



HAL
open science

On Universal Decoding over Memoryless Channels with the Krichevsky-Trofimov Estimator

Henrique K. Miyamoto, Sheng Yang

► **To cite this version:**

Henrique K. Miyamoto, Sheng Yang. On Universal Decoding over Memoryless Channels with the Krichevsky-Trofimov Estimator. 2024 IEEE International Symposium on Information Theory (ISIT), Jul 2024, Athens, Greece. pp.1498-1503, <10.1109/ISIT57864.2024.10619414>. <hal-04720852>

HAL Id: hal-04720852

<https://hal.science/hal-04720852v1>

Submitted on 5 Oct 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire HAL, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

On Universal Decoding over Memoryless Channels with the Krichevsky–Trofimov Estimator

Henrique K. Miyamoto and Sheng Yang
Laboratoire des Signaux et Systèmes (L2S)
Université Paris-Saclay, CNRS, CentraleSupélec
Gif-sur-Yvette, France

Email: {henrique.miyamoto, sheng.yang}@centralesupelec.fr

Abstract—We study the problem of universal decoding over memoryless channels with a decoder based on the Krichevsky–Trofimov estimator. We show that this decoder is random-coding universal for codebooks of any size, i.e., despite being ignorant of the channel in use, it has asymptotically the same random-coding error exponent as the optimal maximum-likelihood decoder for that channel. Then, we incorporate this decoding rule in schemes to decode practical linear block codes and convolutional codes when the channel is unknown to the receiver. Numerical results show that efficient performance can be achieved even for moderate blocklength or constraint length.

I. INTRODUCTION

When communications take place in practice, the channel is often unknown, which precludes the use of the optimal maximum likelihood (ML) decoder. Instead, it is usually estimated by sending a training sequence known to the receiver, which can then study the statistics of the corresponding received sequence to get an estimate of the channel law. This strategy comes, nonetheless, with a trade-off: a sequence too short results in an inaccurate estimation, while a sequence too long harms the communication rate. A possible way to overcome this is to use universal decoders for the family of channels in consideration. These are decoders that, despite having no prior knowledge on the specific channel in use, can asymptotically achieve the same random-coding error exponent as the optimal ML decoder tuned for that channel.

Among known universal decoders for finite-alphabet channels, we mention the maximum mutual information (MMI) [1], [2] and maximum conditional entropy [3] rules for discrete memoryless channels; the decoder based on Lempel–Ziv (LZ) parsing for finite-state channels [3], [4]; and merged decoders for general families of channels [5]. Universal decoding for Gaussian intersymbol interference channels has been studied in [6], [7], and other recent developments for more sophisticated setups include [8]–[11].

Universal decoding has mostly remained a theoretical topic so far. Indeed, a number of potential issues appear when trying to implement them in practice. First, the results are asymptotic with blocklength, and performance with finite blocklength may be far from the promised one. Then, they are stated for random codes, and might not hold for practical codes that have strong algebraic structure. Finally, the computational complexity for implementing universal decoding rules is usually very high: they typically require comparing the universal metric of each

codeword (this has the same complexity as that of an ML decoder that compares each codeword, but becomes unpractical for large codebooks).

An exception to that are the modified stack algorithms proposed in [3], [12], [13] to decode convolutional codes. There, the usual Fano metric, which is consistent with ML decoding, is replaced by universal metrics that do not depend on the channel in use. These metrics lend themselves to asymptotic error probabilities similar to those of stack decoding with the Fano metric and ML decoding with Viterbi algorithm, both of which depend on the channel law.

In this work, we consider a universal decoder based on the (conditional) Krichevsky–Trofimov (KT) estimator [14], which enjoys a simple sequential update expression. It corresponds to the weighting of posterior distributions with Dirichlet priors, and is asymptotically optimal in the minimax sense [15]. A decoding rule based on the KT estimator has been proposed to decode convolutional codes over memoryless channels, and shown to be random-coding universal for codebooks of two codewords in [12]. A Bayesian motivation for its use, in the particular case of binary symmetric channels, was given in [13]. Here, we first extend these results, by giving a more general justification for this rule, and showing that it is in fact random-coding universal with codebooks of any size, for the family of discrete memoryless channels.

Then, we incorporate the KT-based decoder in decoding schemes for practical codes over unknown memoryless channels. For linear block codes of moderate blocklength, our decoder performs exhaustive search on a modified version of the original code. In numerical results with Golay codes, the KT decoder is able to track the block error rate of the ML decoder. For convolutional codes, we implement a modified version of the sequential stack decoding proposed in [3], [12], exploiting the structure of the KT decoder to apply further practical simplifications. Numerical results reveal a performance that is not too far from that of a stack decoder that knows the channel.

This paper is organised as follows. Section II introduces preliminary results and notation on universal decoding and the method of types. Random-coding universality of the KT decoder is derived in Section III. In Section IV, we study decoding schemes for practical codes. Section V concludes the paper with some future perspectives.

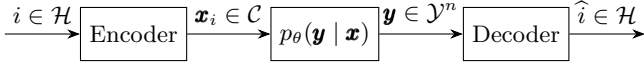


Fig. 1. Communication model.

II. PRELIMINARIES

A. Universal Decoding

Consider the communication model depicted in Fig. 1: a code $\mathcal{C} := \{\mathbf{x}_1, \dots, \mathbf{x}_M\} \subseteq \mathcal{B}_n$, of blocklength n and rate $R = (\log M)/n$, has its codewords selected from the input set $\mathcal{B}_n \subseteq \mathcal{X}^n$. A message $i \in \mathcal{H} := \{1, 2, \dots, M\}$ is uniformly selected, encoded as codeword $\mathbf{x}_i \in \mathcal{C}$, and transmitted through the channel p_θ , producing the received sequence $\mathbf{y} \in \mathcal{Y}^n$. The channel belongs to a parametric family $\mathcal{F} := \{p_\theta(\mathbf{y} | \mathbf{x}) : \theta \in \Theta\}$ with input and output alphabets, respectively, $\mathcal{X} := \{1, \dots, |\mathcal{X}|\}$ and $\mathcal{Y} := \{1, \dots, |\mathcal{Y}|\}$.

The decoder implements a rule $\phi: \mathcal{Y}^n \rightarrow \mathcal{H}$ that maps the channel output $\mathbf{y} \in \mathcal{Y}^n$ to a guess $\hat{i} = \phi(\mathbf{y})$ of the original message. The maximum *a posteriori* (MAP) rule $\phi_{\text{MAP}}(\mathbf{y}) = \arg \max_{i \in \mathcal{H}} p_\theta(\mathbf{x}_i | \mathbf{y})$ minimises the probability of error, and coincides with the ML rule $\phi_{\text{ML}}(\mathbf{y}) = \arg \max_{i \in \mathcal{H}} p_\theta(\mathbf{y} | \mathbf{x}_i)$ when messages are equiprobable. Both require knowing the channel distribution p_θ . A *universal decoder*, on the other hand, is a decoder that, despite having no prior information on the channel p_θ (other than the family \mathcal{F} to which it belongs), can asymptotically achieve the same random-coding error exponent as the optimal ML decoder tuned for p_θ .

When considering random codes, we suppose that the codewords in \mathcal{C} are chosen independently and uniformly among the sequences in \mathcal{B}_n . Let $P_{\theta, \phi}(\text{error})$ denote the average error probability (over messages and codes) when decoder ϕ is used in channel p_θ . We adopt the following definitions from [5].

Definition 1: A sequence of decoders $(u_n)_{n \in \mathbb{N}}$ is said to be *random-coding (weakly) universal* for the family of channels $\mathcal{F} = \{p_\theta(\mathbf{y} | \mathbf{x}) : \theta \in \Theta\}$ and the sequence of input sequences $(\mathcal{B}_n)_{n \in \mathbb{N}}$, if $\lim_{n \rightarrow \infty} \frac{1}{n} \log \left(\frac{P_{\theta, u_n}(\text{error})}{P_{\theta, \text{ML}}(\text{error})} \right) = 0$, for all $\theta \in \Theta$; and *random-coding strongly universal*, if $\lim_{n \rightarrow \infty} \sup_{\theta \in \Theta} \frac{1}{n} \log \left(\frac{P_{\theta, u_n}(\text{error})}{P_{\theta, \text{ML}}(\text{error})} \right) = 0$.

We borrow the following lemma from [5, Eq. (25)] (see also [3, Corollary 1]), which will be helpful in proving universality.

Lemma 1 ([5, Eq. (25)]): Consider two decoders of the form $\phi_j(\mathbf{y}) = \arg \max_{i \in \mathcal{H}} f_j(\mathbf{x}_i, \mathbf{y})$, for functions $f_j: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, $j \in \{1, 2\}$. Their average error probabilities are related by

$$\frac{P_{\theta, \phi_2}(\text{error})}{P_{\theta, \phi_1}(\text{error})} \leq \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}_n \times \mathcal{Y}^n} \frac{|\mathcal{E}_{\phi_2}(\mathbf{x}, \mathbf{y})|}{|\mathcal{E}_{\phi_1}(\mathbf{x}, \mathbf{y})|}, \quad (1)$$

where $\mathcal{E}_{\phi_j}(\mathbf{x}, \mathbf{y}) := \{\mathbf{x}' \in \mathcal{B}_n : f_j(\mathbf{x}', \mathbf{y}) \geq f_j(\mathbf{x}, \mathbf{y})\}$.

The next lemma, inspired by [12], gives an upper bound on $|\mathcal{E}_\phi(\mathbf{x}, \mathbf{y})|$ when the decoder metric is a probability $\hat{p}(\cdot | \mathbf{y})$ on \mathcal{B}_n . It formalises the fact that there can only be so many sequences $\mathbf{x}' \in \mathcal{B}_n$ with probability larger than $\hat{p}(\mathbf{x} | \mathbf{y})$.

Lemma 2: Let $\mathbf{x} \in \mathcal{B}_n$ and $\mathbf{y} \in \mathcal{Y}^n$. Consider a decoder of the form $\phi(\mathbf{y}) = \arg \max_{i \in \mathcal{H}} \frac{1}{n} \log \hat{p}(\mathbf{x}_i | \mathbf{y})$, with $\hat{p}(\cdot | \mathbf{y})$ a probability distribution over $\mathcal{B}_n \subseteq \mathcal{X}^n$. Then,

$$|\mathcal{E}_\phi(\mathbf{x}, \mathbf{y})| \leq \frac{1}{\hat{p}(\mathbf{x} | \mathbf{y})}. \quad (2)$$

Proof: Let \mathbf{X}' be a uniform random variable over \mathcal{B}_n . We have, by Markov's inequality,

$$\begin{aligned} |\mathcal{E}_\phi(\mathbf{x}, \mathbf{y})| &= \left| \left\{ \mathbf{x}' : \frac{1}{n} \log \hat{p}(\mathbf{x}' | \mathbf{y}) \geq \frac{1}{n} \log \hat{p}(\mathbf{x} | \mathbf{y}) \right\} \right| \\ &= |\mathcal{B}_n| \cdot \mathbb{P}[\hat{p}(\mathbf{X}' | \mathbf{y}) \geq \hat{p}(\mathbf{x} | \mathbf{y})] \\ &\leq |\mathcal{B}_n| \frac{\sum_{\mathbf{x}' \in \mathcal{B}_n} \frac{1}{|\mathcal{B}_n|} \hat{p}(\mathbf{x}' | \mathbf{y})}{\hat{p}(\mathbf{x} | \mathbf{y})} = \frac{1}{\hat{p}(\mathbf{x} | \mathbf{y})}. \end{aligned}$$

B. Method of Types

We briefly recall definitions from the method of types in order to introduce our notation. The *type* of a sequence $\mathbf{x} \in \mathcal{X}^n$ is the probability distribution $\pi_{\mathbf{x}}$ given by the relative frequency of symbols, i.e., $\pi_{\mathbf{x}}(x) = a_{\mathbf{x}}(x)/n$, where $a_{\mathbf{x}}(x)$ is the number of times that symbol $x \in \mathcal{X}$ appears in sequence \mathbf{x} . We denote $\mathbf{a}_{\mathbf{x}} := (a_{\mathbf{x}}(x))_{x \in \mathcal{X}}$ the vector of counts. The *type class* $\mathcal{T}^{(n)}(\pi_{\mathbf{x}})$ is the set of all sequences in \mathcal{X}^n that have type $\pi_{\mathbf{x}}$. We denote $H(\pi_{\mathbf{x}}) := -\sum_{x \in \mathcal{X}} \pi_{\mathbf{x}}(x) \log \pi_{\mathbf{x}}(x)$.

Given sequences $\mathbf{x} \in \mathcal{X}^n$ and $\mathbf{y} \in \mathcal{Y}^n$, denote $\mathbf{a}_{\mathbf{x}, \mathbf{y}} := (a_{\mathbf{x}, \mathbf{y}}(x, y))_{(x, y) \in \mathcal{X} \times \mathcal{Y}}$, where $a_{\mathbf{x}, \mathbf{y}}(x, y)$ is the number of times that $(x, y) \in \mathcal{X} \times \mathcal{Y}$ appears in the joint sequence $\mathbf{x} \otimes \mathbf{y} := ((x_1, y_1), \dots, (x_n, y_n)) \in (\mathcal{X} \times \mathcal{Y})^n$. Note that $a_{\mathbf{y}}(y) = \sum_{x \in \mathcal{X}} a_{\mathbf{x}, \mathbf{y}}(x, y)$. The *joint type* is $\pi_{\mathbf{x}, \mathbf{y}}(x, y) = a_{\mathbf{x}, \mathbf{y}}(x, y)/n$, and the *marginal type*, $\pi_{\mathbf{y}}(y) = a_{\mathbf{y}}(y)/n$. The *conditional type* $\pi_{\mathbf{x}|\mathbf{y}}$ is given by $\pi_{\mathbf{x}|\mathbf{y}}(x | y) := \pi_{\mathbf{x}, \mathbf{y}}(x, y)/\pi_{\mathbf{y}}(y)$, if $\pi_{\mathbf{y}}(y) \neq 0$, and 0 otherwise. We then have

$$\pi_{\mathbf{x}|\mathbf{y}}(\mathbf{x} | \mathbf{y}) := \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \pi_{\mathbf{x}|\mathbf{y}}(x | y)^{a_{\mathbf{x}, \mathbf{y}}(x, y)}.$$

We adopt the notations

$$H(\pi_{\mathbf{x}|\mathbf{y}} | \pi_{\mathbf{y}}) := -\sum_{y \in \mathcal{Y}} \pi_{\mathbf{y}}(y) \sum_{x \in \mathcal{X}} \pi_{\mathbf{x}|\mathbf{y}}(x | y) \log \pi_{\mathbf{x}|\mathbf{y}}(x | y), \quad (3)$$

and $I(\pi_{\mathbf{x}} : \pi_{\mathbf{y}}) := H(\pi_{\mathbf{x}}) + H(\pi_{\mathbf{y}}) - H(\pi_{\mathbf{x}, \mathbf{y}})$.

III. RANDOM-CODING UNIVERSALITY

Let us consider the class of discrete memoryless channels (DMCs). Denoting $\mathbf{x} := x_1^n := x_1 x_2 \dots x_n \in \mathcal{X}^n$ and, analogously, $\mathbf{y} := y_1^n \in \mathcal{Y}^n$, the channel law can be written as

$$p_\theta(\mathbf{y} | \mathbf{x}) = \prod_{i=1}^n p_\theta(y_i | x_i) = \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} p_\theta(y | x)^{a_{\mathbf{x}, \mathbf{y}}(x, y)}. \quad (4)$$

This family can be parametrised by vectors $\theta \in (\Delta^{(|\mathcal{Y}|-1)})^{|\mathcal{X}|}$, where $\Delta^k \subseteq \mathbb{R}^{k+1}$ denotes the k -dimensional simplex. The next lemma, borrowed from [3, Lemma 1, Eq. (23a)], states a property of the family of channels (4).

Lemma 3 ([3, Lemma 1]): Let $\mathbf{y} \in \mathcal{Y}^n$ and $\mathbf{x} \in \mathcal{B}_n$, with $\mathcal{B}_n = \mathcal{X}^n$ or $\mathcal{B}_n = \mathcal{T}^{(n)}(\pi)$ for some type π on \mathcal{X}^n . Then,

$$\begin{aligned} |\mathcal{E}_{\text{ML}}(\mathbf{x}, \mathbf{y})| &= |\{\mathbf{x}' \in \mathcal{B}_n : p_\theta(\mathbf{y} | \mathbf{x}') \geq p_\theta(\mathbf{y} | \mathbf{x})\}| \\ &\geq 2^{nH(\pi_{\mathbf{x}|\mathbf{y}}|\pi_{\mathbf{y}})}(n+1)^{-|\mathcal{X}||\mathcal{Y}|}. \end{aligned} \quad (5)$$

We are interested in studying a decoding rule based on the KT estimator [14]. The KT estimator for a sequence $\mathbf{x} \in \mathcal{X}^n$ is the probability distribution

$$p_{\text{KT}}(\mathbf{x}) = \frac{\Gamma\left(\frac{|\mathcal{X}|}{2}\right) \prod_{x \in \mathcal{X}} \Gamma\left(a_{\mathbf{x}}(x) + \frac{1}{2}\right)}{\Gamma\left(\frac{1}{2}\right)^{|\mathcal{X}|} \Gamma\left(n + \frac{|\mathcal{X}|}{2}\right)}, \quad (6)$$

where $\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt$ is the Gamma function.

To derive a conditional version of that, note that the posterior distribution can be written as

$$\begin{aligned} p_\theta(\mathbf{x} | \mathbf{y}) &= \frac{p_\theta(\mathbf{y} | \mathbf{x})p(\mathbf{x})}{p_\theta(\mathbf{y})} \\ &= \frac{|\mathcal{B}_n|^{-1}}{p_\theta(\mathbf{y})} \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} p_\theta(y | x)^{a_{\mathbf{x}, \mathbf{y}}(x, y)}, \end{aligned}$$

with $p_\theta(\mathbf{y}) := \sum_{\mathbf{x} \in \mathcal{B}_n} p_\theta(\mathbf{y} | \mathbf{x})p(\mathbf{x})$. We denote $f_\theta(x, y) := p_\theta(y | x)$, so that the normalised version

$$q_\xi(x | y) := \frac{f_\theta(x, y)}{C_\theta(y)}, \quad \text{with } C_\theta(y) := \sum_{x' \in \mathcal{X}} f_\theta(x', y)$$

can be seen as a posterior distribution on \mathcal{X} parametrised by $\xi = \xi(\theta) \in (\Delta^{|\mathcal{X}|-1})^{|\mathcal{Y}|}$. Then, we have

$$\begin{aligned} p_\theta(\mathbf{x} | \mathbf{y}) &= \frac{|\mathcal{B}_n|^{-1}}{p_\theta(\mathbf{y})} \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} C_\theta(y)^{a_{\mathbf{x}, \mathbf{y}}(x, y)} q_\xi(x | y)^{a_{\mathbf{x}, \mathbf{y}}(x, y)} \\ &= \left(\frac{\prod_{y \in \mathcal{Y}} C_\theta(y)^{a_{\mathbf{y}}(y)}}{|\mathcal{B}_n| p_\theta(\mathbf{y})} \right) \underbrace{\prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} q_\xi(x | y)^{a_{\mathbf{x}, \mathbf{y}}(x, y)}}_{(\star)}. \end{aligned}$$

For a fixed $\mathbf{y} \in \mathcal{Y}^n$, maximising $p_\theta(\mathbf{x} | \mathbf{y})$ over \mathbf{x} is thus equivalent to maximising (\star) , where q_ξ can be seen as a memoryless law and some sort of ‘backward channel’. If the channel parameters θ , and therefore ξ , are unknown, we can marginalise it with a choice of prior on $\pi(\xi)$, that is, compute

$$\widehat{p}(\mathbf{x} | \mathbf{y}) \propto \int \left(\prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} q_\xi(x | y)^{a_{\mathbf{x}, \mathbf{y}}(x, y)} \right) d\pi(\xi).$$

Choosing independent Dirichlet priors with parameters $(\frac{1}{2}, \dots, \frac{1}{2})$ for each $q_\xi(\cdot | y)$ corresponds to Jeffreys prior and is asymptotically optimal in the minimax sense [15]. In this case, the integral can be computed in much the same way as in [16, Prop. 2.15], resulting in the conditional version of the KT estimator (also in [12, Eq. (49)]):

$$p_{\text{KT}}(\mathbf{x} | \mathbf{y}) = \frac{\Gamma\left(\frac{|\mathcal{X}|}{2}\right)^{|\mathcal{Y}|} \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \Gamma\left(a_{\mathbf{x}, \mathbf{y}}(x, y) + \frac{1}{2}\right)}{\Gamma\left(\frac{1}{2}\right)^{|\mathcal{X}||\mathcal{Y}|} \prod_{y \in \mathcal{Y}} \Gamma\left(a_{\mathbf{y}}(y) + \frac{|\mathcal{X}|}{2}\right)}. \quad (7)$$

Note that this quantity only depends on the counts $a_{\mathbf{x}, \mathbf{y}}$ and $a_{\mathbf{y}}$, or, equivalently, on the joint and marginal types $\pi_{\mathbf{x}, \mathbf{y}}$ and $\pi_{\mathbf{y}}$. We will be interested in studying the universal decoder

$$\phi_{\text{KT}}(\mathbf{y}) = \arg \max_{i \in \mathcal{H}} \frac{1}{n} \log p_{\text{KT}}(\mathbf{x}_i | \mathbf{y}). \quad (8)$$

First, we provide some results on the conditional KT estimator.

Lemma 4: For $\mathbf{x} \in \mathcal{X}^n$ and $\mathbf{y} \in \mathcal{Y}^n$,

$$\begin{aligned} \log \left(\frac{\prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \pi_{\mathbf{x}|\mathbf{y}}(x | y)^{a_{\mathbf{x}, \mathbf{y}}(x, y)}}{p_{\text{KT}}(\mathbf{x} | \mathbf{y})} \right) \\ \leq \frac{(|\mathcal{X}| - 1)|\mathcal{Y}|}{2} \log n + 2|\mathcal{Y}|. \end{aligned} \quad (9)$$

Proof: Using the expression in (7), we get

$$\begin{aligned} \log \left(\frac{\prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \pi_{\mathbf{x}|\mathbf{y}}(x | y)^{a_{\mathbf{x}, \mathbf{y}}(x, y)}}{p_{\text{KT}}(\mathbf{x} | \mathbf{y})} \right) \\ = \sum_{y \in \mathcal{Y}} \log \left(\frac{\prod_{x \in \mathcal{X}} \pi_{\mathbf{x}|\mathbf{y}}(x | y)^{a_{\mathbf{x}, \mathbf{y}}(x, y)}}{\frac{\Gamma\left(\frac{|\mathcal{X}|}{2}\right) \prod_{x \in \mathcal{X}} \Gamma\left(a_{\mathbf{x}, \mathbf{y}}(x, y) + \frac{1}{2}\right)}{\Gamma\left(\frac{1}{2}\right)^{|\mathcal{X}|} \Gamma\left(a_{\mathbf{y}}(y) + \frac{|\mathcal{X}|}{2}\right)}} \right) \\ \leq \sum_{y \in \mathcal{Y}} \left(\frac{|\mathcal{X}| - 1}{2} \log a_{\mathbf{y}}(y) + 2 \right) \\ \leq \frac{(|\mathcal{X}| - 1)|\mathcal{Y}|}{2} \log n + 2|\mathcal{Y}|, \end{aligned}$$

where, for the first inequality, we apply [16, Thm. 2.16] on each subsequence of \mathbf{x} of length $a_{\mathbf{y}}(y)$ whose symbols jointly appear with the same $y \in \mathcal{Y}$ in $\mathbf{x} \otimes \mathbf{y}$. ■

Lemma 5: Given sequences $x_1^{n+1} = x_1 \cdots x_n x_{n+1} \in \mathcal{X}^{n+1}$ and $y_1^{n+1} = y_1 \cdots y_n y_{n+1} \in \mathcal{Y}^{n+1}$, we have

$$\begin{aligned} p_{\text{KT}}(x_1^{n+1} | y_1^{n+1}) &= \\ &= \left(\frac{a_{x_1^n, y_1^n}(x_{n+1}, y_{n+1}) + \frac{1}{2}}{a_{y_1^n}(y_{n+1}) + \frac{|\mathcal{X}|}{2}} \right) p_{\text{KT}}(x_1^n | y_1^n). \end{aligned} \quad (10)$$

Proof: Write (7) for sequences of length n and $n+1$, and use the fact that $\Gamma(z+1) = z\Gamma(z)$. ■

This nice sequential behaviour is analogous to the one for the simple estimator (6). It means that $p_{\text{KT}}(\mathbf{x} | \mathbf{y})$ can be sequentially computed by keeping track only of the counts $a_{\mathbf{x}, \mathbf{y}}$. In particular, we can write

$$\begin{aligned} p_{\text{KT}}(x_{n+1}^{n+k}, y_{n+1}^{n+k} | a_{x_1^n, y_1^n}) &:= \frac{p_{\text{KT}}(x_1^{n+k} | y_1^{n+k})}{p_{\text{KT}}(x_1^n | y_1^n)} \\ &= \prod_{i=n}^{n+k-1} \left(\frac{a_{x_1^i, y_1^i}(x_{i+1}, y_{i+1}) + \frac{1}{2}}{a_{y_1^i}(y_{i+1}) + \frac{|\mathcal{X}|}{2}} \right). \end{aligned} \quad (11)$$

Remark 1: In [3], the minimum conditional entropy (MCE) decoding rule $\phi_{\text{MCE}}(\mathbf{y}) = \arg \min_{i \in \mathcal{H}} H(\pi_{\mathbf{x}_i|\mathbf{y}} | \pi_{\mathbf{y}})$ was proposed for the family of memoryless channels, and shown

to be strongly random-coding universal, when either $\mathcal{B}_n = \mathcal{X}^n$ or $\mathcal{B}_n = \mathcal{T}^{(n)}(\pi)$. This rule is equivalent to

$$\begin{aligned}\phi_{\text{MCE}}(\mathbf{y}) &= \arg \max_{i \in \mathcal{H}} \pi_{\mathbf{x}_i | \mathbf{y}}(\mathbf{x}_i | \mathbf{y}) \\ &= \arg \max_{i \in \mathcal{H}} \sup_{\xi} q_{\xi}(\mathbf{x}_i | \mathbf{y}),\end{aligned}\quad (12)$$

where the supremum is taken over all memoryless laws $q_{\xi}(x | y)$, making it somewhat similar to the generalised likelihood test [17, p. 2166], with the likelihood replaced by a posterior law. When $\mathcal{B}_n = \mathcal{T}^{(n)}(\pi)$, this rule coincides with the MMI decoder $\phi_{\text{MMI}}(\mathbf{y}) = \arg \max_{i \in \mathcal{H}} I(\pi_{\mathbf{x}_i} : \pi_{\mathbf{y}})$ [2], since $H(\pi_{\mathbf{x}_i | \mathbf{y}} | \pi_{\mathbf{y}}) = H(\pi_{\mathbf{x}_i}) - I(\pi_{\mathbf{x}_i} : \pi_{\mathbf{y}})$. Finally, we note that a sequential expression analogous to (10) holds for the estimator in (12). The updating term, given in the following, has a slightly less simple expression:

$$\begin{aligned}\pi_{x_1^{n+1} | y_1^{n+1}}(x_1^{n+1} | y_1^{n+1}) &= \pi_{x_1^n | y_1^n}(x_1^n | y_1^n) \\ &\times \left(\frac{a_{x_1^n, y_1^n}(x_{n+1}, y_{n+1}) + 1}{a_{y_1^n}(y_{n+1}) + 1} \right)^{a_{x_1^n, y_1^n}(x_{n+1}, y_{n+1}) + 1} \\ &\times \left(\frac{a_{y_1^n}(y_{n+1})}{a_{x_1^n, y_1^n}(x_{n+1}, y_{n+1})} \right)^{a_{x_1^n, y_1^n}(x_{n+1}, y_{n+1})}.\end{aligned}$$

We are now ready to state our main result concerning random-coding universality of the KT decoder:

Theorem 1: Let $\mathcal{B}_n = \mathcal{X}^n$ or $\mathcal{B}_n = \mathcal{T}^{(n)}(\pi)$ the type class of some type π on \mathcal{X}^n . The decoder (8) is strongly universal for the family of memoryless channels (4).

Proof: Using (2) with $\hat{p}(\cdot | \mathbf{y}) = p_{\text{KT}}(\cdot | \mathbf{y})$ gives $|\mathcal{E}_{\text{KT}}(\mathbf{x}, \mathbf{y})| \leq \frac{1}{p_{\text{KT}}(\mathbf{x} | \mathbf{y})}$. Together with (3) and (5), we get

$$\frac{|\mathcal{E}_{\text{KT}}(\mathbf{x}, \mathbf{y})|}{|\mathcal{E}_{\text{ML}}(\mathbf{x}, \mathbf{y})|} \leq (n+1)^{|\mathcal{X}||\mathcal{Y}|} \frac{\prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \pi_{\mathbf{x} | \mathbf{y}}(x | y)^{a_{\mathbf{x}, \mathbf{y}}(x, y)}}{p_{\text{KT}}(\mathbf{x} | \mathbf{y})}.$$

Taking the logarithm on both sides and applying (9) yields

$$\begin{aligned}\log \frac{|\mathcal{E}_{\text{KT}}(\mathbf{x}, \mathbf{y})|}{|\mathcal{E}_{\text{ML}}(\mathbf{x}, \mathbf{y})|} &\leq |\mathcal{X}||\mathcal{Y}| \log(n+1) + \frac{(|\mathcal{X}|-1)|\mathcal{Y}|}{2} \log n + 2|\mathcal{Y}|.\end{aligned}$$

Finally, using (1), we conclude that

$$\begin{aligned}\frac{1}{n} \log \frac{P_{\theta, \text{KT}}(\text{error})}{P_{\theta, \text{ML}}(\text{error})} &\leq \frac{|\mathcal{X}||\mathcal{Y}|}{n} \log(n+1) + \frac{(|\mathcal{X}|-1)|\mathcal{Y}|}{2n} \log n + \frac{2|\mathcal{Y}|}{n},\end{aligned}$$

which goes to 0 as $n \rightarrow \infty$, and does not depend on θ . ■

This means that, in the ensemble of random codes, and asymptotically with blocklength n , the error exponent obtained with the KT decoder is the same as that of the optimal ML decoder. In particular, as far as asymptotic performance is concerned, the KT decoder (8) is on par with the MCE decoder (12). The particular case of this result for codebooks of size $M = 2$ appeared in [12].

IV. DECODING SCHEMES FOR PRACTICAL CODES

A. Linear Block Codes

When applying universal decoders to binary (k, n) -linear block codes, two issues appear due to the code structure. First, it contains the all-zero codeword, which is always chosen by the KT decoder, since it maximises (7), regardless of \mathbf{y} . To avoid this, all the codewords are shifted by a constant sequence known to the receiver. Second, since the code contains pairs of antipodal codewords, the decoder assigns the same metric for each pair, as it does not know whether or not bits are flipped by the channel. To deal with this, we identify antipodal codewords as representing the same message. This reduces the code rate from $\frac{k}{n}$ to $\frac{k-1}{n}$. Another solution could be to send a short training sequence to initialise the decoder.

We simulate binary channels ($|\mathcal{X}| = |\mathcal{Y}| = 2$) and encode our messages using a Golay code ($k = 12, n = 24$) [18, Sec. 1.9], with codewords shifted and antipodal pairs identified. For every block of $B = 10$ messages, a DMC with random cross-over probabilities $p := P_{Y|X}(1 | 0)$ and $q := P_{Y|X}(0 | 1)$ is uniformly drawn with $p, q \sim \mathcal{U}([0, \alpha])$, $\alpha \in [0, 1]$. Decoding is done by comparing the metrics of each codeword, as in (8), which is still feasible for moderate codebook sizes. Instead of computing the KT estimator $p_{\text{KT}}(\mathbf{x}_i | \mathbf{y})$ in (7) using only the counts $\mathbf{a}_{\mathbf{x}_i, \mathbf{y}}$ obtained from the received $\mathbf{y} \in \mathcal{Y}^n$ and codeword $\mathbf{x}_i \in \mathcal{C}$, we use a version of the counts stored in the receiver memory. For each codeword, if decoding is done above a confidence threshold (i.e., if the ratio between first and second highest metrics associated to codewords is greater than a certain threshold), then the counts obtained from the decoded codeword and corresponding received sequence are used to update the counts in the memory.

We compare the performance of the KT decoder with that of the MCE decoder (12), MMI decoder [1], [2], LZ decoder from [3], [4], and an omniscient ML decoder. All decoders (except the ML) are allowed to use a similar confidence threshold to update its rule according to previously decoded messages. In Fig. 2, we compare the block error rate (BLER) as function of the channel meta-parameter α for these schemes. We observe that KT, MCE and MMI essentially track the performance of the ML decoder, despite having no prior information on the channel. The more distant performance of LZ is explained by the slow convergence of this algorithm, which is actually universal for the broader family of finite-state channels.

B. Convolutional Codes

We consider feed-forward convolutional codes of constraint length K : at time i , a binary b -tuple u_i is input to the encoder, and the corresponding output, a binary ν -tuple, is added to $\mathbf{v}_{0,i} \in \mathbb{F}_2^\nu$, resulting in the ν -tuple \mathbf{v}_i ; this is then mapped to $\mathbf{x}_i := (x_{i,1}, \dots, x_{i,l}) = \mathcal{L}(\mathbf{v}_i)$, a sequence in \mathcal{X}^l , see [19, Fig. 5.1]. To simplify, we assume that $|\mathcal{X}|^l = 2^\nu$. Denote $\mathbf{x}_i^j \in \mathcal{X}^{l(j-i+1)}$ the sequence formed by the concatenation of the symbols in $\mathbf{x}_i, \dots, \mathbf{x}_j$. The rate of the encoder is $R = b/l$ bits per channel use. To apply universal decoding to convolutional codes, we follow [3], [12] and use a modified stack decoder.

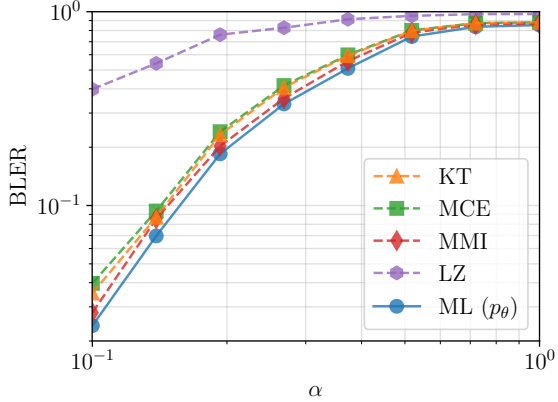


Fig. 2. BLER of different decoders for modified Golay code ($R = 11/24$) in random binary DMCs with cross-over probabilities $p, q \sim \mathcal{U}([0, \alpha])$.

Sequential decoding of convolutional codes achieves reduced complexity by avoiding to compute the metrics of all paths in the trellis, and concentrating instead on the paths with higher metrics. The basic *stack sequential decoding* [19, Ch. 6] (see also [20], [21]) works as follows¹: the algorithm keeps a stack of searched paths (of different lengths), ordered according to their metric. At each step, the path with higher metric is extended by one branch, and replaced by its 2^b successors; then, the stack is sorted by metric. The algorithm ends when the path on top of the stack (higher metric) reaches the end of the trellis. For memoryless channels of known parameters, the usual choice of metric, which is consistent with ML decoding, is the Fano metric [19, Eq. (6.1.2)].

Following [3], [12], when the channel is unknown, we replace the Fano metric by a universal metric. Let $\mathbf{x}_i^j := \mathbf{x}_i \cdots \mathbf{x}_j \in \mathcal{X}^{l(j-i+1)}$ denote an input sequence (path in the trellis) and $\mathbf{y}_i^j := \mathbf{y}_i \cdots \mathbf{y}_j \in \mathcal{Y}^{l(j-i+1)}$ an output sequence of the channel. For the KT estimator, we have the path metric

$$M_{\text{KT}}(\mathbf{x}_i^j, \mathbf{y}_i^j) := \log p_{\text{KT}}(\mathbf{x}_i^j | \mathbf{y}_i^j) + l(j-i+1) [\log |\mathcal{X}| - (R + \Delta)],$$

for $\Delta > 0$. Despite not being additive in the sense that, in general, $M_{\text{KT}}(\mathbf{x}_i^j, \mathbf{y}_i^j) \neq M_{\text{KT}}(\mathbf{x}_i^k, \mathbf{y}_i^k) + M_{\text{KT}}(\mathbf{x}_{k+1}^j, \mathbf{y}_{k+1}^j)$, this metric can be *sequentially computed*, in the sense that

$$M_{\text{KT}}(\mathbf{x}_i^{j+1}, \mathbf{y}_i^{j+1}) = M_{\text{KT}}(\mathbf{x}_i^j, \mathbf{y}_i^j) + [\log |\mathcal{X}| - (R + \Delta)] + \log p_{\text{KT}}(\mathbf{x}_{j+1}, \mathbf{y}_{j+1} | \mathbf{a}_{\mathbf{x}_i^j, \mathbf{y}_i^j}),$$

so long as one keeps track of the counts $\mathbf{a}_{\mathbf{x}_i^j, \mathbf{y}_i^j}$. Storing the counts needs no more than $|\mathcal{X}||\mathcal{Y}| \log(n+1)$ bits, since each count is a value from 0 to n , which is much less than the $n \log(|\mathcal{X}||\mathcal{Y}|)$ bits used to store the whole sequence $\mathbf{x} \otimes \mathbf{y}$. Note that this simplifies implementation when using the KT estimator, and cannot be done, for instance, with the LZ metric [3], due to the changes in the incremental parsing every time a path is extended.

¹We ignore the merging step proposed in [19], as doing so reduces complexity, and does not significantly increase error probability [19, p. 371].

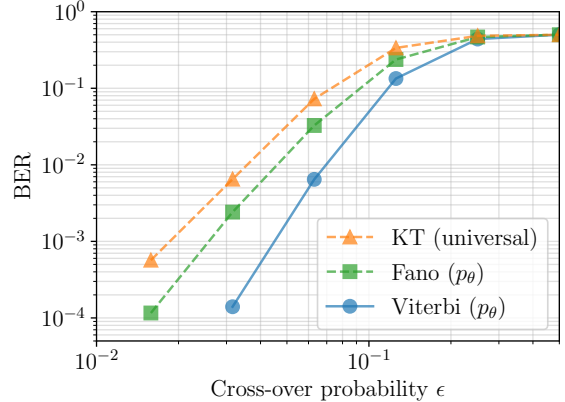


Fig. 3. BER of different decoders for convolutional code ($K = 7, R = 1/2$) in a BSC with cross-over probability ϵ .

Therefore, we can use the basic stack sequential decoding algorithm in much the same way as described previously, except that, now, along with each path \mathbf{x}_i^j in the stack, we store not only its metric $M_{\text{KT}}(\mathbf{x}_i^j, \mathbf{y}_i^j)$, but also the counts $\mathbf{a}_{\mathbf{x}_i^j, \mathbf{y}_i^j}$. As in [12], we limit the number of computations in a given branch of the trellis: when decoding u_i , if more than L_{max} nodes are visited in the branch corresponding to $u_i = u \in \mathbb{F}_2^b$, we declare $\hat{u}_i = u$. Additionally, we propose that, when this occurs, instead of reinitialising the counts to restart the process for u_{i+1} , we keep the previous values in the counts. We remark that [12, Thm. 1] assures an average bit error probability of order $O(K2^{-bK})$, over the ensemble of random convolutional codes, for rates below the cut-off rate.

In our simulations, we set $|\mathcal{X}| = |\mathcal{Y}| = 2$, $L_{\text{max}} = 2^{bK}$, $\Delta = 10^{-5}$, and use the rate-1/2 convolutional code with $K = 7$ determined by the coefficients 133_8 and 171_8 in octal notation. The channel is a binary symmetric channel (BSC) with cross-over probability $\epsilon \in [0, 1]$. In Fig. 3, we plot the bit error rate (BER) as a function of ϵ for stack decoders using the KT and Fano metrics, and for the Viterbi decoder implementing the ML rule. We see that, despite having no prior information about the channel (not even that it is symmetric), the stack decoder with KT metric has performance not too far from that of the stack decoder that knows the channel and uses the Fano metric.

V. CONCLUSION

We have studied universal decoding with the KT estimator over DMCs. We showed that this decoder is random-coding strongly universal for the family of DMCs, and integrated it in decoding schemes for linear block codes and convolutional codes. Numerical results revealed that efficient decoding can be done even for moderate blocklength or constraint length. Future perspectives include investigating analogous decoders for more general channels. If the same proof techniques are to be followed, the main difficulty so far seems to reside in deriving an analogue to Lemma 3 for those channels, which could require new combinatorial results.

REFERENCES

- [1] V. D. Goppa, "Nonprobabilistic mutual information with memory," *Problems Control Inf. Theory*, vol. 4, pp. 97–102, 1975.
- [2] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*, 2nd ed. Cambridge, UK: Cambridge Univ. Press, 2011.
- [3] J. Ziv, "Universal decoding for finite-state channels," *IEEE Trans. Inf. Theory*, vol. 31, no. 4, pp. 453–460, 1985.
- [4] A. Lapidoth and J. Ziv, "On the universality of the LZ-based decoding algorithm," *IEEE Trans. Inf. Theory*, vol. 44, no. 5, pp. 1746–1755, 1998.
- [5] M. Feder and A. Lapidoth, "Universal decoding for channels with memory," *IEEE Trans. Inf. Theory*, vol. 44, no. 5, pp. 1726–1745, 1998.
- [6] N. Merhav, "Universal decoding for memoryless Gaussian channels with a deterministic interference," *IEEE Trans. Inf. Theory*, vol. 39, no. 4, pp. 1261–1269, 1993.
- [7] W. Huleihel and N. Merhav, "Universal decoding for Gaussian inter-symbol interference channels," *IEEE Trans. Inf. Theory*, vol. 61, no. 4, pp. 1606–1618, 2015.
- [8] O. Shayevitz and M. Feder, "Universal decoding for frequency-selective fading channels," *IEEE Trans. Inf. Theory*, vol. 51, no. 8, pp. 2770–2790, 2005.
- [9] R. Averbuch and N. Merhav, "Exact random coding exponents and universal decoders for the asymmetric broadcast channel," *IEEE Trans. Inf. Theory*, vol. 64, no. 7, pp. 5070–5086, 2018.
- [10] N. Merhav, "Universal decoding using a noisy codebook," *IEEE Trans. Inf. Theory*, vol. 64, no. 4, pp. 2231–2239, 2018.
- [11] R. Tamir and N. Merhav, "Universal decoding for the typical random code and for the expurgated code," *IEEE Trans. Inf. Theory*, vol. 68, no. 4, pp. 2156–2168, 2022.
- [12] A. Lapidoth and J. Ziv, "On the decoding of convolutional codes on an unknown channel," *IEEE Trans. Inf. Theory*, vol. 45, no. 7, pp. 2321–2332, 1999.
- [13] J. K. Nelson and A. C. Singer, "Bayesian sequential detection for the BSC with unknown crossover probability," in *Proc. 2006 IEEE Int. Symp. Inf. Theory (ISIT)*, 2006, pp. 640–644.
- [14] R. Krichevsky and V. Trofimov, "The performance of universal encoding," *IEEE Trans. Inf. Theory*, vol. 27, no. 2, pp. 199–207, 1981.
- [15] B. S. Clarke and A. R. Barron, "Jeffreys' prior is asymptotically least favorable under entropy risk," *J. Statist. Planning Inference*, vol. 41, no. 1, pp. 37–60, 1994.
- [16] E. Gassiat, *Universal Coding and Order Identification by Model Selection Methods*. Cham, Switzerland: Springer, 2018.
- [17] A. Lapidoth and P. Narayan, "Reliable communication under channel uncertainty," *IEEE Trans. Inf. Theory*, vol. 44, no. 6, pp. 2148–2177, 1998.
- [18] W. C. Huffman and V. Pless, *Fundamentals of Error-Correcting Codes*. Cambridge, UK: Cambridge Univ. Press, 2003.
- [19] A. J. Viterbi and J. K. Omura, *Principles of Digital Communication and Coding*. Mineola, NY, USA: Dover, 2009.
- [20] F. Jelinek, "Fast sequential decoding algorithm using a stack," *IBM J. Res. Develop.*, vol. 13, no. 6, pp. 675–685, 1969.
- [21] K. S. Zigangirov, "Procedures of sequential decoding," in *Coding and Complexity*, G. Longo, Ed. Vienna, Austria: Springer, 1975, pp. 109–130.