprometer significativamente a informação diagnóstica. A análise também poderia indicar qual dos dois testes seria preferível se apenas um pudesse ser escolhido.

6.4.3. Outras Aplicações (Breve Menção)

A utilidade da Informação Mútua se estende a muitos outros domínios:

- **Processamento de Linguagem Natural (PLN):** Usada em tarefas como alinhamento de palavras em tradução automática (medindo a dependência entre

uma palavra na língua de origem e sua tradução na língua alvo), ou para encontrar collocations (palavras que frequentemente aparecem juntas).

- **Bioinformática:** Análise de sequências genéticas (e.g., para encontrar regiões conservadas ou dependências entre posições de nucleotídeos ou aminoácidos), construção de redes de interação gênica.
- **Neurociência:** Estudo da conectividade funcional no cérebro, medindo a IM entre as atividades de diferentes regiões cerebrais (e.g., a partir de dados de fMRI ou EEG) para inferir como elas processam informação conjuntamente.
- **Finanças:** Análise de dependência entre diferentes ativos financeiros ou indicadores de mercado.

A Informação Mútua é uma medida de dependência mais geral e poderosa do que a correlação linear porque pode capturar quaisquer tipos de relações estatísticas, não se limitando às lineares.³⁰ Por exemplo, o coeficiente de correlação de Pearson mede apenas a força de uma relação linear entre duas variáveis contínuas. Se duas variáveis têm uma relação forte, mas não linear (e.g., uma relação em forma de "U" ou senoidal), a correlação de Pearson pode ser próxima de zero, sugerindo falsamente independência. No entanto, a informação mútua entre elas ainda seria alta, pois conhecer o valor de uma variável ainda reduziria significativamente a incerteza sobre o valor da outra. Essa capacidade de detectar dependências não-lineares e de aplicar-se a diferentes tipos de dados é o que torna a IM uma ferramenta tão fundamental e valiosa em campos que lidam com sistemas complexos, onde as interações raramente são simples ou lineares.

7. Além da Entropia Básica: Divergência KL e Entropia Cruzada

Além das entropias marginal, conjunta e condicional, e da informação mútua, existem outras duas medidas importantes na teoria da informação que são cruciais para comparar distribuições de probabilidade e para aplicações em aprendizado de máquina: a Divergência de Kullback-Leibler e a Entropia Cruzada.

7.1. Divergência de Kullback-Leibler ($D_{KL}(P||Q)$): Medindo a "Distância" entre Distribuições

A **Divergência de Kullback-Leibler** (pronuncia-se "Kullback-Laibler"), denotada como $DKL(P || Q)$, é uma medida de quão diferente uma distribuição de probabilidade Q é de uma distribuição de probabilidade de referência (ou "verdadeira") P .³⁸ Ela também é conhecida como **entropia relativa**.

Para duas distribuições de probabilidade discretas $P=\{p(x_i)\}$ e $Q=\{q(x_i)\}$ definidas sobre o mesmo conjunto de eventos $X=\{x_i\}$, a Divergência KL de Q em relação a P é

definida como:

$$DKL(P \parallel Q) = \sum_{x \in X} p(x) \log_2 \frac{p(x)}{q(x)}$$

(A base do logaritmo pode ser e para nats, ou 2 para bits. Usaremos base 2 aqui para consistência com bits).

Por convenção, $0 \log(0/q) = 0$ e $p \log(p/0) = \infty$ se $p > 0$.

Intuição por trás da Divergência KL:

A $DKL(P \parallel Q)$ pode ser interpretada de algumas maneiras inter-relacionadas:

- **"Perda de informação" ou "Ganho de informação":** Mede a quantidade de informação perdida quando a distribuição Q é usada para aproximar a distribuição verdadeira P . Ou, inversamente, o ganho de informação se passarmos de uma crença a priori Q para uma crença a posteriori P .⁴⁰
- **"Surpresa em excesso":** É a surpresa média adicional (ou inesperada) que experimentamos ao observar eventos que são realmente gerados pela distribuição P , mas que estávamos esperando que fossem gerados pela distribuição Q .³⁹
- **"Bits extras":** Representa o número médio de bits extras necessários para codificar amostras da distribuição P usando um esquema de codificação que foi otimizado para a distribuição Q , em comparação com usar um esquema de codificação otimizado para a própria P .⁴⁰ A entropia $H(P)$ é o número mínimo de bits para codificar P com seu código ótimo. $DKL(P \parallel Q)$ é a penalidade em bits por usar o código de Q .

Propriedades da Divergência KL:

1. **Não-negatividade:** $DKL(P \parallel Q) \geq 0$. A divergência é sempre maior ou igual a zero.
 2. **Identidade das indiscerníveis:** $DKL(P \parallel Q) = 0$ se, e somente se, $P(x) = Q(x)$ para todos os x (ou seja, as distribuições P e Q são idênticas).
 3. **Assimetria:** Em geral, $DKL(P \parallel Q) \neq DKL(Q \parallel P)$.³⁸ Por causa dessa assimetria, a Divergência KL não é uma "distância" métrica no sentido matemático estrito (pois uma métrica de distância deve ser simétrica e obedecer à desigualdade triangular, o que a DKL não faz).
- **Analogia 1 (Usar um Mapa Desatualizado para Navegar na Cidade):**
Imagine que P é o mapa perfeitamente preciso e atualizado de uma cidade. Q é um mapa antigo e desatualizado da mesma cidade.
 $DKL(P \parallel Q)$ mede o quão "ruim" ou "ineficiente" é usar o mapa antigo (Q) para tentar se locomover e encontrar destinos na cidade real (descrita por P). Ela quantifica os "erros extras" que você cometeria em média (e.g., tomar rotas mais longas, acabar em ruas sem saída que não existiam no mapa antigo) devido às discrepâncias entre Q e P . Se você usasse o mapa antigo para navegar em uma cidade que realmente correspondesse a esse mapa antigo (i.e., $DKL(Q \parallel Q)$), não

haveria "erros extras", e a divergência seria zero.

- Analogia 2 (Apostar em um Jogo de Dados com Informações Incorretas sobre sua "Viciação"):

Suponha que um dado de seis faces é viciado, e sua verdadeira distribuição de probabilidade para os resultados $\{1, 2, 3, 4, 5, 6\}$ é P . No entanto, você acredita erroneamente que o dado é honesto, então sua crença sobre a distribuição dos resultados é Q (onde $q(i)=1/6$ para todos os i).

$DKL(P || Q)$ mediria o quão "ruins" ou "mal informadas" seriam suas estratégias de aposta em média, devido à sua crença incorreta Q sobre um dado que na verdade se comporta de acordo com P . Se sua crença Q fosse igual à realidade P , $DKL(P || P)=0$, e suas apostas seriam baseadas na melhor informação possível.

- Exemplo Prático (Avaliação de um Modelo Estatístico ou de Aprendizado de Máquina):

Em muitas tarefas de modelagem, temos um conjunto de dados observados que acreditamos ter sido gerado por alguma distribuição de probabilidade "verdadeira" P . Construímos um modelo estatístico ou de aprendizado de máquina que tenta aprender ou aproximar essa distribuição P . Seja Q a distribuição de probabilidade que nosso modelo aprendeu (ou que ele prevê para os dados).

A Divergência KL, $DKL(P || Q)$, pode então ser usada para medir quão bem o nosso modelo Q está aproximando a distribuição real P . Um valor baixo de $DKL(P || Q)$ indica que o modelo é uma boa aproximação da realidade. Durante o treinamento de alguns modelos generativos, o objetivo é minimizar essa divergência.

A assimetria da DKL é uma característica crucial e informativa. $DKL(P || Q)$ e $DKL(Q || P)$ penalizam diferentes tipos de "erros" de aproximação:

- $DKL(P || Q)$ tende a infinito se $Q(x) \rightarrow 0$ para algum x onde $P(x) > 0$. Isso significa que Q é considerada uma aproximação muito ruim de P se Q atribui uma probabilidade muito baixa (ou zero) a um evento que P considera possível (ou provável). Q "perdeu" um evento importante de P . Nesse caso, o modelo Q precisa "cobrir" todas as regiões onde P tem massa de probabilidade para evitar uma divergência alta.⁴⁰
- $DKL(Q || P)$ tende a infinito se $P(x) \rightarrow 0$ para algum x onde $Q(x) > 0$. Isso significa que P é considerada muito diferente de Q (do ponto de vista de Q) se Q atribui probabilidade a um evento que P considera impossível (ou muito improvável). Q "inventou" um evento que não existe (ou é raro) em P . Nesse caso, o modelo Q tentará não colocar massa de probabilidade onde P é zero. A escolha de qual distribuição usar como P (a referência) e qual como Q (a aproximação) depende

do que se quer enfatizar ou penalizar mais na modelagem.

7.2. Entropia Cruzada ($H(P,Q)$): Custo de Codificar com a Distribuição Errada

A **Entropia Cruzada** (Cross-Entropy), denotada como $H(P,Q)$, é outra medida que compara duas distribuições de probabilidade P e Q sobre o mesmo conjunto de eventos. Ela é definida como o número médio de bits necessários para codificar eventos que são amostrados de uma distribuição verdadeira P , quando o esquema de codificação utilizado é otimizado para uma distribuição (possivelmente diferente) Q .⁴³

A fórmula para a Entropia Cruzada entre distribuições discretas P e Q é:

$$H(P,Q) = -\sum_{x \in X} p(x) \log_2 q(x)$$

Alternativamente, $H(P,Q) = E_P[-\log_2 Q(X)]$, ou seja, o valor esperado, sob a distribuição P , da auto-informação calculada usando as probabilidades de Q .

A Entropia Cruzada está intimamente relacionada com a Divergência KL. A relação é:

$$H(P,Q) = H(P) + D_{KL}(P \parallel Q)$$

onde $H(P) = -\sum p(x) \log_2 p(x)$ é a entropia de Shannon da distribuição verdadeira P .³⁸

Interpretação da Relação:

- $H(P)$ representa o número médio mínimo de bits por evento necessário para codificar dados da fonte P se usarmos o código ótimo para P (este é o limite dado pelo Teorema da Codificação de Fonte).
- $D_{KL}(P \parallel Q)$ representa os "bits extras" ou a "penalidade" média em bits por evento que pagamos por usar um código otimizado para Q em vez do código ótimo para P .
- Portanto, $H(P,Q)$ representa o **número total de bits** médio por evento que realmente usamos quando os dados vêm de P mas são codificados usando o esquema de Q [5]

Referências citadas

1. Claude Shannon | Father of Information Theory, American Engineer ..., acessado em junho 7, 2025, <https://www.britannica.com/biography/Claude-Shannon>
2. Teoria Da Informação | PDF - Scribd, acessado em junho 7, 2025, <https://www.scribd.com/document/30103260/Teoria-da-informacao>
3. Claude Shannon - Wikipedia, acessado em junho 7, 2025, https://en.wikipedia.org/wiki/Claude_Shannon
4. A Mathematical Theory of Communication* - CultureMath, acessado em junho 7, 2025, <https://culturemath.ens.fr/sites/default/files/p3-shannon.pdf>
5. Shannon information and integrated information: message and meaning - ResearchGate, acessado em junho 7, 2025, https://www.researchgate.net/publication/387104730_Shannon_information_and

- [integrated_information_message_and_meaning](#)
6. Considerações Sobre o Conceito de Entropia Na Teoria Da Informação | PDF - Scribd, acessado em junho 7, 2025,
<https://pt.scribd.com/document/495209819/Consideracoes-sobre-o-conceito-de-entropia-na-teoria-da-informacao>
 7. information theory - Intuitive explanation of entropy - Mathematics ..., acessado em junho 7, 2025,
<https://math.stackexchange.com/questions/331103/intuitive-explanation-of-entropy>
 8. Lecture 1: Entropy and Data Compression, acessado em junho 7, 2025,
<https://home.ttic.edu/~dmcalister/ttic101-07/lectures/entropy/entropy.pdf>
 9. Entropia (teoria da informação) de uma moeda de viés desconhecido : r/mathematics, acessado em junho 7, 2025,
https://www.reddit.com/r/mathematics/comments/pe7edx/entropy_information_theory_of_a_coin_of_unknown/?tl=pt-br
 10. entropy, relative entropy, and mutual information, acessado em junho 7, 2025,
https://sgfin.github.io/files/notes/Cover_and_Thomas_ch2_entropy.pdf
 11. UN C O RREC TED PR OOF - ResearchGate, acessado em junho 7, 2025,
https://www.researchgate.net/profile/Roberta-Bianco/publication/338112048_Pupil_responses_to_pitch_deviants_reflect_predictability_of_melodic_sequences/link/s/5e2c854a4585150ee7835dd2/Pupil-responses-to-pitch-deviants-reflect-predictability-of-melodic-sequences.pdf?origin=scientificContributions
 12. (PDF) ENTROPIA DESCOMPLICADA: UM GUIA PELA SEGUNDA LEI DA TERMODINÂMICAENTROPY MADE SIMPLE: A GUIDE TO THE SECOND LAW OF THERMODYNAMICSENTROPIA DESCOMPLICADA: UM GUIA PELA SEGUNDA LEI DA TERMODINÂMICAENTROPIA DESCOMPLICADA: UNA GUÍA POR LA SEGUNDA LEY DE LA TERMODINÁMICA - ResearchGate, acessado em junho 7, 2025,
https://www.researchgate.net/publication/388672993_ENTROPIA_DESCOMPLICADA_UM_GUIA_PELA_SEGUNDA_LEI_DA_TERMODINAMICAENTROPY_MADE_SIMPLE_A_GUIDE_TO_THE_SECOND_LAW_OF_THERMODYNAMICSENTROPIA_DESCOMPLICADA_UM_GUIA_PELA_SEGUNDA_LEI_DA_TERMODINAMICAENTROPIA
 13. Entropia negativa : r/AskPhysics - Reddit, acessado em junho 7, 2025,
https://www.reddit.com/r/AskPhysics/comments/15esk8z/negative_entropy/?tl=pt-br
 14. Entenda como é feita a compactação de arquivos em computadores - Tribuna de Ituverava, acessado em junho 7, 2025,
<https://www.tribunadeituverava.com.br/entenda-como-e-feita-a-compactacao-d-e-arquivos-em-computadores/>
 15. Entropia da Informação - Shannon | PPT - SlideShare, acessado em junho 7, 2025,
<https://pt.slideshare.net/slideshow/entropia-da-informao-shannon/91677177>
 16. Shannon's source coding theorem - Wikipedia, acessado em junho 7, 2025,
https://en.wikipedia.org/wiki/Shannon%27s_source_coding_theorem
 17. en.wikipedia.org, acessado em junho 7, 2025,
https://en.wikipedia.org/wiki/Shannon%27s_source_coding_theorem#:~:text=Nam

[ed%20after%20Claude%20Shannon%20the.symbol\)%20is%20less%20than%20the](#)

18. Huffman Coding | Greedy Algo-3 | GeeksforGeeks, acessado em junho 7, 2025, <https://www.geeksforgeeks.org/huffman-coding-greedy-algo-3/>
19. Huffman Code | Brilliant Math & Science Wiki, acessado em junho 7, 2025, <https://brilliant.org/wiki/huffman-encoding/>
20. Joint entropy - chemeuropa.com, acessado em junho 7, 2025, https://www.chemeuropa.com/en/encyclopedia/Joint_entropy.html
21. THEORETICAL NEUROSCIENCE I Lecture 17: Conditional entropy ..., acessado em junho 7, 2025, https://bernstein-network.de/wp-content/uploads/2021/02/18_Lecture-18-Application-to-neurons.pdf
22. Information Theory: Entropy, Markov Chains, and Huffman Coding, acessado em junho 7, 2025, https://math.nd.edu/assets/275279/leblanc_thesis.pdf
23. Unsupervised Learning: Information Theory Recap - Swyx, acessado em junho 7, 2025, <https://www.swyx.io/unsupervised-learning-information-theory-recap-4iem>
24. Unsupervised Learning: Information Theory Recap - DEV Community, acessado em junho 7, 2025, <https://dev.to/swyx/unsupervised-learning-information-theory-recap-4iem>
25. Entropy and Information Theory - Stanford Electrical Engineering, acessado em junho 7, 2025, <https://ee.stanford.edu/~gray/it.pdf>
26. Information Theory Tutorial: Mutual Information - YouTube, acessado em junho 7, 2025, <https://www.youtube.com/watch?v=d7AUaut6hso>
27. Understanding Mutual Information - Home - Matthew Kowal, acessado em junho 7, 2025, <https://mkowal2.github.io/posts/2020/01/understanding-mi/>
28. Youtube/Feature_Engineering/Feature Selection using Mutual Information - Tutorial 6.ipynb at main - GitHub, acessado em junho 7, 2025, https://github.com/atulpatelDS/Youtube/blob/main/Feature_Engineering/Feature%20Selection%20using%20Mutual%20Information%20-%20Tutorial%206.ipynb
29. (PDF) The Mutual Information Diagram for Uncertainty Visualization, acessado em junho 7, 2025, https://www.researchgate.net/publication/264543763_The_Mutual_Information_Diagram_for_Uncertainty_Visualization
30. An Intuitive View on Mutual Information - Towards Data Science, acessado em junho 7, 2025, <https://towardsdatascience.com/an-intuitive-view-on-mutual-information-db0655535f84/>
31. Getting Started With Mutual Information As Feature Selection ..., acessado em junho 7, 2025, <https://www.kaggle.com/discussions/getting-started/563669>
32. InfoMat: Leveraging Information Theory to Visualize and Understand Sequential Data, acessado em junho 7, 2025, <https://www.mdpi.com/1099-4300/27/4/357>
33. Information diagram - Wikipedia, acessado em junho 7, 2025, https://en.wikipedia.org/wiki/Information_diagram
34. mutual_info_classif — scikit-learn 1.7.0 documentation, acessado em junho 7,

2025,

https://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.mutual_info_classif.html

35. The Information Theoretic Perspective on Medical Diagnostic Inference - PMC, acessado em junho 7, 2025, <https://pmc.ncbi.nlm.nih.gov/articles/PMC6993929/>
36. "My nose is running." "Are you also coughing?": Building A Medical Diagnosis Agent with Interpretable Inquiry Logics - IJCAI, acessado em junho 7, 2025, <https://www.ijcai.org/proceedings/2022/0592.pdf>
37. Information Theory and Medical Decision Making - IOS Press Ebooks, acessado em junho 7, 2025, <https://ebooks.iospress.nl/pdf/doi/10.3233/SHTI190108>
38. Kullback-Leibler (KL) Divergence - Soupage IT Solutions, acessado em junho 7, 2025, <https://soupageit.com/ai-glossary/kullback-leibler-divergence-explained/>
39. Primers • Loss Functions - Vinija Jain, acessado em junho 7, 2025, <https://vinija.ai/concepts/loss/>
40. Variational Bayes and The Mean-Field Approximation | Bounded ..., acessado em junho 7, 2025, <https://bjlkeng.io/posts/variational-bayes-and-the-mean-field-approximation/>
41. [D] A Short Introduction to Entropy, Cross-Entropy and KL-Divergence : r/MachineLearning - Reddit, acessado em junho 7, 2025, https://www.reddit.com/r/MachineLearning/comments/7vhmp7/d_a_short_introduction_to_entropy_crossentropy/
42. [Q] Understanding intuitive difference between KL divergence and ..., acessado em junho 7, 2025, https://www.reddit.com/r/statistics/comments/10im0wr/q_understanding_intuitive_difference_between_kl/
43. A Gentle Introduction to Cross-Entropy for Machine Learning ..., acessado em junho 7, 2025, <https://machinelearningmastery.com/cross-entropy-for-machine-learning/>