

Data Compression With and Without Deep Probabilistic Models

Lecture 2 (28 April 2022)

Recap from last lecture:

- source coding vs channel coding
- source-channel separation

- symbol codes:

↳ message $\underline{x} = (x_1, x_2, \dots, x_k)$ with $x_i \in \mathcal{X}$ $\forall i \in \{1, \dots, k\}$

↳ code book $C: \mathcal{X} \rightarrow \{0, 1, \dots, B-1\}^*$ (usually: base $B=2$)

↳ code $C^*: \mathcal{X}^* \rightarrow \{0, 1, \dots, B-1\}^*$ with $C^*(\underline{x}) := C(x_1) \parallel C(x_2) \parallel \dots \parallel C(x_k)$

"alphabet" (finite or countably infinite set)

concatenation

Recap from tutorial:

probabilistic model
of the data source

length of code word $C(x)$ (in bits)

- Def. "expected code word length": $L_c := \mathbb{E}_p[\ell_c(x)] = \sum_{x \in \mathcal{X}} p(x) \ell_c(x)$
- Def. "prefix free symbol code C" (or "prefix code" for short):
no code word $C(x)$ is the prefix of another code word $C(x')$
- Def. "uniquely decodable symbol code C^* : C^* is injective
- prefix free \Rightarrow unique decodability; but inverse is not necessarily true
- Huffman coding: algorithm that takes a probabilistic model p (on a finite alphabet) and generates a prefix code that is "tailored" for this probabilistic model.

Today: Source Coding Theorem

Two fundamental truths about lossless compression ("good news and bad news"):

- bad news: Consider a data source that produces symbols with probability distribution p . Then, there is a fundamental lower bound $H[p]$, and no uniquely decodable compression code can reach an expected code word length L that is lower than $H[p]$.

$$H_B[p] \leq L_c (= \mathbb{E}_p[\ell_c(x)]) \quad \forall \text{ uniquely decodable codes } C$$

- good news: For every data source, there exists a prefix-free (and thus uniquely decodable) code (the so-called "Shannon code") whose expected code word length approaches the fundamental lower bound $H[p]$ with an overhead of less than 1 bit per symbol.

$$H_B[p] \leq L_{c_{\text{Shannon}}} < H_B[p] + 1$$

- bonus: For finite alphabets, the Huffman coding algorithm always produces an optimal symbol code (i.e., a symbol code with lowest possible expected code word length L)

$$H_B[p] \leq L_{c_{\text{Huff}}} \leq L_{c'} \quad \forall \text{ un. dec. } C' \Rightarrow H_B[p] \leq L_{c_{\text{Huff}}} \leq L_{c_{\text{Shannon}}} < H[p] + 1$$

Kraft-McMillan Theorem

(a) \forall B-ary **uniquely decodable** symbol codes:

$$\sum_{x \in \mathcal{X}} \frac{1}{B^{\ell_c(x)}} \leq 1 \quad (\text{where } \ell_c = |C(x)|) \quad \text{"Kraft inequality"}$$

(Interpretation: we can't make code words arbitrarily short. If we shorten one code word by one bit, then we may have to make some other code word(s) longer or else our code can no longer be uniquely decodable)

(b) \forall functions $\ell: \mathcal{X} \rightarrow \mathbb{N}_0$ that satisfy Kraft ineq.

\exists B-ary **prefix code** C with $|C(x)| = \ell(x) \quad \forall x \in \mathcal{X}$
 ← "length of $C(x)$ "

Corollary: \forall unig. dec. sym. codes C :

\exists prefix code C' with $|C'(x)| = |C(x)| \quad \forall x \in \mathcal{X}$

\Rightarrow When searching for an optimal symbol code, it suffices to consider only prefix codes.

(Actually, we don't really have to search directly for an optimal symbol code. It suffices to search for an optimal assignment of code word lengths $\ell(x)$ that satisfy the Kraft inequality.

Once we have that, we can construct a prefix code with these code word lengths, see below.)

Proof of the Kraft-McMillan Theorem

Lemma: Let $s \in \mathbb{N}_0$, C unig. dec. symbol code,

$$Y_s := \{x \in \mathcal{X}^* \text{ with } |C^*(x)| = s\} \quad \leftarrow \text{"length of"}$$

Then: $|Y_s| \leq B^s$

(Proof: $\exists B^s$ distinct bit strings of length s)

\Rightarrow if $|Y_s| > B^s$ then $\exists x, x' \in \mathcal{X}^*$ with $x \neq x'$ but $C^*(x) = C^*(x')$
 \rightarrow not possible because C is unig. dec., i.e., $C^*: \mathcal{X}^* \rightarrow \{0, \dots, B-1\}^*$ is injective

Proof of part (a):

• Let $k \in \mathbb{N}$

$$\bullet r^k = \left(\sum_{x \in \mathcal{X}} B^{-\ell_c(x)} \right)^k$$

$$= \underbrace{\left(\sum_{x_1 \in \mathcal{X}} B^{-\ell_c(x_1)} \right) \left(\sum_{x_2 \in \mathcal{X}} B^{-\ell_c(x_2)} \right) \dots \left(\sum_{x_k \in \mathcal{X}} B^{-\ell_c(x_k)} \right)}_{k \text{ factors}}$$

$$= \sum_{\substack{x_1 \in \mathcal{X}, x_2 \in \mathcal{X}, \\ \dots, x_k \in \mathcal{X}}} B^{-\ell_c(x_1)} B^{-\ell_c(x_2)} \dots B^{-\ell_c(x_k)} = \sum_{x \in \mathcal{X}^k} B^{-\sum_{i=1}^k \ell_c(x_i)}$$

Claim:

$$\forall \text{ B-ary uniquely decodable symbol codes:} \\ r := \sum_{x \in \mathcal{X}} \frac{1}{B^{\ell_c(x)}} \leq 1 \quad (\text{where } \ell_c = |C(x)|)$$

(i) assume (for now) that \mathcal{X} is finite. $\Rightarrow \exists l_{\max} < \infty$ s.t. $l_c(x) \leq l_{\max} \quad \forall x \in \mathcal{X}$

$$r^k = \sum_{x \in \mathcal{X}^k} B^{-\sum_{i=1}^k l_c(x_i)} = \sum_{s=0}^{k l_{\max}} |Y_s| B^{-s} \leq k l_{\max} + 1$$

$\underbrace{\leq B^s}_{\leq 1} \text{ (lemma)}$

$$\Rightarrow r^k \leq k l_{\max} + 1 \quad \forall k \in \mathbb{N}$$

$$\Rightarrow \underbrace{\frac{r^k - 1}{k}}_{\substack{\rightarrow \infty \\ \text{if } r > 1}} \leq \underbrace{l_{\max}}_{\substack{\text{const (indep. of } k) \\ < \infty}} \quad \forall k \in \mathbb{N} \quad \Rightarrow \boxed{r \leq 1}$$

(ii) if \mathcal{X} is countably infinite: with restriction, assume that $\mathcal{X} = \mathbb{N}$

$$\text{Then } r = \sum_{x \in \mathcal{X}} \frac{1}{B^{l_c(x)}} \underset{\substack{\uparrow \\ \text{all terms are } \geq 0}}{=} \sum_{x=1}^{\infty} \frac{1}{B^{l_c(x)}} = \lim_{n \rightarrow \infty} \sum_{x=1}^n \frac{1}{B^{l_c(x)}} \leq 1$$

$\underbrace{\quad}_{\substack{\text{finite alphabet} \\ \{1, \dots, n\} \Rightarrow (i) \text{ applies}}}$

Proof of part (b) of the Kraft-McMillan Theorem:

Constructive proof, i.e., we show existence of such a prefix code by providing an explicit algorithm that constructs it for any l .

Claim:

\forall functions $l: \mathcal{X} \rightarrow \mathbb{N}_0$ that satisfy Kraft ineq. $\sum_{x \in \mathcal{X}} B^{-l(x)} \leq 1$ (length of $l(x)$)
 \exists B-ary prefix code C with $|C(x)| = l(x) \quad \forall x \in \mathcal{X}$

Algorithm (*):

- **Input:** function $l: \mathcal{X} \rightarrow \{0, \dots, B-1\}^*$ that satisfies Kraft ineq.: $\sum_{x \in \mathcal{X}} B^{-l(x)} \leq 1$

- **Output:** prefix code $C: \mathcal{X} \rightarrow \{0, \dots, B-1\}^*$ with $|C(x)| = l(x) \quad \forall x \in \mathcal{X}$

- Steps:

• sort symbols in $\mathcal{X} = \{x, x', x'', \dots\}$ s.t. $l(x) \geq l(x') \geq l(x'') \geq \dots$

• initialize $\xi = 1$

• for each $x \in \mathcal{X}$ in above order

↳ update $\xi \leftarrow \xi - B^{-l(x)}$ (\Rightarrow claim $\xi \in [0, 1)$)

↳ write $\xi = (0.\underbrace{???}_{l(x)} \dots)_B$

↳ set $C(x)$ to the first $l(x)$ bits after "0."
 (pad with trailing zeros if necessary)

Claim: The resulting code book C is prefix free. (Proof: Problem Set 2)

Example: Simplified game of Monopoly

(B=2)

x	p(x)	$-\log_2 p(x)$	$\ell(x)$	$c(x)$	$\ell'(x)$	$c'(x)$
2	1/9	3.17	① 4	1111	3	111
3	2/9	2.17	③ 3	110	2	10
4	3/9 = 1/3	1.58	⑤ 2	01	2	01
5	2/9	2.17	④ 3	101	2	00
6	1/9	3.17	② 4	1110	3	110
$H_2[p] \approx 2.20 \text{ bits}$			$L_c = \frac{26}{9} \approx 2.89$		$L_{c'} = \frac{20}{9} \approx 2.22$	

$\xi = 1$

$\cdot x=2: \xi \leftarrow 1 - 2^{-4} = (0.1111)_2$

$\cdot x=6: \xi \leftarrow (0.1111)_2 - 2^{-4} = (0.1110)_2$

$\cdot x=3: \xi \leftarrow (0.111)_2 - 2^{-3} = (0.110)_2$

$\cdot x=5: \xi \leftarrow (0.110)_2 - 2^{-3} = (0.101)_2$

$\cdot x=4: \xi \leftarrow (0.101)_2 - 2^{-2} = (0.011)_2$

Check that Kraft inequality holds for ℓ :

$$r = \sum_{x \in \mathcal{X}} 2^{-\ell(x)} = \dots = \frac{5}{8} \leq 1 \quad \checkmark$$

Questions: (1) Can we efficiently find the optimal code word lengths $\ell(x)$ that satisfy the Kraft McMillan inequality and that lead to the lowest expected code word length L ?

→ Huffman coding

(2) Can we estimate the optimal expected code word length L without having to find the whole table of optimal code word lengths $\ell(x)$?

→ Optimization problem: minimize $L := \sum_{x \in \mathcal{X}} p(x) \ell(x)$ for some fixed p over all $\ell: \mathcal{X} \rightarrow \mathbb{N}$ that satisfy Kraft ineq.

→ Spoiler: we'll find $H_B[p] \leq L_{opt} < H_B[p] + 1$

entropy $H_B[p] = \mathbb{E}_p[-\log_B p(x)] = -\sum_{x \in \mathcal{X}} p(x) \log_B p(x)$

To address question (2), we use the following strategy:

- We derive a lower bound on L .
- We show that there exists a valid assignment of code word lengths that approaches the lower bound with less than one bit of overhead.

(i) • relaxed opt. problem: minimize $L = \sum_{x \in \mathcal{X}} p(x) \ell(x)$ over all positive real valued fcts $\ell: \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ that satisfy

$$r := \sum_{x \in \mathcal{X}} B^{-\ell(x)} \leq 1$$

⇒ treat $\ell(x) \forall x \in \mathcal{X}$ as indep. variables & enforce constraint with Lagrange mult.

⇒ find stationary point

$$A := \sum_{x \in \mathcal{X}} p(x) \ell(x) + \lambda \left(\sum_{x' \in \mathcal{X}} B^{-\ell(x')} - 1 \right)$$

$$B^{-\ell(x')} = \exp(\ln(B^{-\ell(x')})) = \exp(-\ell(x') \cdot \ln B)$$

$$\forall x: 0 = \left. \frac{\partial A}{\partial l(x)} \right|_{l^*} = p(x) - (\lambda \ln B) B^{-l^*(x)}$$

$$\Rightarrow l^*(x) = -\log_B p(x) + \alpha \quad \text{with constant } \alpha = \log_B (\lambda \ln B)$$

$$0 = \left. \frac{\partial A}{\partial \lambda} \right|_{l^*} = \sum_{x \in \mathcal{X}} B^{-l^*(x)} - 1 = B^{-\alpha} \underbrace{\sum_{x \in \mathcal{X}} p(x)}_{=1} - 1 = \underbrace{B^{-\alpha} - 1}_{=0} \Rightarrow \alpha = 0$$

$$\Rightarrow l^*(x) = -\log_B p(x) \quad \text{"information content of symbol } x \text{ under the model } p"$$

$\Rightarrow l^*$ minimizes the expected code word length for under the relaxed constraint. Thus, for any other $l: \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ that satisfies the Kraft inequality, we have:

$$\sum_{x \in \mathcal{X}} p(x) l(x) \geq \sum_{x \in \mathcal{X}} p(x) l^*(x)$$

in particular, it holds for integer valued l that satisfy Kraft inequ

$\Rightarrow \forall$ uniquely dec. symbol codes C :

$$\underbrace{\sum_{x \in \mathcal{X}} p(x) l_c(x)}_{L_C} \geq \underbrace{-\sum_{x \in \mathcal{X}} p(x) \log_B p(x)}_{H_B[p]} = \mathbb{E}_p[-\log_B p(x)]$$

$$L_C \geq H_B[p] \quad \leftarrow \text{entropy}$$

(ii) How closely can we approach this lower bound (taking into account that $l(x)$ must be integer)?
 \rightarrow Answer: within an overhead of less than 1 bit per symbol.

$$\text{Proof: choose } l(x) := \lceil l^*(x) \rceil = \lceil -\log_B p(x) \rceil \geq l^*(x) \quad \forall x \in \mathcal{X}$$

$$\Rightarrow \sum_{x \in \mathcal{X}} B^{-l(x)} \leq \sum_{x \in \mathcal{X}} B^{-l^*(x)} = 1 \quad \rightarrow \text{satisfies Kraft ineq.}$$

$$\Rightarrow L_{\text{Shannon}} = \mathbb{E}_p[\lceil -\log_B p(x) \rceil] < \mathbb{E}_p[-\log_B p(x) + 1]$$

$$\Rightarrow H_B[p] \leq L_{\text{Shannon}} < H_B[p] + 1$$

"source coding theorem"

Note: If we use these code word lengths $l(x)$ and apply Algorithm (*), then the resulting prefix code C is called the "Shannon Code for the probability distribution p ".

\rightarrow see code C in the "Simplified Game of Monopoly" example above;
 more examples on Problem Set 2.