# Data Compression With and Without Deep Probabilistic Models

## Recap From Last Lecture:

- Entropy: fundamental lower bound for expected code word length $L_c$ of any symbol code C:

$$H[p] \leq L_C \qquad \forall \text{ uniquely decodable } C$$

- Shannon code: reaches lower bound within less than 1 bit of overhead (per symbol)

$$H[p] \leq L_{C_{Shannon}} < H[p] + 1 \qquad H[p] = \mathbb{E}_p[-\log_2 p(x)]$$

- bonus: a Shanon code satisfies the above guarantee not only in expectation, but even individually
  for each symbol

$$\rightarrow \text{information content: } -\log_2 p(x)$$

## Recap From Last Problem Set:

Huffman Coding:
- conceptually simple algorighm (takes probability distribution as input and returns a code book as output)

- claim: Huffman Coding builds an <u>optimal</u> code book (i.e., it minimizes the expected code word length)
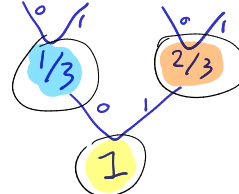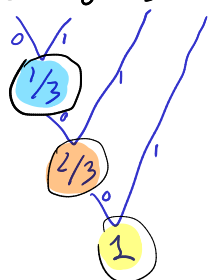
$$\rightarrow \text{Proof: today}$$

- While the code words and even their individual lengths may not be uniquely defined
  (due to ties during execution of the algorithm), the <u>expected</u> code word length is
  independent of how one breaks a tie:



expected code word length:

$$L = \sum_{x \in \mathcal{X}} p(x)\ell(x) = \frac{1}{6} \times 3 + \frac{1}{6} \times 3 + \frac{1}{3} \times 2 + \frac{1}{3} \times 1 = 2$$

$$L = \frac{1}{3} \times 1 + \frac{2}{3} \times 1 + 1 \times 1 = 2$$

*alternatively:*

expected code word length:

$$L = 2$$

$$L = \frac{1}{3} \times 1 + \frac{2}{3} \times 1 + 1 \times 1 = 2$$

## Today:

- proof of optimality of Huffman coding

- theoretical groundwork for more powerful (machine learning based) probabilistic models

# Optimality of Huffman Coding

Goal: find an optimal (uniquely decodable) symbol code for a given probability distribution p,
i.e., one with the lowest expected code word length L.

Reminder: - Among all optimal uniquely decodable symbol codes for a given p, there is at least
one prefix-free code.
→ Why?    → Kraft - Mc Millan

- We define "optimality" here as minimizing the expected code word length. This is
appropriate for many applications, but there are also use cases of data compression
where one should optimize different metrics.
→ Examples:   → real-time comm. & safety -crit.

**The Huffman algorithm for finite alphabets:** see Problem Set 1

**Theorem:** The Huffman algorithm constructs an optimal symbol code.
More precisely: assume we have
- finite alphabet $\mathcal{X}$ with $|\mathcal{X}| \geq 2$
- prob. dist. $p: \mathcal{X} \to [0,1]$ with $p(x) > 0 \quad \forall x \in \mathcal{X}$    $\Big\}$ (✱)

Then:

$\forall$ uniq. dec. symb. codes $C$ on $\mathcal{X}$ that minimize exp. code
word length w.r.t. $p$: $\exists$ a Huffman code $C_H$ with same
code word lengths, i.e., $|C(x)| = |C_H(x)| \quad \forall x \in \mathcal{X}$

**Reminder:** We may assume, without loss of generality, that C is a prefix-free code (due to the
corollary to the Kraft-McMillan Theorem).

**Lemma 1:** Assume again (*), and let C be an optimal (w.r.t. p) prefix code; let's sort the symbols s.t.

$$p(x_1) \leq p(x_2) \leq p(x_3) \leq \dots$$

We break ties by code word lengths (descendingly):

$$\text{if } p(x_i) = p(x_{i+1}) \quad \text{then } \underline{\ell_C(x_i) \geq \ell(x_{i+1})}$$

(if there are still ties after this, break them arbitrarily)

Then: (i) $\ell_C(x_1) \geq \ell_C(x_2) \geq \ell_C(x_3) \geq \dots$

(ii) $\ell_C(x_1) = \ell_C(x_2)$

**Proof of Lemma 1:**

(i) by contradiction: assume $\exists i$ with $\ell_c(x_i) < \ell_c(x_{i+1})$

we have $p(x_i) \gtrsim p(x_{i+1})$

if $p(x_i) = p(x_{i+1})$ then $\ell_c(x_i) \geq \ell_c(x_{i+1})$ ⨇

$\Rightarrow p(x_i) < p(x_{i+1})$ and $\ell_c(x_i) < \ell_c(x_{i+1})$

$\Rightarrow C$ is not optimal because we could swap

$C(x_i)$ with $C(x_{i+1})$

(would reduce $L$)

(ii) proof by contradiction, building on (i): ← we know $\ell_c(x_i) \geq \ell_c(x_{i+1})$ $\forall i$

(to show: $\ell_c(x_1) = \ell_c(x_2)$ for $C$ optimal prefix code)

assume $\ell_c(x_1) > \ell_c(x_2) \geq \ell_c(x_3) \geq \ell_c(x_4) \geq \ldots$

$\Rightarrow \ell_c(x_1) > \ell_c(x') \quad \forall x' \in X_1$

<u>claim:</u> $C$ can't be an optimal prefix code because we could drop the last bit of $C(x_1)$ and we'd still have a prefix code

e.g. $C(x_1) = \text{"0110"} \Rightarrow$
$=: \gamma$

$\begin{cases} \cdot \text{ if } \exists x': C(x') \text{ is prefix of } \gamma \\ \quad \text{then}: C(x') \text{ is also prefix of } C(x_1) \\ \quad \Rightarrow C \text{ is not prefix free (contradiction)} \\ \cdot \text{ if } \exists x': \gamma \text{ is a prefix of } C(x') \\ \quad \text{then}: \ell_c(x') \geq |\gamma| = \ell_c(x_1) - 1 \\ \quad \Rightarrow \ell_c(x') = |\gamma| \\ \quad \Rightarrow C(x') = \gamma \Rightarrow C(x') \text{ prefix of } C(x_1) \end{cases}$

(reduce $L$ by $p(x_1) > 0$)

$L = \sum_{x \in X} p(x) \ell_c(x)$

**Lemma 2:** Assume again (*), and let C be an optimal (w.r.t. p) prefix code. Then $\exists x, x' \in X$ with $x \neq x'$ and:

(i) $\ell_c(x) = \ell_c(x') \geq \ell_c(\tilde{x}) \quad \forall \tilde{x} \in X$; and

(ii) $C(x)$ & $C(x')$ only differ on last bit

**Proof of Lemma 2:** Assume that such a pair does not exist. But, from Lemma 1, we know:

$\exists x \neq x'$ that satisfy (i)
↑ ↑
$x_1$ $x_2$ (in Lemma 1)

Claim: either C is not optimal because we can drop the last bit of C(x) without violating the properties of a prefix code, or there exists a different pair of symbols that satisfy both (i) and (ii)

Proof of the claim:

Let $\gamma := C(x)$ with last bit dropped

$\forall \tilde{x} \neq x$

- if $C(\tilde{x})$ is prefix of $\gamma$ then
  $C(\tilde{x})$ is also prefix of $C(x)$ $\Rightarrow$ contradiction
- if $\gamma$ is prefix of $C(\tilde{x})$ then
  $\ell_C(\tilde{x}) \geq |\gamma| = \underbrace{\ell_C(x)}_{\text{longest codeword}} - 1$

$C(x) = "011010"$
$\phantom{C(x) =}\underbrace{\phantom{"01101}}_{=: \gamma}$

$C(\tilde{x}) = "0110"$

$C(x) = "\underset{\underset{\text{longest}}{\uparrow}}{0}11010"$
$\phantom{C(x) = }\underbrace{\phantom{"0110}}_{=: \gamma}$
codeword

$C(\tilde{x}) = "011011"$

$\Rightarrow$ { either: $C(\tilde{x})$ is one bit shorter than $C(x)$ and equal to $\gamma$
  (but then $C(\tilde{x})$ is prefix of $C(x)$, and thus $C$ is not prefix free)
  or: $C(\tilde{x})$ has same length as $C(x)$, and $\gamma$ is one bit shorter
  than both and prefix of both $\Rightarrow$ (ii) holds

**Recap:**

**Theorem:** The Huffman algorithm constructs an optimal symbol code.
More precisely: assume we have
- finite alphabet $\mathcal{X}$ with $|\mathcal{X}| \geq 2$
- probability distribution $p: \mathcal{X} \to [0,1]$ with $p(x) > 0$ $\forall x \in \mathcal{X}$ } (*)

Then:

$\forall$ uniq. dec. sym. codes $C$ on $\mathcal{X}$ that minimize the expected code
word length $L$ w.r.t. $p$: $\exists$ a Huffman code $C_H$ with the same codeword lengths,
i.e., $|C(x)| = |C_H(x)|$ $\forall x \in \mathcal{X}$.

**Lemma 1:** Assume again (*), and let $C$ be an optimal (w.r.t. p) prefix code; let's sort the symbols s.t.

$p(x_1) \leq p(x_2) \leq p(x_3) \leq \ldots$

We break ties by code word lengths (descendingly): $\ell(x_i) := |C(x_i)|$

if $p(x_i) = p(x_{i+1})$ then $\ell(x_i) \geq \ell(x_{i+1})$

(if there are still ties after this, break them arbitrarily)

Then: (i) $\ell(x_1) \geq \ell(x_2) \geq \ell(x_3) \geq \ldots$
(ii) $\ell(x_1) \overset{\uparrow}{=} \ell(x_2)$

**Lemma 2:** Assume again (*), and let $C$ be an optimal (w.r.t. p) prefix code. Then $\exists x, x' \in \mathcal{X}$ with $x \neq x'$ and:

(i) $\ell(x) = \ell(x') \geq \ell(\tilde{x})$ $\forall \tilde{x} \in \mathcal{X}$; and
(ii) $C(x)$ & $C(x')$ only differ in last bit

## Proof of the Theorem (optimality of Huffman coding):
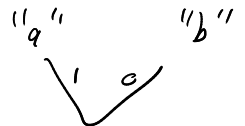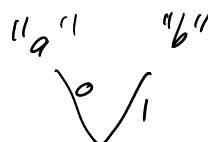
$\to$ by induction over $|\mathcal{X}|$

- base case: $|\mathcal{X}| = 2$
  $\hookrightarrow \exists$ only two optimal prefix codes:
  
  $C("a") = 0$
  $C("b") = 1$
  
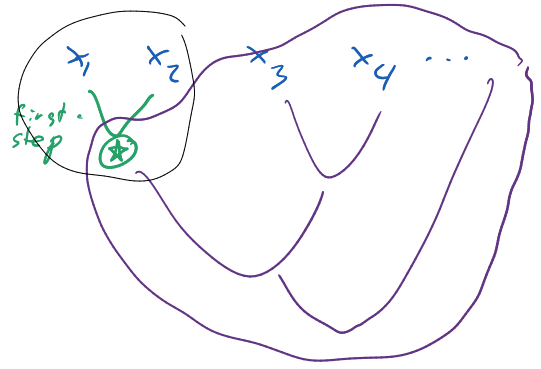  and
  
  $C("a") = 1$
  $C("b") = 0$

- induction step: $|\mathcal{X}| > 2$ assuming that theorem holds for $\forall |\mathcal{X}'| = |\mathcal{X}| - 1$

↳ from Lemma 2: $\exists x \neq x'$ with longest code words that differ only on last bit

↳ if $p(x)$ & $p(x')$ aren't among the 2 lowest probs then apply Lemma 1:

$\exists x_1 \neq x_2$ with lowest probs a also longest code words

⟹ construct a prefix code $C'$ from $C$ by swapping

only differ on last bit $(\widetilde{C(x)}, C(x'))$    with    $(C(x_1), C(x_2))$    $p(x_1)$ & $p(x_2)$ are lowest prob

all have the same (longest) length ($\Rightarrow$ swapping them doesn't change $\ell(x)$ for any $x$)

$\Rightarrow$ In $C'$: $x_1, x_2$: 
- have lowest prob
- their code words differ only last bit

Def: $\widetilde{\mathcal{X}} := (\mathcal{X} \smallsetminus \{x_1, x_2\}) \cup \{⍟\}$

$\cdot \; \tilde{p}(\tilde{x}) = \begin{cases} p(\tilde{x}) & \text{if } \tilde{x} \in \mathcal{X} \\ p(x_1) + p(x_2) & \text{if } \tilde{x} = ⍟ \end{cases}$

$\cdot \; \widetilde{C}(\tilde{x}) = \begin{cases} C'(\tilde{x}) & \text{if } \tilde{x} \in \mathcal{X} \\ C'(x) \text{ with last bit dropped} & \text{if } \tilde{x} = ⍟ \end{cases}$
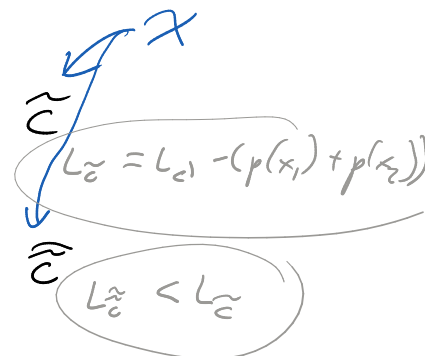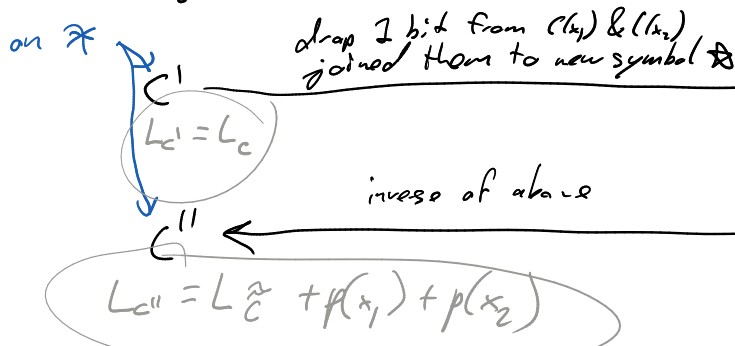
↳ Claim: $\widetilde{C}$ is an optimal prefix code on $\widetilde{\mathcal{X}}$ (with respect to $\tilde{p}$)

Proof: if it isn't an optimal prefix code then there exists a better prefix code $\widetilde{\widetilde{C}}$ on $\widetilde{\mathcal{X}}$

$\Rightarrow$ can construct a prefix code $C''$ on $\mathcal{X}$ by inverting above steps:

$$C''(x_1) := \widetilde{\widetilde{C}}(⍟) \, \| \, \text{"0"}$$
$$C''(x_2) := \widetilde{\widetilde{C}}(⍟) \, \| \, \text{"1"}$$

on $\mathcal{X}$    drop 1 bit from $C(x_1)$ & $C(x_2)$ joined them to new symbol $⍟$

$C'$    $\longrightarrow$    $\widetilde{C}$

$L_{C'} = L_C$    $L_{\widetilde{C}} = L_{C'} - (p(x_1) + p(x_2))$

inverse of above

$C''$    $\longleftarrow$    $\widetilde{\widetilde{C}}$

$L_{C''} = L_{\widetilde{\widetilde{C}}} + p(x_1) + p(x_2)$    $L_{\widetilde{\widetilde{C}}} < L_{\widetilde{C}}$

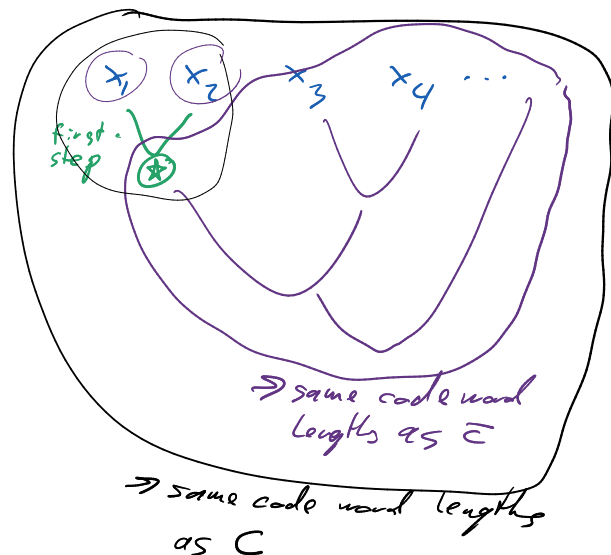$\Rightarrow L_c > L_{c''} \Rightarrow C$ was not optimal (contradiction)

Thus, $\tilde{C}$ is indeed an optimal prefix code on $\tilde{\mathcal{X}}$ (which has size $|\tilde{\mathcal{X}}| = |\mathcal{X}| - 1$ ).

$\Rightarrow$ induction hypothesis holds

$\Rightarrow \exists$ Huffman code on $\tilde{\mathcal{X}}$ with code word lengths $\tilde{C}$

$\hookrightarrow$ Recall that $x_1$ and $x_2$ (which are "contracted" in the definition of C) are two symbols with lowest probability.

$\Rightarrow$ Running Huffman algorithm on $\mathcal{X}$ also contracts contracts $x_1$ and $x_2$ in the first step. The remaining steps of the algorithm then construct a prefix code with the same code word lengths as $\tilde{C}$ on $\tilde{\mathcal{X}}$ by induction assumption.

$x_1$ $x_2$ $x_3$ $x_4$ $\cdots$

First step

$\Rightarrow$ same code word lengths as $\tilde{C}$

$\Rightarrow$ same code word lengths as C

# Remarks and Outlook:

- Huffman coding is still widely used in practice (e.g., in the "deflate" compression method used in zip/gzip and for compressed HTTP streams, in PNG, in most JPEGs, ...)

- However, Huffman coding is only an optimal symbol code. In Problem 2.4 of the current problem set (discussed tomorrow), your task is to think about the limitations of symbol codes. In the next lecture, we will start discussing so-called stream codes, which outperform Huffman coding (especially in the regime of low entropy per symbol, which is relevant for modern machine learning based data compression methods).

- On the next week's problem set (Problem Set 4), you will then use our implementation of Huffman Coding (from Problem Set 2) and you'll combine it with a machine learning model that you'll train yourself. The two components (model and entropy coding algorithm) together will result in a fully functioning (albeit ridiculously slow) deep learning based compression method for English text.

**Science Department**

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

# Probabilistic Models, Random Variables, and Correlations

Robert Bamler · 5 May 2022

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

## Quantifying Modeling Errors: The Kullback-Leibler Divergence

▶ **Qualitatively:** better probabilistic models $\Rightarrow$ better compression performance

▶ **Goal:** *quantify* loss in compression performance due to imperfect probabilistic models

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

## Reminder: Optimal Compression Performance

Consider general lossless compression setup (i.e., no longer restricted to symbol codes)

↳ Problem 2.4

▶ discrete message space $\mathcal{X}$

▶ some data source generates a message $\mathbf{x} \in \mathcal{X}$ with probability $p_{\text{data}}(\mathbf{x})$

▶ encoder $C$ maps $\mathbf{x}$ injectively to a bit string $C(\mathbf{x}) \in \{0, 1\}^*$

▶ Def: "bit rate" $R_C(\mathbf{x}) := |C(\mathbf{x})|$, i.e., length (in bits) of compressed representation

$\Rightarrow$ if $C$ is the *optimal* code for $p_{\text{data}}$ then: $R_C(\mathbf{x}) = -\log_2 p_{\text{data}}(\mathbf{x}) + \varepsilon \quad \forall \mathbf{x} \in \mathcal{X}$
(see Problem 2.4 on Problem Set 2)

$\Rightarrow$ $\boxed{\text{\textit{optimal expected} bit rate: } \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}}\left[R_{C_{\text{optimal for } p_{\text{data}}}}(\mathbf{x})\right] = H[p_{\text{data}}] + \varepsilon}$

**Problem: In practice, we don't know $p_{\text{data}}$.**

- E.g., consider the probability distribution $p_{\text{data}}$ for videos that you might take with your phone's camera.
  - huge message space $\mathcal{X}$ (all possible HD videos);
  - it's inconceivable to know $p_{\text{data}}$ exactly.
- We might, however, have some empirical samples ("training set") from $p_{\text{data}}$.
  $\Rightarrow$ Use these samples to fit a model $p_{\text{model}}$ that approximates the (unknown) $p_{\text{data}}$.
- Then, consider a compression code $C$ that is optimal for $p_{\text{model}}$:
  - bit rate of a given message $\mathbf{x}$: $\quad R_C(\mathbf{x}) = -\log_2 p_{\text{model}}(x) \qquad \forall \mathbf{x} \in \mathcal{X}$
  - *expected* bit rate: $\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}}[R_C(\mathbf{x})] = \mathbb{E}_{x \sim p_{\text{data}}}[-\log_2 p_{\text{model}}(x)] = H(p_{\text{data}}, p_{\text{model}})$

     "cross entropy"
- Idea: fit optimal $p_{\text{model}}$ by minimizing $H(p_{\text{model}}, p_{\text{data}})$

---

## Entropy, Cross Entropy, and Kullback-Leibler Divergence

- Compare:
  - true entropy of the data source: $H[p_{\text{data}}] = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}}[-\log_2 p_{\text{data}}(\mathbf{x})]$
    $\rightarrow$ fundamental lower bound of expected bit rate; :( can't evaluate
  - entropy of the model: $H[p_{\text{model}}] = \mathbb{E}_{\mathbf{x} \sim p_{\text{model}}}[-\log_2 p_{\text{model}}(\mathbf{x})]$
    $\rightarrow$ not so relevant for data compression
  - Cross entropy between data source and model: $H(p_{\text{data}}, p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}}[-\log_2 p_{\text{model}}(\mathbf{x})]$
    $\rightarrow$ practically achievable bit rate; :) can <u>estimate</u> based on samples from $p_{\text{data}}$

- **Def. "Kullback Leibler" divergence** := bit rate overhead due to modeling errors

  $$\boxed{D_{\text{KL}}(p_{\text{data}} \,||\, p_{\text{model}}) := H(p_{\text{data}}, p_{\text{model}}) - H[p_{\text{data}}] \geq 0} \quad \text{(see Problem 3.2)}$$

---

## Needed: Expressive Probabilistic Models

So far: $\mathbf{x} \in \mathcal{X}^*$ and $p_{\text{model}}(\mathbf{x}) = (p_{\text{length}}(k)) \prod_{i=1}^{k} p(x_i)$.

I.e., symbols were assumed to be "i.i.d." ("*independent* and *identically distributed*")

- *"identically distributed:"* $p$ is the same probability distribution for all $i \in \{1, \ldots, k\}$
  - We can easily overcome this limitation: $p_{\text{model}}(\mathbf{x}) = (p_{\text{length}}(k)) \prod_{i=1}^{k} p_i(x_i)$
  - Construct an individual code book $C_i$ (optimized for $p_i$) for each $i \in \{1, \ldots, k\}$.
  - Easy to see: if all $C_i$ are prefix codes then the concatenation of code words
    $C_1(x_1) \,||\, C_2(x_2) \,||\, \ldots \,||\, C_k(x_k)$ is still uniquely decodable.
- *"independent:"* the probability distribution $p_i$ does not change if we change the value of some symbol $x_j$ with $j \neq i$.
  - simplistic assumption: e.g., in English text, $p_i(\text{'u'})$ increases considerably if $x_{i-1} = \text{'q'}$.
  - This limitation is more difficult to overcome. $\rightarrow$ *correlations*