

# Effects of noise on the grammar of languages

Gerrit Bauch\*

May 29, 2024

## Abstract

We study a signaling game of common interest in which a stochastic noise is perturbing the communication between an informed sender and an uninformed receiver. Despite this inhibiting factor, efficient communication is possible for any kind of noise and improves upon babbling unless the noisy channel is uninformative. Endowing a compositional message space with the Hamming distance, we explore the impact of a well-known noise channel from information theory on the grammatical structure of efficient communication. Under noise, relabeling of cells cannot be arbitrary, but has to assign distant messages to the most distant states. The more noisy the channel, the less frequent messages are used that describe states closer to the pooling action. Efficient communication under noise can be learned through the forces of evolution, but not every equilibrium is stable.

Keywords: cheap talk, noisy communication, language formation, Voronoi language

## 1 Introduction

In many situations our communication is flawed by errors having various origins. A person may stammer or slip their tongue, background noise may make it harder to understand or the recipient may suffer from a hearing impairment. It is thus natural to assume that any kind of communication is imperfect and prone to error. This noise is to be taken into account by both the speaker and the recipient in order to come to a proper understanding and buffer minor errors. This chapter proposes

---

\*For inspiring discussions and comments I thank Andreas Blume, Gerhard Jäger, Frédéric Koessler, Frank Riedel and Joel Sobel. Financial support by the DFG via grant Ri 1128-9-1 is gratefully acknowledged.  
gerrit.bauch@uni-bielefeld.de

a simple cheap talk game of common interest with a stochastic noisy channel. The state space is infinite, but behavioral limitations restrict the agents to a finite message space. Irrespective of the noise, communication is helpful. Indeed, efficient languages exist for arbitrary noise channels and strictly reduce the joint expected loss of the agents below the one if they did not communicate. We analyze a natural class of noisy channels that captures the idea that close words are more likely to be mistaken. Bayesian updates are possible even for events with zero probability, explaining why slight stammers or spelling mistakes do not disturb a proper understanding. We are interested in the structural rules of efficient languages, i.e., their *grammar*. The following conclusions can be made under a quadratic loss on a Euclidean space. The sender wants to induce a maximal spread of optimal receiver actions. To minimize their loss, the sender assigns words to different convex clusters of states. These clusters have sharp boundaries, reducing vagueness in their language. Being highly *compositional*, common languages have a specific structure which makes some kind of errors more likely to happen than others. Endowing the message space with a metric that respects its compositionality, we employ a noisy channel from information theory. Close words are more easily confused. If noise is present, the sender ideally labels states that lead to a high loss when confused by distinct words. If necessary, the sender reduces their frequency of using words that describe average states in favor of stressing extreme ones. Agents can learn locally efficient communication by means of evolution. However, not all equilibria are stable.

This chapter contributes to the economic literature on cheap talk games with common interest and a stochastic noise. Our benchmark model follows the one of Voronoi languages, Jäger et al. (2011), featuring a sender who is restricted to a finite message space to describe a state out of a continuum. Messages are pooled into cells, giving them a geometric structure that can be interpreted as the grammar of a language, cf. Jäger (2007), Gärdenfors (2004). Noisy communication has already been studied as a generalization of many influential paper. For instance, Blume et al. (2007) extend the seminal work of Crawford & Sobel (1982) to a noisy talk. The authors show that noise can improve welfare under conflict of interest. A similar observation was made by Myerson (1991). Another example of adding noise to communication is Jeitschko & Normann (2012) who extend the famous labor market model of Spence (1978). They find that under stochastic signaling subjects' strategies are closer to equilibrium play. Nowak & Krakauer (1999) show that evolution favors restricting to finitely many messages if signals can be misunderstood. Deriving an equilibrium concept for codes, Hernández & von Stengel (2014) bridge the gap between classical information theory and game theory. Being limited to broad terms by bounded rationality, Cremer et al. (2007) study efficient communication in which the receiver faces decoding cost that in-

crease in the breadth of the word, i.e., the number of states covered. If this loss depends on the state and the action taken rather than just the breadth, Sobel (2015) recovers convexity of the states lumped together for each message. Martel et al. (2019) argue for coarse communication to arise even in the absence of conflict or bounded rationality. Rubinstein (1989) shows that optimal strategies can differ significantly if the common knowledge assumption is disturbed by errors in communication. Communication can also be impaired if agents are ignorant or do not share the same vocabulary, Blume & Board (2013).

The remaining structure of the paper is as follows. Section 2 introduces the formal model. The best reply of the receiver is analyzed in Section 3. The sender's best response and the existence of efficient equilibria are explored in Section 4. Section 5 presents a concrete noisy channel with desirable properties. Structural implications for efficiency under a quadratic loss are given in Section 6. Section 7 shows that efficient communication can be learned over time. A brief summary is given in Section 8. The appendix contains proofs and calculations.

## 2 Model and notation

We adapt the setting of Voronoi languages, cf. Jäger et al. (2011). There are two players, a sender and a receiver who engage in a cheap talk game. Let  $T \subsetneq \mathbb{R}^L$ ,  $L \in \mathbb{N}_{\geq 1}$ , be a convex and compact set representing states of the world. We think of an element  $t \in T$  as an observation the sender has made and wants to inform the receiver about. Nature draws the state according to a common prior distribution described by an atomless measure  $\mu_0$  on  $T$  that is absolutely continuous w.r.t. the Lebesgue measure with a strictly positive and continuous density function  $f_0$ . The sender can communicate by choosing a message  $v$  to be sent to the receiver. The message space  $W$  available to the sender is finite, making perfect revelation of the state impossible. Communication is further frustrated by introducing stochastic error or noise  $\varepsilon: W \rightarrow \Delta(W)$  that may confound the sent message. It is possible that not the intended word  $v$  is being received, but instead an erroneous message  $w$  with probability  $\varepsilon(w|v)$ . The error admits the notion of a *Markov kernel* by interpreting  $\varepsilon: 2^W \times W \rightarrow [0, 1]$  where  $W$  is endowed with the discrete  $\sigma$ -algebra. Having observed a message  $w$ , the receiver takes an action  $\alpha(w) \in T$ , assigning it a representative type. Following the *cooperative principle* of Grice (1975), communication is first and foremost a cooperative effort. We thus assume that the agents have a common interest to match the interpretation  $\alpha(w)$  and the state  $t$ . To this end, we endow  $T$  with a norm  $\|\cdot\|$  and weigh the norm difference between the state and the interpretation by a strictly convex and strictly increasing function  $\ell: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$ . If  $t$  is the state of the world and  $s$  is the action

taken by the receiver, the loss for both parties is thus  $\ell(\|t - s\|)$ . Since the sender does not know the state prior to the play, agents face an expected loss according to  $\mu_0$ . A (pure) strategy of the sender is a ( $\mu_0$ -measurable)  $\pi: T \rightarrow W$ , called a *communication device*. The receiver's (pure) strategy is given by an *interpretation (map)*  $\alpha: W \rightarrow T$ . The expected joint loss the agents face is thus

$$L(\pi, \alpha) := \mathbb{E}_{\mu_0}[\mathbb{E}_{\varepsilon(\cdot | \pi(t))}[\ell(\|t - \alpha(w)\|)]] \quad (1)$$

$$= \int_T \sum_{w \in W} \varepsilon(w | \pi(t)) \cdot \ell(\|t - \alpha(w)\|) \mu_0(dt). \quad (2)$$

A strategy profile  $(\pi, \alpha)$  is referred to as a *language*. A language describes how information is articulated and processed. The sender and the receiver aspire to use a language that minimizes the loss of communication. The solution concept employed is that of perfect Bayesian Nash equilibria which we subsequently also refer to as *noise equilibrium* following Blume et al. (2007). Altering  $\pi$  on null sets does not change the expected loss. Furthermore, if every message is flawlessly transmitted, i.e.,  $\varepsilon(w | v) = \mathbb{1}_v(w)$  is the indicator function, the expected loss and hence the analysis reduces to the one in Jäger et al. (2011). We briefly motivate our assumptions. Compactness of  $T$  and continuity of  $\ell$  ensure integrability. Convexity of  $T$  guarantees all best replies to be in  $T$  while convexity and monotonicity of  $\ell$  makes the receiver's best replies unique.

In this chapter, we are interested in how the agents best respond to their peer's behavior, how they can achieve efficient communication and whether or not they can learn to reach a better understanding over time.

### 3 Induced beliefs and the receiver's best reply

A first step towards understanding a game is to pin down the best replies of both players aiming at characterizing equilibria. As in games of common interest a strategy profile leading to an efficient outcome entails mutual best replies and is thus an equilibrium, we get a better understanding of efficient languages. To this end assume that the receiver is aware of the sender's communication device  $\pi: T \rightarrow W$ . The knowledge of  $\pi$  allows the receiver to form expectations about the distribution of the words they are going to receive. Formally, the number

$$\lambda^\pi(w) := \mathbb{E}_{\mu_0}[\varepsilon(w | \pi(t))] = \int_T \varepsilon(w | \pi(t)) \mu_0(dt). \quad (3)$$

specifies the expected probability with which the receiver will observe the word  $w$  if the sender uses the communication device  $\pi$ .

A message  $w$  is received either if it was actually sent ( $\varepsilon(w|w) > 0$ ) or by error ( $\varepsilon(w|v) > 0$  for some  $v$ ). If  $w$  is received with positive probability  $\lambda^\pi(w) > 0$ , the receiver can use their knowledge about  $\pi$  to re-assess their informational environment. Formally, they use Bayes Rule to update their prior belief. The resulting posterior  $\mu_w^\pi$  is characterized by its density function

$$f_w^\pi(t) := \frac{f_0(t) \cdot \varepsilon(w|\pi(t))}{\lambda^\pi(w)}. \quad (4)$$

Knowing  $\pi$  and receiving the signal  $w$ , the receiver re-evaluates the chances that the sender is of some type  $t$  by applying  $\mu_w^\pi$ . If  $\lambda^\pi(w) = 0$ , we set  $\mu_w^\pi := \mu_0$  by convention.

It is worth noting that the set of induced posterior beliefs  $\{\mu_w^\pi\}_{w \in W}$  can be interpreted as a decomposition of the prior belief  $\mu_0$ . We can interpret  $\lambda^\pi$  as a distribution over posterior beliefs with support on the finite set  $\{\mu_w^\pi\}_{w \in W} \subset \Delta(W)$ . It is well-known and useful that the weighted average of the posterior beliefs sums up to the prior belief, i.e.,

$$\sum_{w \in W} \lambda^\pi(\mu_w^\pi) \cdot \mu_w^\pi = \mu_0. \quad (5)$$

More precisely, for any random variable  $X: T \rightarrow \mathbb{R}$  we have

$$\mathbb{E}_{\lambda^\pi}[\mathbb{E}_{\mu_w^\pi}[X]] = \mathbb{E}_{\mu_0}[X]. \quad (6)$$

This property is referred to as *Bayes-Plausibility*, cf. Kamenica & Gentzkow (2011). Intuitively, the sender can only induce posterior beliefs that are in expectation (w.r.t.  $\lambda^\pi$ ) the prior belief  $\mu_0$  by means of a communication device. Using the notion of  $\lambda^\pi$  and  $\mu_w^\pi$  allows us to write the expected loss in the following equivalent way

$$\begin{aligned} L(\pi, \alpha) &= \mathbb{E}_{\mu_0}[\mathbb{E}_{\varepsilon(\cdot|\pi(t))}[\ell(\|t - \alpha(w)\|)]] & (7) \\ &= \int_T \sum_{w \in W} \varepsilon(w|\pi(t)) \cdot \ell(\|t - \alpha(w)\|) \mu_0(dt) \\ &= \sum_{w \in W} \lambda^\pi(w) \cdot \int_T \lambda^\pi(w)^{-1} \cdot \varepsilon(w|\pi(t)) \cdot \ell(\|t - \alpha(w)\|) \mu_0(dt) \\ &= \mathbb{E}_{\lambda^\pi}[\mathbb{E}_{\mu_w^\pi}[\ell(\|t - \alpha(w)\|)]]. & (8) \end{aligned}$$

Expression (7) describes the expected loss as a weighted sum of the deficits that occur due to the error for each realized type  $t$ . Instead, expression (8) aggregates the expected losses under each of the posteriors and weights them according to the

probability with which the posterior is induced. The latter expression proves useful in characterizing the receiver's loss minimizing interpretation for any (induced) posterior belief.

**Lemma 3.1.** *Having any belief  $\mu \in \Delta(T)$  with positive density  $f$ , the receiver minimizes their expected loss by choosing the unique response*

$$\hat{s} \in \arg \min_{s \in T} \mathbb{E}_\mu[\ell(\|t - s\|)]. \quad (9)$$

*Proof.* All proofs are delegated to the appendix.  $\square$

Henceforth, denote by  $\hat{\alpha}(\mu)$  the unique minimizer of a receiver holding the belief  $\mu$ . We can thus define  $\hat{\alpha}(w) := \hat{\alpha}(\mu_w^\pi)$  if  $\pi$  is understood and refer to  $\hat{\alpha}$  as the unique best reply of the receiver. There are thus two ways of thinking of the receiver's best reply  $\hat{\alpha}(w)$  to a word. Firstly, their response is based on the heard word and its implicit meaning and interpretation. Secondly, the received word changes the listener's belief about the state of the world leading to their response. Furthermore, having a unique solution to the minimization problem, the receiver has no incentive to play a proper mixed strategy in equilibrium.

In a setting of common interest, we expect communication to serve the purpose of reducing misunderstandings, i.e., the expected loss. In particular, communication should foster understanding and result in a lower loss compared to the situation where individuals cannot exchange information. In the latter situation, the receiver does not receive a message from the sender and thus cannot update their belief. Sticking to the information at hand, the receiver seeks to minimize the loss given  $\mu_0$ . We call  $\alpha_{\text{pool}} := \alpha(\mu_0) = \arg \min_{s \in T} \mathbb{E}_{\mu_0}[\ell(\|t - s\|)]$  the *pooling action*. If the receiver applies the pooling action, the *pooling loss*  $L_{\text{pool}} := \mathbb{E}_{\mu_0}[\ell(\|t - \alpha_{\text{pool}}\|)]$  realizes. The following result characterizes precisely those communication devices that improve upon the pooling loss if the receiver plays a best response.

**Proposition 3.2.**  *$L(\pi, \hat{\alpha}) \leq L_{\text{pool}}$  for any communication device  $\pi$ . The inequality is strict if and only if there is a word  $w \in W$  with  $\lambda^\pi(w) > 0$  and  $\hat{\alpha}(w) \neq \alpha_{\text{pool}}$ .*

As long as the receiver knows the communication device  $\pi$ , the presence of signals cannot be detrimental to the communication. If furthermore there is a word sent with positive probability that does want the receiver to not take the pooling action, communication strictly reduces the expected loss and vice versa. The following corollary states two readily verifiable conditions under which communication does not improve over the pooling loss.

**Corollary 3.3.** *If  $\pi$  is constant or if  $v \mapsto \varepsilon(\cdot \mid v)$  is constant, then  $\mu_w^\pi = \mu_0$  for all  $w \in W$ . Consequently,  $\hat{\alpha} \equiv \alpha_{\text{pool}}$  and thus  $L(\pi, \hat{\alpha}) = L_{\text{pool}}$ .*

Although easy and intuitive, Corollary 3.3 demonstrates that there are two sources that can lead to non-beneficial communication. Firstly, the communication device may not be meaningful, i.e., received messages never offer additional information. Secondly, if the error channel is *uninformative*, i.e., does not convey any information: If everything is equally likely to be received, no matter what has been sent, the receiver cannot infer any additional information. While the latter problem is to be considered an exogenous problem of the environment, the first one lies within the power of the sender. On a warning notice, we stress that not all non-constant communication devices improve upon the pooling loss, not even if it changes the informational environment of the receiver.

**Example 3.4.** Let  $T = [-\frac{1}{2}, \frac{1}{2}]$  with uniform prior  $\mu_0$ . Assume the sender wants to inform the receiver about whether the state is close to the center  $\alpha_{\text{pool}} = 0$  of the interval or not, using the words  $C$  and  $NC$  by following the communication device depicted in Figure 1. Assume the noisy channel confuses the two messages with probability  $p < \frac{1}{2}$ , i.e, if message  $C$  ( $NC$ ) is sent, message  $NC$  ( $C$ ) is received with probability  $p$ . If the receiver gets the message  $C$ , their posterior is given by its density

$$f_C^\pi(t) = 2 \cdot \begin{cases} 1 - p & , t \in [-\frac{1}{4}, \frac{1}{4}], \\ p & , \text{otherwise} \end{cases} . \quad (10)$$

The noise channel is thus informative and makes the receiver believe that the state is more likely to lie in the center than outside of it, changing their posterior belief after the arrival of new information. However, if the expected loss is quadratic, i.e.,  $\ell \circ \|t - s\| = (t - s)^2$ , the receiver best responds to both messages with the pooling action  $\alpha_{\text{pool}} = 0$  which constitutes a noise equilibrium. The expected loss is  $L_{\text{pool}}$  even though posteriors other than the prior belief are induced. Consequently, merely changing the receiver's beliefs is not enough to reduce the expected loss below  $L_{\text{pool}}$ . One needs to make the receiver do different things for different messages. The situation is illustrated in Figure 1.

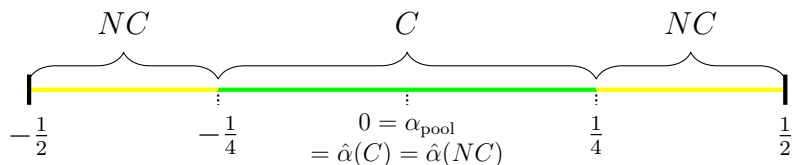


Figure 1: The language from Example 3.4 with two words that induces different posterior beliefs, but only leads to the constant pooling action. Changing beliefs is not sufficient to reduce expected loss below the pooling one.

In cheap talk games, there always are equilibria that attain the pooling loss. Assume the sender is using a constant communication device, thus provoking the constant pooling action on the receiver’s part, resulting in the pooling loss. Since the interpretation map of the receiver is constant and thus independent of any message received, the sender does not have an incentive to deviate from their constant communication device either, establishing a so-called *babbling equilibrium*<sup>1</sup>

## 4 Efficient languages

In the previous section we have established that communication can only make the agents better off in equilibrium despite the inhibiting noise. However, the expected loss cannot be reduced to zero since full separation of states is impossible due to having more states than messages. That raises the questions of what is the most efficient way of communicating and (how) can it be achieved?

While the agents cannot re-negotiate their strategies during the play, say, after the state was revealed to the sender, we can think of the two players meeting before the game. Using a meta-language they discuss their strategies prior to the play in search of an efficient communication, cf. Jäger et al. (2011). A language  $(\pi, \alpha)$  is called *efficient* if it minimizes the expected loss  $L(\pi, \alpha)$  over all possible languages. In the presented setting, this requires both agents to minimize the occurring loss over their own action sets, i.e., to play a best response each. Any efficient language is thus a noise equilibrium. Unsurprisingly, noise equilibria may improve upon the pooling loss without being efficient, see Section 6.2. We now give a simple example of an efficient language.

**Example 4.1.** Let  $T = [-\frac{1}{2}, \frac{1}{2}]$  with uniform prior  $\mu_0$ . Assume the sender has two available messages, i.e.,  $W = \{A, B\}$ . The noisy channel confounds a sent message with probability  $p$ . Assume a quadratic loss, i.e., agents lose  $(t - s)^2$  if the state is  $t$  and action  $s$  is taken. Figure 2 depicts an efficient language for every  $p \leq \frac{1}{2}$ . The sender uses their messages efficiently by cutting the interval into a left (sending  $A$ ) and a right one (sending  $B$ ). That way, the best responses of the receiver move away from the pooling action. If  $p$  increases, the best responses come closer to the pooling one as long as  $p \leq \frac{1}{2}$ . At  $p = \frac{1}{2}$  the noisy channel becomes uninformative.

---

<sup>1</sup>The name “babbling equilibrium” is best explained with mixed strategies. Assume the sender randomizes according to a fixed distribution  $G$  over messages for all  $t$ . As the state of the world and the randomly picked message are completely independent from one another, the communication device is dubbed “babbling”. The received messages on the receiver’s end will thus also be independent of the state. Gaining no insights from any message, they optimally reply with the pooling action, in turn making the sender indifferent between all their communication devices.



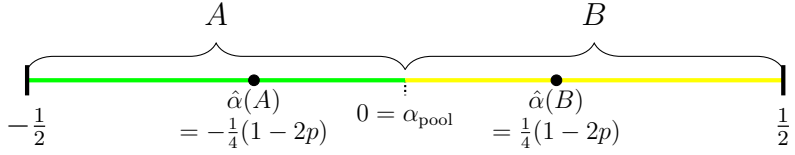


Figure 2: An efficient language in the setting of Example 4.1. The sender uses different words trying to reveal on which side of the interval the true state is situated. The closer the probability  $p$  of confusing the two words goes to  $\frac{1}{2}$ , the closer the optimal receiver responses become to the pooling action.

Efficient languages, such as the one in Example 4.1, can be explicitly computed by characterizing and internalizing a best response of the sender for any interpretation map. In contrast to the receiver's best response, the sender's best response is not unique. For instance, they can always perturb their strategy on a null set without altering the expected loss. More importantly for an economic interpretation however, the sender can be indifferent between sending two or more different words at their interim stage. The reason for this is that the corresponding actions taken by the receiver in response amount to the same loss. More precisely, fix any interpretation  $\alpha: W \rightarrow T$  of the receiver. Focusing on the sender's interim behavior<sup>2</sup> a type  $t$ -sender is indifferent between sending any word  $v \in W$  out of the (non-empty) set

$$\arg \min_{v'} \sum_{w \in W} \varepsilon(w | v') \cdot \ell(\|t - \alpha(w)\|). \quad (11)$$

Applying reasonable choices in the presence of multiple minimizers (see the proof of Theorem 4.2 for details) we can derive a measurable partition  $C^\alpha = \{C_v^\alpha\}_{v \in W}$  of  $T$ , where each  $C_v^\alpha$  consists only of types where  $v$  is a minimizing response of the sender given  $\alpha$ , i.e., an element of (11). Any such partition defines a best reply communication device by letting  $\pi^{C^\alpha}(t) = v$  if and only if  $t \in C_v^\alpha$ . It is worth mentioning that the set of types where different words can serve as a minimizing interim response, may not be a null set<sup>3</sup>. Hence, the two best replies constructed as explained above may differ perceptibly. Fortunately, the indifference sets have measure zero for the Euclidean norm as long as interpretations differ, cf. Proposition 6.4. Irrespective of the particular loss function or noise employed, we positively answer the question of existence of efficient languages.

**Theorem 4.2.** *Efficient languages exist.*

<sup>2</sup>Since the marginal distributions have full support, ex-ante and interim behavior coincide.

<sup>3</sup>For the maximum norm, the set of types that are indifferent between sending two different word might have positive measure, see, e.g., Figure 1 in Jäger et al. (2011).

Even if agents were unable to find or agree on an efficient equilibrium play, they should not use mixed strategies, i.e., randomizing their behavior. Intuitively, any randomization blurs proper understanding under common interest, thereby increasing the expected loss. If the receiver mixed between responses to a message  $w$ , the loss in communication would increase for states close to one of these receiver actions. Likewise, if the sender mixed between words that do not induce the same expected loss, the overall expected loss would increase. We give a formal explanation for why introducing randomness into the signaling or interpretation procedure makes coordination harder to achieve. Firstly note that the receiver always strictly favors a pure strategy  $\alpha: W \rightarrow T$  over a mixed one  $\tau: W \rightarrow \Delta(T)$  a posteriori by Lemma 3.1. Secondly, let  $\sigma: T \rightarrow \Delta(W)$  be a mixed strategy of the sender, specifying a probability  $\sigma(w | t)$  of playing  $w \in W$  if they are of type  $t \in T$ . The expected loss is thus given by

$$L(\sigma, \alpha) = \int_T \sum_{v \in W} \sigma(v | t) \cdot \sum_{w \in W} \varepsilon(w | v) \cdot \ell(\|t - \alpha(w)\|) \mu_0(dt). \quad (12)$$

Sending any pure  $v \in \arg \min_{v'} \sum_w \varepsilon(w | v') \cdot \ell(\|t - \alpha(w)\|)$  if the state is  $t$  is weakly decreasing the expected loss, even strictly if  $\sigma(\cdot | t)$  assigns a positive measure to any  $\tilde{v}$  not in (11). We thus have proven the following proposition.

**Proposition 4.3.** *For any mixed language  $(\sigma, \tau)$  there is a pure language  $(\pi, \alpha)$  with weakly smaller loss  $L(\pi, \alpha) \leq L(\sigma, \tau)$ .*

*The inequality can be made strict if  $\tau$  is not a pure strategy or the support of  $\sigma(\cdot | t)$  contains a word not in (11) for positive mass of states  $t$ .*

## 5 A metric-dependent noise channel

So far, the stochastic noise could be arbitrarily chosen. In the following we study a noise that perturbs communication in a natural way. Messages that are similar are more easily mistaken than ones that are dissimilar from one another. For instance, the English words “flank” and “plank” can be more easily misunderstood than “flank” and “igloo”<sup>4</sup>. In order to formalize this idea, we have to endow the message space with sufficient structure and a measure of how similar its messages are. Our approach follows the models used in information theory that are inspired by the following theory in linguistics. Human languages are *compositional*<sup>5</sup>, i.e.,

<sup>4</sup>The words taken for this example are inspired by the solutions to the *Wordle web-based word game* from July 20, May 17, May 23 2023, respectively.

<sup>5</sup>While implicitly being used in the seminal work of Frege (1892), cf. Miller (2007) for a modern reference, and even sometimes being called *Frege’s principle of compositionality*, Frege arguably never explicitly formalized the concept himself, cf. Pelletier (2001). Although com-

they are constructed from smaller building blocks. For instance, text is made up of sentences, sentences consist of words, words are sequences of letters, syllables, phonemes. The agents have a set  $\mathcal{A}$  with at least two elements at hand from which they can construct a sequence of length  $n$  to form their message space  $W = \mathcal{A}^n$ . For easier reference, we call  $\mathcal{A}$  an *alphabet* with elements that could represent *letters*, syllables or phonemes. A sequence of letters  $w \in W$  is a *word*. Restricting to a fixed length of words is common in information theory and has proven a solid baseline for coding theory Roth (2006). As a measure of distance between two words we use the *Hamming distance*, defined by

$$d: W \times W \rightarrow \mathbb{N}_0, ((w_k)_k, (v_k)_k) \mapsto \#\{k \in \{1, \dots, n\} \mid w_k \neq v_k\}. \quad (13)$$

Words are considered farther away from one another the more letters in the order of appearance differ. The Hamming distance was first studied by Hamming (1950) and is in fact a metric on  $W$ . It plays a crucial role in many applied fields related to information theory, especially coding theory and telecommunication, cf. Roth (2006).

The noisy channel we are now introducing behaves well with the Hamming distance and is best understood defining it on letters first. Let  $\tilde{\varepsilon}: \mathcal{A} \rightarrow \Delta(\mathcal{A})$  denote the function

$$a \mapsto \tilde{\varepsilon}(\cdot \mid a), \quad \tilde{\varepsilon}(b \mid a) := \begin{cases} 1 - p & , b = a, \\ \frac{p}{\#\mathcal{A} - 1} & , b \neq a \end{cases}, \quad (14)$$

where the exogenous parameter  $p \in [0, 1]$  is the *crossover probability*, i.e., the probability of wrongly transmitting one intended letter  $a$ . In case of an error each of the other  $\#\mathcal{A} - 1$  symbols is assumed to be equally likely received. It is a well-known noise that is used to model error transitions in telecommunication, data storage, but also finds application in DNA heritage of cell-divisions, cf. MacKay (2002) and Cover & Thomas (2006)).

Depending on the concrete scenario, errors can feature different levels of correlation. We follow the branch of literature that assumes independent occurrences of errors for each letter. To this end we can extend the  $\#\mathcal{A}$ -ary symmetric error channel with crossover probability  $p$  to  $W = \mathcal{A}^n$  by gluing  $n$  independent copies of  $\tilde{\varepsilon}$  together. The result is called  *$\#\mathcal{A}$ -ary symmetric channel of length  $n$  with crossover probability  $p$* . Using the Hamming distance, the transition probabilities are given by

$$\varepsilon(w \mid v) := (1 - p)^{n-d(w,v)} \cdot \left( \frac{p}{\#\mathcal{A} - 1} \right)^{d(w,v)}. \quad (15)$$

---

elling, deriving meaning from compositionality alone is a debatable practice among linguists and challenged by other concepts, such as *contextuality*, cf. Szabó (2022).

The probability  $\varepsilon(w | v)$  that  $w$  is received if  $v$  is sent only depends on the Hamming distance  $d(w, v)$  and the crossover probability  $p$ . Especially,  $\varepsilon(w | v) = \varepsilon(v | w)$ . For fixed  $v \in W$  and  $d \in \{0, \dots, n\}$  there are precisely  $\binom{n}{d} \cdot (\#\mathcal{A} - 1)^d$  different words  $w$  with  $d = d(w, v)$  in  $W$ . Hence,  $\varepsilon$  follows a binomial distribution on the set of families  $\{w \in W | d(w, v) = d\}_{d=0}^n$ . Being mainly interested in the transition probability from  $v$  to a particular  $w$  rather than to such a family of words,  $\varepsilon$  is itself the natural probability distribution to consider. For convenience, we write  $m := \#\mathcal{A} - 1$  and  $\tilde{p} := \frac{p}{(1-p)^m}$  and can rewrite  $\varepsilon(w | v) = (1-p)^n \cdot \tilde{p}^{d(w,v)}$  which is often convenient. The following immediate statements capture the nice interplay between the noise channel  $\varepsilon$  and the Hamming-distance for different crossover probabilities  $p$ .

**Remark 5.1.** Suppose  $v \in W$  is the word sent.

- (i) If  $p = \tilde{p} = 0$ , there is no noise, i.e.,  $v$  is received with probability 1.
- (ii) If  $0 < p < \frac{m}{m+1}$ , i.e.,  $0 < \tilde{p} < 1$ ,  $\varepsilon(w | v)$  is decreasing in  $d(w, v)$ . Especially it is most likely to receive  $v$ .
- (iii) If  $p = \frac{m}{m+1}$ , i.e.,  $\tilde{p} = 1$ ,  $\varepsilon$  is uninformative.
- (iv) If  $\frac{m}{m+1} < p < 1$ , i.e.,  $1 < \tilde{p} < \infty$ ,  $\varepsilon(w | v)$  is increasing in  $d(w, v)$ . Especially, words with maximum distance  $d(w, v) = n$  are most likely received.
- (v) If  $p = 1$ , i.e.,  $\tilde{p} = \infty$ , only words with maximum distance  $d(w, v) = n$  are received.

For  $p = 0$  there is no noise in communication and our model coincides with Jäger et al. (2011). If  $0 < p < \frac{m}{m+1}$  the  $\#\mathcal{A}$ -ary symmetric channel of length  $n$  captures our intuition that words that are close in the Hamming distance are more likely to be confounded. This property is most pronounced for  $p \approx 0$  and vanishes completely at the *uninformativeness bound*  $p = \frac{m}{m+1}$ . In fact, we find a continuous and monotonic loss in the amount of information that can be transmitted through the noisy channel for increasing  $p$  measured by *Shannon entropy*, cf. Section 5.2. At the uninformativeness bound,  $\varepsilon(\cdot | v)$  is the constant uniform distribution on  $W$  for all  $v$ , making the sent and received word independent, recall also Corollary 3.3.

Interestingly, for  $\frac{m}{m+1} < p$  informative communication can take place again, even less efficient than before. The receiver now suspects the received word to stem most likely from one among those having the maximal distance to the it. Since this set contains more than one word if  $m > 1$ , the receiver cannot pin down a single most probable word and is thus less confident about their guess of

the originally sent word.<sup>6</sup> Our discussion of entropy in Section 5.2 quantifies this observation.

## 5.1 Bayesian updates and limit cases

In many realistic scenarios, errors in communication may be present, but not too abundant. In the following, we will thus focus on the case  $0 < p < \frac{m}{m+1}$  for which errors occur with positive probability, the noisy channel is not uninformative and has the property that close words in the Hamming distance are more easily misunderstood.<sup>7</sup> In the following, we will derive the receiver's posterior beliefs under  $\varepsilon$ . We also study how the noisy channel behaves in the limit cases if the crossover probability  $p$  goes to zero or the unformativeness bound and find suitable updates even for events of probability zero.

Since  $\varepsilon(w | v) > 0$  for all  $w, v$  the receiver can always use Bayes rule to update their prior belief. A short calculation yields posteriors beliefs with densities of the form

$$f_w^\pi(t) = f_0(t) \cdot \left( \int_T \tilde{p}^{d(w, \pi(t')) - d(w, \pi(t))} \mu_0(dt') \right)^{-1} \quad (16)$$

for a received word  $w$ . Note that the integrand is continuous in  $p \in [0, 1)$ . The employed error channel allows us to study the limit cases for the receivers posterior belief for  $p \rightarrow 0$  and  $p \rightarrow \frac{m}{m+1}$ .

**Proposition 5.2.** *Let  $\pi$  be known to the receiver who observes  $w \in W$ . Then the following properties hold.*

(i)  $\lim_{p \rightarrow \frac{m}{m+1}} f_w^\pi(t) = f_0(t).$

(ii) Let  $d^* := \min \{d \in \{0, \dots, n\} \mid \mu_0(\{t' \mid d(w, \pi(t')) = d\}) > 0\}$ .

(a) If  $d^* < d(w, \pi(t))$  then  $\lim_{p \rightarrow 0} f_w^\pi(t) = 0$ .

(b) If  $d^* = d(w, \pi(t))$  then  $\lim_{p \rightarrow 0} f_w^\pi(t) = f_0(t) \cdot \mu_0(\{t' \mid d(w, \pi(t')) = d(w, \pi(t))\})^{-1}$ .

(c) If  $d^* > d(w, \pi(t))$  the limit of the posterior belief for  $p \rightarrow 0$  is not defined.

---

<sup>6</sup>For binary channels, i.e.,  $m = 1$ , there is exactly one word with maximal distance  $n$  to a fixed  $v$ . The roles of the two letters (bits) simply switch and entropy is symmetric around the unformativeness bound  $p = \frac{1}{2}$ .

<sup>7</sup>We restrict to the intuitive case  $p < \frac{m}{m+1}$  where errors not abound. However, an analysis for  $p \rightarrow 1$  can be deduced with analogous results.

The first statement simply says that there is a smooth transition of the beliefs towards the common prior if the error channel gets uninformative. The second part deals with the behavior of the posteriors if the crossover probability goes to zero. In the presented continuous case, it is important to keep track of null sets but can still be intuitively explained. To start with, when receiving  $w$  the receiver determines the closest words to  $w$  that are sent the sender with positive probability. Let the according distance be  $d^*$  and let receiver contemplate about the state of the world being  $t$ .

If  $d^* < d(w, \pi(t))$ , the likelihood that the state is  $t$  goes to zero if  $p \rightarrow 0$  since there are events with positive probability in which words with  $d(w, v) = d^*$  are sent. If  $d^* = d(w, \pi(t))$ , then there is no event with positive probability in which words strictly closer to  $w$  are sent than  $\pi(t)$ . Taking the limit  $p \rightarrow 0$ , the receiver will discard states in which words even farther away from  $w$  than  $d^*$  are sent. Consequently, the receiver concludes that the true state  $t'$  fulfills  $d^* = d(w, \pi(t'))$ . The posterior is consequently given as stated in Proposition 5.2 (ii)(b). Remarkably, this is true even if  $w$  is not expected to be sent with positive probability and the Bayesian update for  $p = 0$  is undefined.<sup>8</sup> Our intuitive interpretation is as follows. A receiver who hears the word “orange” concludes that “orange” must have been the word sent even if they believe the error to be arbitrarily small. However, the case  $d^* > d(w, \pi(t))$  makes it impossible to apply Bayes rule as the receiver neither expects  $w$  to be sent with positive probability, nor do they believe that another word with distance  $d(w, \pi(t))$  has been sent.<sup>9</sup>

## 5.2 Entropy

Ever since the seminal work of Shannon (1948), the most important measure for the “amount of information” a communication channel can transport is given by (*Shannon*) *entropy*. In the following we determine the entropy of the  $q$ -ary symmetric channel of length  $n$  and crossover probability  $p$ . Intuitively, the more noise the harder it is for the receiver to confidently decode an observed message. We confirm the natural guess that the noise is maximal at the uninformativeness bound  $p = \frac{m}{m+1}$  and strictly monotonically increasing for both  $p \nearrow \frac{m}{m+1}$  and  $p \searrow \frac{m}{m+1}$ . While there is no noise for  $p = 0$ , information cannot become perfect for  $p = 1$  unless for a binary alphabet, i.e.,  $m = 1$ .

Formally, the entropy of a discrete probability measure  $P$  on a finite set  $X$  is

---

<sup>8</sup>Ortoleva (2012) provide a theoretical model to update given null events.

<sup>9</sup>If one considers the analogue model with a finite state space and the prior belief has full support, the case (c) disappears. Especially, the limit for  $p \rightarrow 0$  is always defined.

defined as

$$H(P) = - \sum_{x \in X} P(x) \cdot \log(P(x)), \quad (17)$$

with the convention  $0 = P(x) \cdot \log(P(x))$  if  $P(x) = 0$ . The base choice of the logarithm is a question of normalization and usually chosen to be  $\#X$  which is convenient for our setting. The value  $H(P)$  is interpreted as the average of the *information content*  $-\log(P(x))$  and attributed to describe how surprising the observation of an element  $x$  is given its probability  $P(x)$ . Some important properties of the entropy function  $H$  include non-negativity, strict concavity in the probability distribution  $P$  with the maximum being attained at the uniform distribution on  $X$  and symmetry in the order of the elements.

Turning to our metric-dependent error channel,  $\varepsilon(\cdot | v)$  defines a probability distributions on the set  $W$  for any fixed sent word  $v$  and any error probability  $p \in [0, 1]$  (that is suppressed by the notation). As the choice of  $v$  only leads to a permutation of the probabilities across  $W$  and by symmetry of  $H$ ,  $H(\varepsilon(\cdot | v))$  does not depend on  $v$  and  $H_\varepsilon(p) = H(\varepsilon(\cdot | v))$  is well-defined. The following proposition summarizes characterizing properties of  $H_\varepsilon(p)$ .

**Proposition 5.3.** *The entropy of the  $\#\mathcal{A}$ -ary symmetric error channel of length  $n$  with crossover probability  $p$  is*

$$H_\varepsilon(p) = -n \cdot (p \log(p) + (1 - p) \log(1 - p)) + n \cdot p \log(m). \quad (18)$$

*It is a concave in  $p \in [0, 1]$  with its unique maximum being attained at  $p = \frac{m}{m+1}$  with value  $\log(\#W)$ . Moreover,  $H_\varepsilon(0) = 0$  and  $H_\varepsilon(1) = n \log(m)$ .*

The proposition gives a quantitative view on the observations made in Remark 5.1. Choosing the base  $\#W$  for the logarithm, we can interpret  $H_\varepsilon(p)$  as the percentage of noise of the considered channel. For  $p = \frac{m}{m+1}$  entropy is maximal and equal to 1, which we interpret as each message being equally likely received. The error channel thus conveys no information. For  $p = 0$  entropy is zero, showing that there is no noise. Each piece of information is perfectly transmitted and can be correctly decoded. If  $p = 1$  entropy is  $\log(m)/\log(m+1)$ , telling us how much information is lost. For  $m = 1$ , i.e., a binary alphabet, this expression is again zero which makes sense since the roles of the letters simply swap, see also footnote 6 on page 13. In contrast, if  $m > 1$ , although information can be recovered, there will be noise left nevertheless. Received words cannot be unambiguously decoded. In between those extreme cases, due to concavity, we have a monotonic increase of entropy towards the uninformative bound  $p = \frac{m}{m+1}$  from both sides. For  $p \nearrow$  communication gets hindered more and more on  $[0, \frac{m}{m+1}]$ , while afterwards communication gets facilitated again.

Entropy is the quantitative measure for the amount of information that can be achieved through communication under noise. The following example illustrates the relation between an increasing noise in terms of entropy and the corresponding expected loss for a fixed communication device with a non-binary alphabet.

**Example 5.4.** Let  $T = [-\frac{1}{2}, \frac{1}{2}]$  with uniform  $\mu_0$ . Let  $W = \{L, M, R\}$  and consider  $\pi: T \rightarrow W, \pi([-\frac{1}{2}, 0]) = L, \pi((0, \frac{1}{2}]) = R$ . Figure 3 depicts the best responses of the receiver, the associated expected loss and the entropy of  $\varepsilon$ . At the uninformative bound  $p = \frac{2}{3}$  actions are pooling and expected loss and entropy are maximal. The receiver can perfectly decode  $L$  and  $R$  as  $R$  and  $L$  respectively if  $p = 1$ . However, for  $p = 1$ , they now also received the erroneous message  $M$  with positive probability, but cannot perfectly recover the originally sent message in that case. Consequently, communication under  $p = 1$  is worse than for  $p = 0$ . In numerical terms, we find that  $L(1)/L(2/3) = 0.625 \approx 0.63 \approx H_\varepsilon(1)$ , thus entropy roughly captures the expected loss.

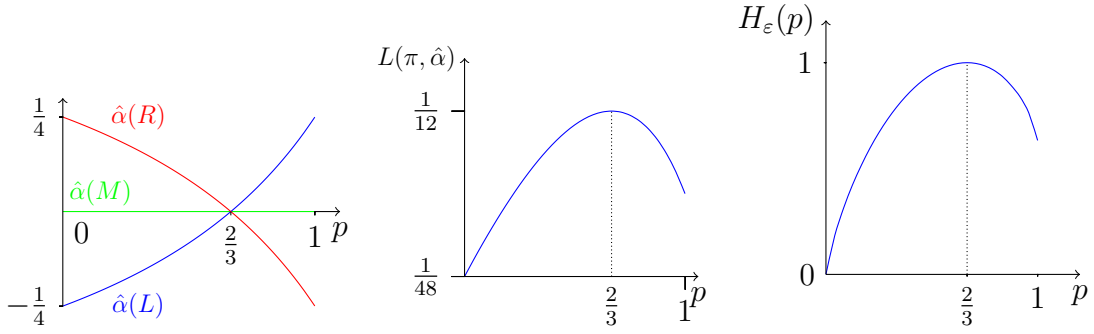


Figure 3: Optimal responses, expected loss and entropy for the communication device in Example 5.4. For three letters the uninformative bound is  $p = \frac{2}{3}$ . Actions coincide with the pooling action  $\alpha_{\text{pool}} = 0$ , expected loss and entropy are maximal. For  $p = 1$ , optimal receiver responses to the messages used with positive probability can be recovered up to re-labeling. Observing  $L$  can be perfectly decoded to have been the original message  $R$ . However,  $M$  will now be received with positive probability, leading to a loss in information in comparison to  $p = 0$ .

## 6 The grammar of efficient languages

The most prominent loss function in the study of communication games is the quadratic loss. It is used as the prime example in the field of cheap talk games ever



since the seminal of Crawford & Sobel (1982), allowing for analytical tractability. We analyze the effects of noise on the geometric structure of languages under a quadratic loss. In the field of linguistics, the set of structural rules of a language is more prominently known as its *grammar*. Seeking for efficiency, we replicate the following grammatical patterns of the agents' best responses under noise. A speaker distributes their words in a way that maximizes a certain spread of their induced interpretations. The range of meanings that a word represents is convex. Words must not be vague. Noise necessitates to reserve distant words to states that result in a huge loss if confused. The more noise the clearer the speaker will stress words that should not be confused, reducing or even forsaking the use of words that are easily confused.

For now, let  $\varepsilon$  be an arbitrary noisy channel. Throughout this section we consider a quadratic loss, i.e, we let  $\|\cdot\| = \|\cdot\|_2$  be the Euclidean norm induced by the scalar product  $\langle \cdot, \cdot \rangle$  and  $\ell(x) = x^2$ . Under this assumption, the best response of the receiver can be written down explicitly.

**Lemma 6.1.** *Having any (posterior) belief  $\mu$ , with continuous density  $f > 0$ , the receiver's unique best interpretation is given by*

$$\hat{\alpha}(\mu) = \mathbb{E}_\mu[t]. \quad (19)$$

*Given a communication device  $\pi$  and the induced best responses, the expected loss is*

$$L(\pi, \hat{\alpha}) = \mathbb{E}_{\lambda^\pi}[\mathbb{E}_{\mu_w^\pi}[\|t - \mathbb{E}_{\mu_w^\pi}[t']\|_2^2]] = \mathbb{E}_{\mu_0}[\|t\|_2^2] - \mathbb{E}_{\lambda^\pi}[\|\hat{\alpha}(w)\|_2^2]. \quad (20)$$

*Furthermore, the receiver on average plays the pooling action*

$$\mathbb{E}_{\lambda^\pi}[\hat{\alpha}(w)] = \alpha_{\text{pool}}. \quad (21)$$

Lemma 6.1 proves useful for pinning down analytical properties of efficient languages. As a first application we show that profitable communication is always possible, as long as the error channel is not uninformative.

**Proposition 6.2.** *If  $\varepsilon$  is not uninformative there is  $(\pi, \alpha)$  with  $L(\pi, \alpha) < L_{\text{pool}}$ .*

In the following subsections, we want to elicit grammatical, i.e., geometric and structural, properties of efficient communication.

## 6.1 Convexity of meaning

Employing a quadratic loss provides a new tractable view on what makes up an efficient language. By expression (20), the expected loss is the difference between a term that depends only on the state space and its measure and a weighted sum of

the square norms of the induced interpretations. Minimizing the loss of a language is thus equivalent to maximizing the spread of induced squared interpretations. The expected loss is the difference of the pooling loss and the (non-negative) difference of this spread and the square of the pooling action, precisely

$$L(\pi, \hat{\alpha}) = L_{\text{pool}} - (\mathbb{E}_{\lambda^\pi}[\|\hat{\alpha}(\mathbf{w})\|_2^2] - \|\alpha_{\text{pool}}\|_2^2). \quad (22)$$

Since  $\alpha_{\text{pool}} = \sum_{\mathbf{w}} \lambda^\pi(\mathbf{w}) \hat{\alpha}(\mathbf{w})$ , an efficient language is a decomposition of  $\alpha_{\text{pool}}$  that maximizes the weighted sum of induced square norm interpretations among all communication devices. Summarized as a structural consequence, efficient languages are *as separating as possible*.

It is worthwhile to discuss this insight by revisiting Example 3.4 and Example 4.1. Even though the communication device in Example 3.4 results in different posterior beliefs, they do not accomplish a spread of induced actions. In contrast, the best responses in Example 4.1 achieve the maximal spread that the noisy channel allows. An apparent difference of the two communication devices is that the cells in the efficient example are convex, while this is not always the case for the inferior one. Indeed, non-convexity of cells is an indicator that the sender does not play a best response to the receivers actions under a quadratic loss, as the next proposition proves.

**Proposition 6.3.** *The set of states for which sending  $v \in W$  is a (the unique) best reply given  $\alpha$  is a closed (open) convex set.*

Proposition 6.3 implies further structural rules on efficient languages. If the sender plays a best reply to the receiver’s interpretation map  $\alpha$ , an induced cell corresponds to the indicative meaning of the corresponding word. Each cell thus forms a category in the sense of Lewis (1969). Our result thus generalizes the linguistic idea that *(simple) words have convex categories* to communication under noise, cf. Gärdenfors (2004), Jäger (2007) and Jäger et al. (2011).

Starting at which temperature does one say that it is “hot”? Certainly, most people agree that 100°F (38°C) is hot, but there is hardly a precise threshold separating “hot” from “not hot” temperatures. Words like “hot”, “tall” or “many” are *vague*, they lack a precise definition, cf. Sorensen (2023) and Lipman (2009).<sup>10</sup> Under a quadratic loss and the metric-dependent noise from Section 5, the sender does not use vague words. If they say “hot” this has a crisp meaning. Formally, if two words  $v, v'$  are no *synonyms*, i.e.,  $\alpha(v) \neq \alpha(v')$ , sending one of the two words is almost surely strictly preferred. As a result, the induced cells have sharp boundaries.

---

<sup>10</sup>Vagueness is often associated with the so-called *sorites paradox*: When makes taking away grains of sand from a heap stop us referring to the remaining sand as a heap?

**Proposition 6.4.** *Let  $\varepsilon$  be the  $q$ -ary symmetric channel of length  $n$  with crossover probability  $p$  and  $\alpha$  an interpretation map of the receiver.*

*If  $\alpha(v) \neq \alpha(v')$ , the set of states for which the sender is indifferent between sending  $v$  and  $v'$  is a null set for all but at most  $n$  values of  $p \in [0, 1]$ .*

*If in addition  $p = 0$ , the sender is almost surely never indifferent, while for the unformativeness bound  $p = \frac{m}{m+1}$  they always are.*

## 6.2 The shape of cells and their labeling

In this section, we study an extensive example of a two-dimensional state space with words of length two. Four classes of noise equilibria are given, each characterized by two different properties. Cells shapes can be either quadratic or triangular and the labeling can respect far away states by using distinct words or not. Our results suggest that efficient communication under noise needs to be clear in distinguishing distinct states by words that are not easily confused and use a cell structure minimizing the length of the boundary between cells.

Let  $T = [-\frac{1}{2}, \frac{1}{2}]^2$  be endowed with the uniform distribution  $\mu_0 \sim \mathcal{U}(T)$ , especially  $\alpha_{\text{pool}} = (0, 0)$ . The sender uses an alphabet with two letters and can send words of length two  $W = \{A, B\}$ . Communication is frustrated by a binary symmetric channel of length two and crossover probability  $p$ . Assume  $p < \frac{1}{2}$  so that the error function is not uninformative and it is less likely to confound words that are farther away from one another. Consider the communication devices depicted in Figure 4 which constitute in fact noise equilibria. A detailed derivation of all formulae in Table 1 and explanation of the claims is given in appendix Section 9.3.

Languages 1a and 1b employ quadratic tessellations, while 2a and 2b use triangular shaped cells. Note that neighboring cells in 1a and 2a are associated with words that are close to one another (Hamming distance 1) and words that are less easily confused are used for cells that are far away from one another. In contrast, each cell in 1b and 2b has a neighbor the assigned word of which has distance 1 and one with distance 2.

Figure 4 plots the loci of best responses of the receiver are drawn from  $p = 0$  to the unformativeness bound  $p = \frac{1}{2}$ . The precise formulae are summarized in Table 1. As is clear from Proposition 5.2, the receiver's interpretations start at the center of each cell and continuously move to the pooling action. Interestingly, the convergence to the pooling action differs significantly. While the loci in 1a and 2a are straight lines, the ones for 1b and 2b bend towards the cell that uses the word with distance 1 to the own one. The intuition for this is straightforward: Look w.l.o.g. at  $\hat{\alpha}(AA)$ : Within the type-a languages, the cells of  $AB$  and  $BA$  absorb the same amount of mistakes from and towards  $AA$ , linearly in  $p$ . The word  $BB$  does so quadratically but also point-symmetrically w.r.t. the pooling

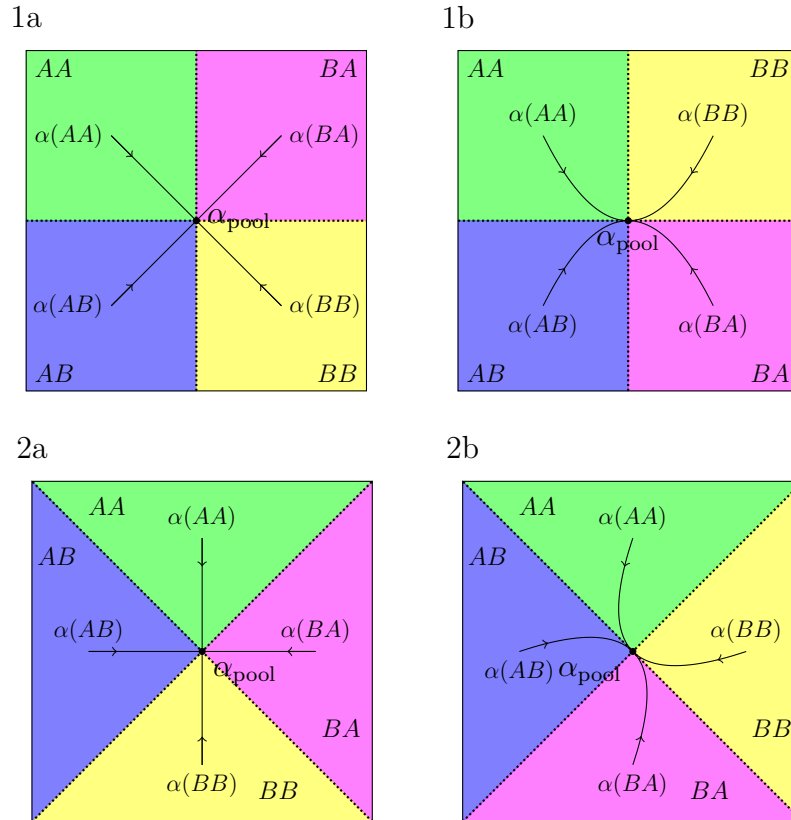


Figure 4: Four communication devices on  $T$ . The black lines depict the loci of the receiver's best responses for  $p \in [0, \frac{1}{2}]$ , moving from the center of the resp. cells to the pooling action. For every  $p$  they constitute a noise equilibrium respectively. While the best responses in 1a and 2a move to  $\alpha_{\text{pool}}$  in a straight line, they are pulled closer to the cell of their neighbor with closer words in 1b and 2b.

Case	$\alpha(AA)$	$\alpha(AB)$	$L(\pi, \alpha)$
1a	$\frac{1}{4}(-1 + 2p, 1 - 2p)$	$\frac{1}{4}(-1 + 2p, -1 + 2p)$	$\frac{1}{6} - \frac{1}{8}(1 - 2p)^2$
1b	$\frac{1}{4}(-1 + 2p, 1 - 4p + 4p^2)$	$\frac{1}{4}(-1 + 2p, -1 + 4p - 4p^2)$	$\frac{1}{6} - \frac{1}{8}(1 - 2p)^2(1 - 2p + 2p^2)$
2a	$\frac{1}{3}(0, -1 + 2p)$	$\frac{1}{3}(-1 + 2p, 0)$	$\frac{1}{6} - \frac{1}{9}(1 - 2p)^2$
2b	$\frac{1}{3}(-p + 2p^2, 1 - 3p + 2p^2)$	$\frac{1}{3}(-1 + 3p - 2p^2, p - 2p^2)$	$\frac{1}{6} - \frac{1}{9}(1 - 2p)^2(1 - 2p + 2p^2)$

Table 1: Explicit formulae for the optimal actions and the expected loss of the languages from Figure 4. The remaining formulae can be found in the appendix Section 9.3 or by symmetry.

action. Increasing  $p$  thus shifts the interpretation on a straight line towards  $\alpha_{\text{pool}}$ . In contrast, for the type-b languages, the locus of  $\hat{\alpha}(AA)$  is again pulled linearly by the cells of  $AB$  and  $BA$ , while quadratically by the one of  $BB$ , which is the weaker force since  $p < \frac{1}{2} < 1$ . As it is less likely to confuse  $AA$  and  $BB$ , the border between their cells is less permeable for mistakes.

Let us take a look at their performance and focus on the nomination of the cells first. Each type-a language outperforms their type-b ones for each noise level, see Figure 5. The labeling of the cells thus is an important determinant of the efficiency of a language. In order to reduce the harm from miscommunication, the sender should use similar words to describe similar states and use distinct words to describe states that should not be mistaken.

Turning towards the structure of the cells employed, we find that the quadratic ones outperform the triangular shaped ones, i.e., type-1 languages have a lower expected loss than type-2 ones. The squares provide a more compact structure and have less points near and on the indifference levels than the triangle shapes.<sup>11</sup> As a result the interpretations for type-1 languages stay farther away from the pooling action than their resp. type-2 peers, leading to a more separable decomposition of the prior belief.

Interestingly, language 2a results in a lower loss than 1b for a crossover probability exceeding  $p = \frac{1}{2} - \sqrt{2} \approx 6\%$ , suggesting that the labeling of the cells can be more important than the choice of stable cell structures. However, when analyzing language formation we find both type-1 languages to be stable outcomes of evolution whereas type-2 languages are unstable, cf. Section 7.

<sup>11</sup>The total interior border length of type-1 languages is  $1 + 1 = 2$ , while for type-2 languages this is  $\sqrt{2} + \sqrt{2} > 2.8$ .

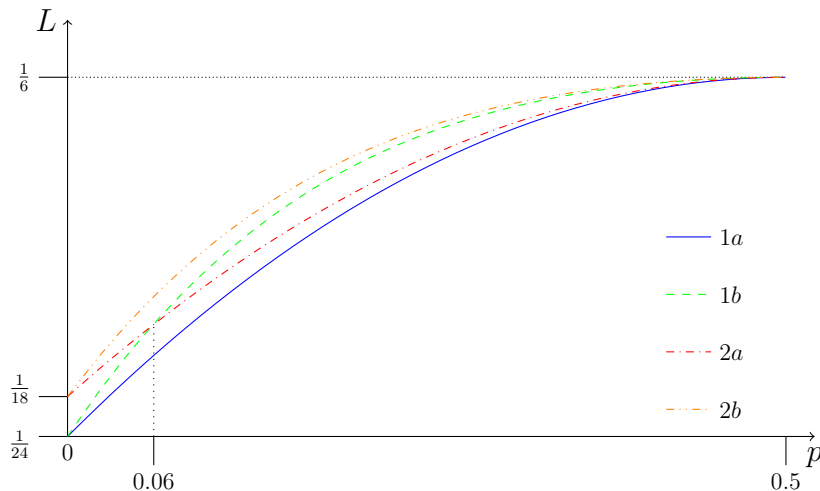


Figure 5: Expected loss of the four languages from Figure 4 for different crossover probabilities  $p$ . For  $p = 0$  the assignments of words to cells is irrelevant and 1a and 1b have a lower loss than 2a and 2b. At the uninformative bound  $p = \frac{1}{2}$  the pooling loss realizes for all languages. For all noise levels, 1a (2a) has a lower loss than 1b (2b). Interestingly, for  $p > 6\%$  2a has a lower loss than 1b.

### 6.3 Efficient languages with four words

Only a few words are really necessary to get the main idea of a sentence. The remaining words are either decorative, concerned with details or redundant. Emphasizing key words becomes more important the more likely communication is noisy. The following example indicates that the sender has an incentive to stress the words that describe the states which are farthest away from one another and would thus lead to a high loss if confused. This property is achieved by enlarging the cells on the outskirts. As a consequence, the cells close to the pooling action shrink. In a sense, the sender is willing to give up precision over the whole state space to ensure that extreme states are not mistaken.

Let  $T = [-\frac{1}{2}, \frac{1}{2}]$  be endowed with the uniform distribution. The word space is given by  $W = \{A, B\}^2$  and the noise is given by a binary symmetric channel of length two with crossover probability  $p$ . For different values of  $p$ , we can categorize efficient languages. Some of them are depicted in Figure 6. Again, note that extreme states are articulated by sending words that are not easily mistaken. However, in contrast to the example in Section 6.2, the low dimension of the current state does not allow to sustain its structure for an increasing noise.

Imagine a speaker describes their friend the height of a person they just have met. They use words in the scheme “very tiny” ( $AA$ ), “not too tiny” ( $AB$ ), “not

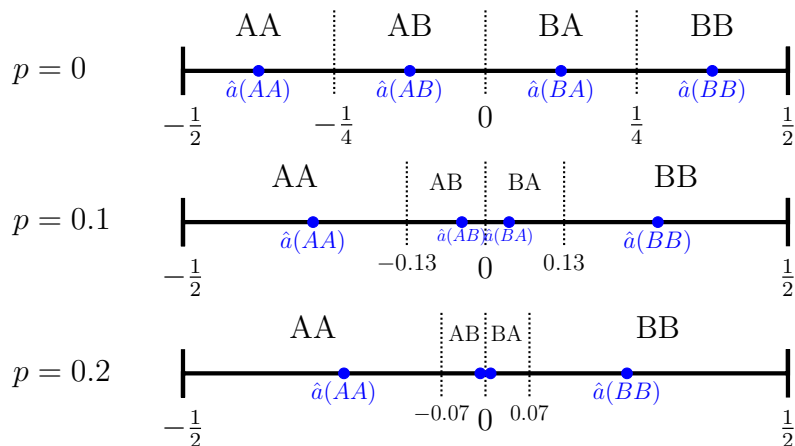


Figure 6: Efficient languages in the setting of Section 6.3 for different crossover probabilities  $p$ . Words of maximal distance are used to articulate the boundaries and their cells enlarge with increasing error.

too tall” ( $BA$ ) and “very tall” ( $BB$ ) to set their height in relation to the average height (indicated by the state 0). Note that we can interpret the letter  $A$  as indicating “tiny” and  $B$  as “tall”. If a letter is repeatedly used, its indicative meaning is stressed, if both are used, its meaning is attenuated by the second letter.<sup>12</sup> Perhaps in contrast to the reader’s intuition, our rational sender does not like vagueness and thus always has precise definition of the ranges of height addressed when using a word. If there is no noise, the speaker efficiently communicates the height by equally splitting the state space (if heights are equally distributed). If noise increases, the speaker puts more emphasis describing the extreme cases than on average sized persons. In the limit, the speaker will not use words to describe an average person at all and the meaning of the opposite words  $AA$  and  $BB$  becomes simply “tiny” and “tall”.

## 7 Evolution

Natural languages have formed over time by the laws of evolution and are perpetually changing. The development from indistinct shrieks to elaborate speech has continuously improved coordination among human beings. Evolutionary game theory has proved useful in deriving qualitative properties of evolution in biology by means of simple models, cf. Smith & Price (1973), Hofbauer & Sigmund (1998). Its main idea is as follows. The population of a species at each point in time is

<sup>12</sup>Of course, exchanging the words  $AB$  and  $BA$  in the example does not qualitatively change the example.

a distribution over groups of different characteristics. Individuals are randomly paired and play a game. Their expected excess payoff is proportional to their fitness. If a group performs better than the average their subsequent share in the population increases. In the following, we apply evolutionary game theory to communication under noise. As it turns out, individuals can learn how to use communication to improve coordination even in the presence of noise.

Formally, we describe an individual by endowing them with both a communication device  $\pi$  out of the set of all measurable communication devices  $\Sigma$  and an interpretation map  $\alpha \in T^{\#W}$  for a fixed word space  $W$ . The strategy space  $\Sigma \times T^{\#W}$  is a complicated space, even more so when considering populations, i.e., probability distributions over strategies. The technical foundation has fortunately been established for certain dynamics and extends to our setting. These include the replicator, cf. Oechssler & Riedel (2001), Cressman et al. (2006), payoff monotone, cf. Heifetz et al. (2007), and Brown-von-Neumann-Nash dynamics, cf. Hofbauer et al. (2009).

We proceed along the lines of Jäger et al. (2011), considering a symmetric version of the cheap talk game studied so far. When two individuals meet, a fair coin toss decides about their roles of sender or receiver. An individual using  $(\pi, \alpha)$  and meeting another one using  $(\pi', \alpha')$  amounts to an expected loss of

$$\Lambda((\pi, \alpha), (\pi', \alpha')) = \frac{1}{2}L(\pi, \alpha') + \frac{1}{2}L(\pi', \alpha) \quad (23)$$

for both. Describing a population of individuals by a probability distribution  $P$  on  $\Gamma := \Sigma \times T^{\#W}$  the expected loss is generalized to

$$\Lambda(P, Q) := \int_{\Gamma} \int_{\Gamma} \Lambda((\pi, \alpha), (\pi', \alpha')) P(d\pi, d\alpha) Q(d\pi', d\alpha'). \quad (24)$$

Individuals are able to learn efficient communication under noise.

**Proposition 7.1.** *The following assertions hold.*

- (i) *The symmetrized loss function is a Lyapunov function for the replicator, regular and payoff monotone and the Brown-von-Neumann-Nash dynamics.*
- (ii) *Locally optimal languages are Lyapunov stable w.r.t. the replicator, regular and payoff monotone and Brown-von-Neumann-Nash dynamics.*

Figure 7 gives a tractable numerical illustration of evolution by means of the best reply dynamics in the setting of Section 6.2. In this case, a sender and a receiver meet every day and play our communication game. In the beginning, both had a random strategy, but after each encounter they learn the strategy of their peer and play a best response to that at the next time. As the figure shows, their



language quickly converges to a noise equilibrium. The depicted equilibrium is a local optimum that improves upon the pooling loss, but is not efficient, compare Section 6.2. Indeed, numerical simulations suggest that the languages 1a and 1b are stable and even attractors of the best reply dynamics, while languages 2a and 2b prove to be unstable.

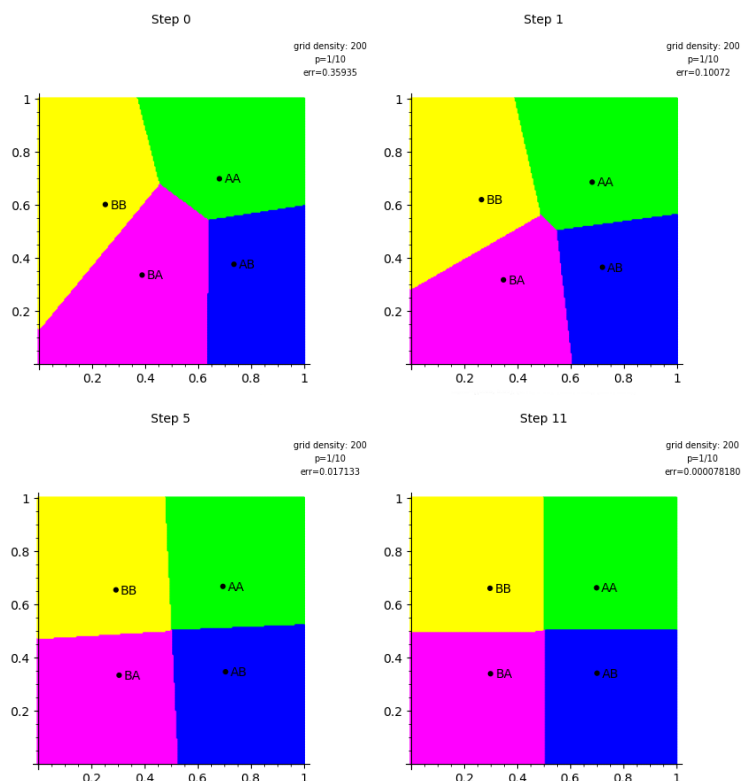


Figure 7: Numerical simulation of the best reply dynamics in the setting of Section 6.2 (shifted to the unit square). Starting with random interpretations, the agents take turn in playing a best reply to the previous strategy of their peer. Convergence to equilibrium is measured by the norm distance between two subsequent interpretation maps. Agents quickly learn one of the noise equilibria 1a or 1b from Figure 4 which, in contrast to 2a and 2b, appear to be stable and attractors.

## 8 Conclusion

Within our daily routine, errors in our communication are ubiquitous. Despite this inhibiting factor, humans have learned to communicate efficiently through the cause of their existence. This chapter formalizes noisy communication as a

cheap talk game of common interest and studies its structural properties, i.e., the grammar of communication. Efficient languages exist and can be learned by evolutionary dynamics. The sender optimally induces receiver actions, the spread of which should be as maximal as possible. As in non-noisy cases, the set of states which are referred to by a single word is convex and has sharp boundaries to reduce inefficient vagueness. If noise is present, examples illustrate that the sender can reduce the expected loss by using distant words to describe distant states. That way, they minimize the relatively high loss that results from confusing states that are very distinct. The more noisy communication is, the more important this becomes. If the dimension of the state space does not allow for errors to balance out, the communication device in place must adjust for higher errors by emphasizing extreme states at the cost of describing average states.

## 9 Appendix

### 9.1 Finite state space

While assuming a rich state space is compelling if we think of describing, say, colors, we can also think of scenarios in which states are finite, for instance telling somebody to go left or right. In fact, one can extend the presented framework by allowing for point masses of  $\mu_0$  on the state space. If  $\mu_0$  only has point masses, we can thus model a finite state space. In fact, many of the results we have seen carry over nicely, for instance the uniqueness of sender's best replies and the existence of efficient equilibria. The classical framework of Shannon (1948) indeed considered *finite sources* which is the natural assumption in information theory where data is discrete. One of its most important fields for all of our digital data transmission is *coding theory*, cf. Roth (2006). A code is a subset  $C \subset \{0, 1\}^n$  of binary sequences of a fixed length  $n$ . Before the transmission starts, agents declare  $C$  and agree that the sender only sends (code-)words  $c \in C$ . At first it might seem surprising that codes typically are a proper subset of all possible words  $\{0, 1\}^n$ , but this is no coincidence. Including redundancy in communication enables the receiver to detect and probably even correct errors resulting from noise in the transmission channel. In the following, we illustrate the economic reasoning behind this using a finite state space with as many words as there are types.

Let  $T = \{-2, -1, 1, 2\}$  be a discrete state space with uniform prior  $\mu_0$ . The message space is  $W = \{A, B\}^2$ . Communication is noisy and the error channel is given by a binary symmetric channel of length two with crossover probability  $p$ . Let  $\text{supp}(\pi) = \pi(T)$  be the support of a communication device  $\pi$ . For the sake of the argument, assume the receiver can still play 'hedging' actions, i.e., any point on the convex hull of  $T$ . The expected loss is quadratic. One can argue that,

up to isometry, there are unique best communication devices for each support of cardinality 2, 3, 4 if the receiver plays their best response. We depict these in Table 2.

com.dev.	$\pi(-2)$	$\pi(-1)$	$\pi(1)$	$\pi(2)$
$\pi_2$	AA	AA	BB	BB
$\pi_3$	AA	AA	AB	BB
$\pi_4$	AA	AB	BA	BB

Table 2: The set of communication devices that are optimal given the cardinality of their support if the receiver plays a best response. Here  $W = \{A, B\}^2$  and  $T = \{-2, -1, 1, 2\}$ .

A straightforward calculation determines the expected losses given each of the communication devices from Table 2, see Figure 8 for a plot.

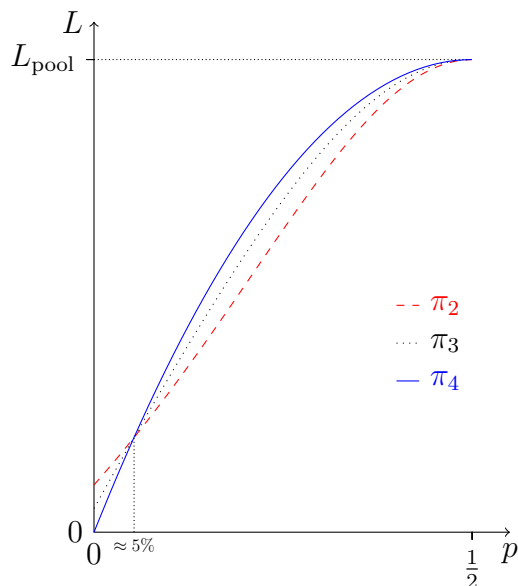


Figure 8: Expected loss of the communication devices in Table 2 under the receiver's best responses. The communication device using four words is the best for small values of  $p$ , while the one using only two words results in the lowest loss starting at  $p \approx 5\%$ . The agents are never better off using only three words.

The example provides an economic reasoning why codes in information theory feature redundancy. If there was no noise, the sender can only perfectly reveal the state to the receiver if they employ as many messages as there are states. However, if there is noise, all the receiver's actions get pulled towards the pooling action.

This effect is more pronounced the more words are used. As a result, there will be a threshold when using fewer words results in a lower expected loss, see Figure 8, justifying redundancy of a code.

## 9.2 Proofs & calculations

*Proof of Lemma 3.1.* The non-emptiness of the best reply set is given by continuity of the integrand and compactness of  $T$ . We will now prove that the function  $T \rightarrow \mathbb{R}$ ,  $s \mapsto \mathbb{E}_\mu[\ell(\|t - s\|)]$  is strictly convex, implying uniqueness of the minimizer. To this end, assume there are two distinct minimizers  $s_1, s_2$  and let  $\lambda \in (0, 1)$ . By the triangle inequality and convexity of  $\ell$  we have

$$\begin{aligned} & \mathbb{E}_\mu[\ell(\|t - (\lambda s_1 + (1 - \lambda)s_2)\|)] \\ & \leq \mathbb{E}_\mu[\ell(\lambda\|t - s_1\| + (1 - \lambda)\|t - s_2\|)] \end{aligned} \quad (25)$$

$$\begin{aligned} & < \mathbb{E}_\mu[\lambda\ell(\|t - s_1\|)] + \mathbb{E}_\mu[(1 - \lambda)\ell(\|t - s_2\|)] \\ & = \mathbb{E}_\mu[\ell(\|t - \hat{s}\|)]. \end{aligned} \quad (26)$$

Since the strictness of the inequality in (26) is not too obvious, we give some more details. By strict convexity of  $\ell$ , we have a strict inequality within the integrand for all  $t$  with  $\|t - s_1\| \neq \|t - s_2\|$ . Consequently, it suffices to prove that the set  $M := \{t \mid \|t - s_1\| \neq \|t - s_2\|\}$  has positive  $\mu$ -mass. By positive definiteness of the norm we have  $s_1 \in M$ . Since norms are continuous (even Lipschitz),  $M$  contains an open environment (in  $T$ ) of  $s_1$ . Since  $\mu$  is absolutely continuous w.r.t. the Lebesgue measure we find  $\mu(M) > 0$  to conclude the proof.  $\square$

*Proof of Proposition 3.2.* We have

$$L(\pi, \hat{\alpha}) = \mathbb{E}_{\lambda^\pi} [\mathbb{E}_{\mu_w^\pi} [\ell(\|t - \hat{\alpha}(w)\|)]] \quad (27)$$

$$\leq \mathbb{E}_{\lambda^\pi} [\mathbb{E}_{\mu_w^\pi} [\ell(\|t - \alpha_{\text{pool}}\|)]] \quad (28)$$

$$= \mathbb{E}_{\mu_0} [\ell(\|t - \alpha_{\text{pool}}\|)] = L_{\text{pool}}, \quad (29)$$

using the defining property for the inequality and Bayes-Plausibility to condense the expectations. The inequality is strict if and only if there is a word  $w$  with  $\lambda^\pi(w) > 0$  and  $\hat{\alpha}(w) \neq \alpha_{\text{pool}}$  as the resp. minimizers are unique.  $\square$

*Proof of Corollary 3.3.* If  $\pi$  is constant, say  $\pi \equiv v$ , then  $\varepsilon(w \mid \pi(t)) = \varepsilon(w \mid v)$  is constant and equal to  $\lambda^\pi(w)$  for any  $w \in W$ . This implies  $f_w^\pi = f_0$  and  $\mu_w^\pi = \mu_0$ .

If  $v \mapsto \varepsilon(\cdot \mid v)$  is constant, we have  $K_w := \varepsilon(w \mid \pi(t))$  is independent of  $t$  for any communication device  $\pi$ . Consequently, also  $\lambda^\pi(w) = K_w$ , implying  $f_w^\pi(t) = f_0(t)$  and thus  $\mu_w^\pi = \mu_0$ .  $\square$

*Proof of Theorem 4.2.* We follow the proof of Lemma 1 in Jäger et al. (2011) which can be adjusted to incorporate any noisy channel.

Consider any pure strategy  $\alpha: W \rightarrow T$  of the receiver. Then, a type  $t$ -sender may optimally send any word  $v$  out of

$$\arg \min_{v' \in W} \sum_{w \in W} \varepsilon(w | v') \cdot \ell(\|t - \alpha(w)\|). \quad (30)$$

Note that the set of minimizers is non-empty as  $W$  is finite. Now, fix any strict ordering  $\leq_W$  on  $W$  and define a partition of  $T$  by setting

$$C_v^\alpha := \left\{ t \mid v \text{ is smallest w.r.t. } \leq_W \text{ in } \arg \min_{v' \in W} \sum_{w \in W} \varepsilon(w | v') \cdot \ell(\|t - \alpha(w)\|) \right\} \quad (31)$$

for each  $v \in W$ .  $C_v^\alpha$  is (Lebesgue-)measurable as all involved functions are continuous in  $t$  and it is the set difference of a closed set from a finite union of closed sets. To see this, start by collecting all  $t$  for which  $v$  is a minimizer, which is a closed set. Now, for all  $v' \leq_W v$  that are also minimizers take away the indifference sets, which are themselves closed, to obtain  $C_v^\alpha$ .

We define a measurable function  $\pi: T \rightarrow W$  by  $\pi(t) = v \iff t \in C_v^\alpha$  which represents one possible best reply of the sender.

Using any such choice we can internalize a sender's best reply and re-write the joint loss minimization as a function depending only on  $\alpha$ , namely

$$\min_{\alpha} \int_T \min_v \left\{ \sum_{w \in W} \varepsilon(w | v) \cdot \ell(\|t - \alpha(w)\|) \right\} \mu_0(dt). \quad (32)$$

We identify any strategy  $\alpha: W \rightarrow T$  with a point in  $T^N$ ,  $N := \#W$ . Thus, by Lebesgue's dominant convergence theorem, it suffices to prove continuity of the integrand in  $\alpha$  for any fixed  $t$ . But this is obvious as the pointwise minimum of finitely many continuous functions is again continuous.  $\square$

*Proof of Remark 5.1.* The case  $p = 0$  is clear. For  $0 < \tilde{p} < 1$  we have that

$$\varepsilon(w | v) > \varepsilon(w' | v) \iff \tilde{p}^{d(w,v)} > \tilde{p}^{d(w',v)} \iff d(w, v) < d(w', v). \quad (33)$$

The other cases follow similarly.  $\square$

*Proof of Proposition 5.2.* (i) Follows immediately from continuity of the integrand in  $\tilde{p}$  and Lebesgue's theorem.

(ii) We split the integral in three parts by disassembling the type space  $T$  into the three stets defined by  $\{t' \mid d(w, \pi(t')) \sim^* d(w, \pi(t))\}$  for  $\sim^* \in \{<, =, >\}$ .

(a) The set  $\{t' \mid d(w, \pi(t')) < d(w, \pi(t))\}$  has positive probability and the negative exponent  $d(w, \pi(t')) - d(w, \pi(t))$  will let the integral go to infinity as  $p \rightarrow 0$ .

(b) The set  $\{t' \mid d(w, \pi(t')) < d(w, \pi(t))\}$  has probability zero and can be neglected. For  $p \rightarrow 0$ , the integral over  $\{t' \mid d(w, \pi(t')) > d(w, \pi(t))\}$  will vanish as the exponent of  $\tilde{p}$  is strictly positive. What is left of the overall integral is  $\int_{\{t' \mid d(w, \pi(t')) = d(w, \pi(t))\}} \mu_0(dt') = \mu_0(\{t' \mid d(w, \pi(t')) = d(w, \pi(t))\})$  which is strictly positive by assumption.

(c) Ignoring the integral over null sets, the limit  $p \rightarrow 0$  makes the integral go to 0 making the limit meaningless.

□

*Proof of Proposition 5.3.* We start by calculating the entropy

$$\begin{aligned} & H(\varepsilon(\cdot \mid \mathbf{v})) \\ &= - \sum_{w \in W} \varepsilon(w \mid \mathbf{v}) \cdot \log(\varepsilon(w \mid \mathbf{v})) \\ &= - \sum_{w \in W} (1-p)^{n-d(w, \mathbf{v})} \cdot \left(\frac{p}{m}\right)^{d(w, \mathbf{v})} \cdot \log \left( (1-p)^{n-d(w, \mathbf{v})} \cdot \left(\frac{p}{m}\right)^{d(w, \mathbf{v})} \right) \end{aligned} \quad (34)$$

$$= - \sum_{d=0}^n \binom{n}{d} \cdot m^d \cdot (1-p)^{n-d} \cdot \left(\frac{p}{m}\right)^d \cdot \log \left( (1-p)^{n-d} \cdot \left(\frac{p}{m}\right)^d \right) \quad (35)$$

$$= - \sum_{d=0}^n \binom{n}{d} \cdot (1-p)^{n-d} \cdot p^d \cdot \log \left( (1-p)^{n-d} \cdot p^d \right) \quad (36)$$

$$+ \log(m) \cdot \sum_{d=0}^n \binom{n}{d} \cdot d \cdot (1-p)^{n-d} \cdot p^d \quad (37)$$

$$= - \sum_{d=0}^n \binom{n}{d} \cdot (1-p)^{n-d} \cdot p^d \cdot \log \left( (1-p)^{n-d} \cdot p^d \right) + np \log(m) \quad (38)$$

$$= - np \cdot \log(p) - n(1-p) \cdot \log(1-p) + np \log(m) \quad (39)$$

$$= n \cdot (H((p, 1-p)) + p \log(m)). \quad (40)$$

During the calculation we used the functional equation of the logarithm and formulae occurring often when dealing with binomial distributions, e.g., its mean. This function is concave in  $p$  since the entropy  $H((p, 1-p)) = -p \log(p) - (1-p) \log(1-p)$

over a binary source with probability  $p$  is. The maximum is attained for the uniform distribution, i.e., if  $p = \frac{m}{m+1}$  by Remark 5.1 and yields  $H(\mathcal{U}(W)) = \log(\#W)$ . The other assertions follow readily from the calculated expression.  $\square$

*Formulae for Example 5.4.* For  $p \in [0, 1]$  the optimal response  $\hat{\alpha}$  of the receiver is given by

$$\hat{\alpha}(L) = \frac{-2 + 3p}{8 - 4p}, \quad \hat{\alpha}(M) = 0, \quad \hat{\alpha}(R) = -\hat{\alpha}(L) \quad (41)$$

and the expected loss given  $p$  can be calculated to be

$$L(\pi, \hat{\alpha})(p) = \frac{1}{12} - 2^{-5} \cdot \frac{(-2 + 3p)^2}{2 - p}. \quad (42)$$

$\square$

*Proof of Lemma 6.1.* The receiver's minimization problem reads as

$$\min_{\alpha \in T} \mathbb{E}_\mu[\|t - \alpha\|_2^2] = \int_T \sum_{k=1}^L (t_k - \alpha_k)^2 \mu(dt). \quad (43)$$

Using the Leibniz rule we check the first and second order conditions for each  $k$  and obtain the unique local and global minimum by choosing

$$\hat{\alpha}_k(\mu) = \int_T t_k \mu(dt) = \mathbb{E}_\mu[t_k]. \quad (44)$$

Plugging  $\hat{\alpha}(\mu) = \mathbb{E}_\mu[t]$  back into the expected loss and using the scalar product  $\langle \cdot, \cdot \rangle$  we get

$$\mathbb{E}_\mu[\|t - \mathbb{E}_\mu[t]\|_2^2] = \mathbb{E}_\mu[\|t\|_2^2] - \|\mathbb{E}_\mu[t]\|_2^2, \quad (45)$$

which is the trace norm of the variance matrix, i.e., the sum of the variances over each dimension. If a communication device  $\pi$  is given the expected loss is

$$\mathbb{E}_{\lambda^\pi}[\mathbb{E}_{\mu_w^\pi}[\|t - \mathbb{E}_{\mu_w^\pi}[t']\|_2^2]] = \mathbb{E}_{\lambda^\pi}[\mathbb{E}_{\mu_w^\pi}[\|t\|_2^2] - \|\mathbb{E}_{\mu_w^\pi}[t]\|_2^2] \quad (46)$$

$$= \mathbb{E}_{\mu_0}[\|t\|_2^2] - \mathbb{E}_{\lambda^\pi}[\|\mathbb{E}_{\mu_w^\pi}[t]\|_2^2] \quad (47)$$

where Bayes-Plausibility has been used. Finally, applying Bayes-Plausibility one more time we observe

$$\mathbb{E}_{\lambda^\pi}[\hat{\alpha}(w)] = \mathbb{E}_{\lambda^\pi}[\mathbb{E}_{\mu_w^\pi}[t]] = \mathbb{E}_{\mu_0}[t] = \hat{\alpha}(\mu_0) = \alpha_{\text{pool}}. \quad (48)$$

$\square$

*Proof of Proposition 6.2.* Without loss of generality, we assume  $\alpha_{\text{pool}} = 0$  by translating the state space and the measure by  $-\alpha_{\text{pool}}$ . Since  $\varepsilon$  is not uninformative, there are two words  $\mathbf{v}, \mathbf{v}'$  with  $\varepsilon(\mathbf{v} | \mathbf{v}) \neq \varepsilon(\mathbf{v} | \mathbf{v}')$ . Assume that  $\varepsilon(\mathbf{v} | \mathbf{v}) > 0$ , applying similar arguments with swapped roles in the following otherwise.

Note that for any normal vector  $\vec{n} \in \mathbb{R}^L$  the corresponding  $L - 1$  dimensional hyperplane  $H := \{t \mid \langle \vec{n}, t \rangle = 0\}$  separated  $T$  into two disjoint convex cells  $C_{\mathbf{v}} := T \cap H_+ = \{t \in T \mid \langle \vec{n}, t \rangle > 0\}$  and  $C_{\mathbf{v}'} = T \cap H_- = \{t \in T \mid \langle \vec{n}, t \rangle \leq 0\}$ . Since  $\mu_0$  is absolutely continuous w.r.t. to the Lebesgue measure, we can choose  $\vec{n}$  in a way that both cells have positive  $\mu_0$ -measure and  $0 = \alpha_{\text{pool}} \in C_{\mathbf{v}'}$ . Note that  $0 \neq \mu_0(C_{\mathbf{v}}) \cdot \int_{C_{\mathbf{v}}} t \mu_0(dt) = \mathbb{E}_{C_{\mathbf{v}}, \mu_0}[t] \in C_{\mathbf{v}}$  by convexity of  $C_{\mathbf{v}}$  and Lemma 6.1 and thus  $0 \neq \int_{C_{\mathbf{v}}} t \mu_0(dt)$ .

Now define  $\pi(t) = \mathbf{v}$  if  $t \in C_{\mathbf{v}}$  and  $\pi(t) = \mathbf{v}'$  otherwise. Since we have  $\lambda^\pi(\mathbf{v}) = \varepsilon(\mathbf{v} | \mathbf{v})\mu_0(C_{\mathbf{v}}) + \varepsilon(\mathbf{v} | \mathbf{v}')\mu_0(C_{\mathbf{v}'}) > 0$ , it is sufficient to show  $\hat{\alpha}(\mathbf{v}) \neq 0 = \alpha_{\text{pool}}$  to prove  $L(\pi, \hat{\alpha}) > L_{\text{pool}}$  by Proposition 3.2. We use Lemma 6.1 again and evaluate

$$0 \neq \hat{\alpha}(\mathbf{v}) \cdot \lambda^\pi(\mathbf{v}) = \int_T \varepsilon(\mathbf{v} | \pi(t)) \cdot t \mu_0(dt) \quad (49)$$

$$\iff 0 \neq \varepsilon(\mathbf{v} | \mathbf{v}) \int_{C_{\mathbf{v}}} t \mu_0(dt) + \varepsilon(\mathbf{v} | \mathbf{v}') \int_{C_{\mathbf{v}'}} t \mu_0(dt). \quad (50)$$

If  $\varepsilon(\mathbf{v} | \mathbf{v}') = 0$ , we conclude  $\hat{\alpha}(\mathbf{v}) \neq 0$  as  $\varepsilon(\mathbf{v} | \mathbf{v}) > 0$  and  $\int_{C_{\mathbf{v}}} t \mu_0(dt) \neq 0$ . If  $\varepsilon(\mathbf{v} | \mathbf{v}') > 0$  we divide by  $\varepsilon(\mathbf{v} | \mathbf{v}')$  and find

$$0 = \alpha_{\text{pool}} = \mathbb{E}_{\mu_0}[t] = \int_{C_{\mathbf{v}}} t \mu_0(dt) + \int_{C_{\mathbf{v}'}} t \mu_0(dt), \quad (51)$$

$$0 \neq \frac{\varepsilon(\mathbf{v} | \mathbf{v})}{\varepsilon(\mathbf{v} | \mathbf{v}')} \cdot \int_{C_{\mathbf{v}}} t \mu_0(dt) + \int_{C_{\mathbf{v}'}} t \mu_0(dt). \quad (52)$$

As  $\varepsilon(\mathbf{v} | \mathbf{v}) \neq \varepsilon(\mathbf{v} | \mathbf{v}')$ , we conclude  $\hat{\alpha}(\mathbf{v}) \neq 0$ .  $\square$

*Proof of Proposition 6.3.* Revisit the argmin-set of sender (11) and recall that  $\|x - y\|_2^2 = \|x\|_2^2 - 2\langle x, y \rangle + \|y\|_2^2$ . In state  $t$ , the sender strictly prefers to send  $\mathbf{v}$  instead of  $\mathbf{v}'$  if and only if

$$\sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \mathbf{v}) \cdot \|t - \alpha(\mathbf{w})\|_2^2 < \sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \mathbf{v}') \cdot \|t - \alpha(\mathbf{w})\|_2^2 \quad (53)$$

$$\iff \sum_{\mathbf{w}} (\varepsilon(\mathbf{w} | \mathbf{v}') - \varepsilon(\mathbf{w} | \mathbf{v})) \cdot (-2\langle t, \alpha(\mathbf{w}) \rangle + \|\alpha(\mathbf{w})\|_2^2) > 0 \quad (54)$$

By linearity of the scalar product, convexity and the topological properties become clear. For the weak preference substitute the proper inequality accordingly.  $\square$



*Proof of Proposition 6.4.* From (54) the sender is indifferent between sending  $\mathbf{w}$  and  $\mathbf{v}'$  if the state is  $t$  if and only if

$$0 = \sum_{\mathbf{w}} (\varepsilon(\mathbf{w} | \mathbf{v}') - \varepsilon(\mathbf{w} | \mathbf{v})) \cdot (-2 \langle t, \alpha(\mathbf{w}) \rangle + \|\alpha(\mathbf{w})\|_2^2) \quad (55)$$

$$\begin{aligned} \iff 0 = -2 \cdot \left\langle t, \sum_{\mathbf{w}} (\varepsilon(\mathbf{w} | \mathbf{v}') - \varepsilon(\mathbf{w} | \mathbf{v})) \cdot \alpha(\mathbf{w}) \right\rangle \\ + \sum_{\mathbf{w}} (\varepsilon(\mathbf{w} | \mathbf{v}') - \varepsilon(\mathbf{w} | \mathbf{v})) \cdot \|\alpha(\mathbf{w})\|_2^2. \end{aligned} \quad (56)$$

Note that if  $\vec{x} := \sum_{\mathbf{w}} (\varepsilon(\mathbf{w} | \mathbf{v}') - \varepsilon(\mathbf{w} | \mathbf{v})) \cdot \alpha(\mathbf{w}) \neq 0$  the solution set to (56) is the translation of the  $L - 1$  dimensional hyperplane perpendicular to the vector  $\vec{x}$  by a particular solution (if it exists, otherwise it is the empty set) and thus a null set in  $\mathbb{R}^L$  w.r.t.  $\mu_0$ . If  $\vec{x} = 0$  we can only have indifference if also  $\sum_{\mathbf{w}} (\varepsilon(\mathbf{w} | \mathbf{v}') - \varepsilon(\mathbf{w} | \mathbf{v})) \cdot \|\alpha(\mathbf{w})\|_2^2 = 0$ . If this is the case, any  $t$  in  $T$  (even  $\mathbb{R}^L$ ) solves (56). Otherwise, (56) is equivalent to

$$0 = \sum_{\mathbf{w}} (\varepsilon(\mathbf{w} | \mathbf{v}') - \varepsilon(\mathbf{w} | \mathbf{v})) \cdot \alpha(\mathbf{w}) \quad (57)$$

$$\iff 0 = \sum_{\mathbf{w}} \left( \tilde{p}^{d(\mathbf{w}, \mathbf{v}')} - \tilde{p}^{d(\mathbf{w}, \mathbf{v})} \right) \cdot \alpha(\mathbf{w}) =: Q(\tilde{p}), \quad (58)$$

where  $Q$  is a (vector-valued) polynomial of degree at most  $n$  in  $\tilde{p}$ . Evidently,  $Q$  has a zero in  $\tilde{p} = 1$ , i.e.,  $p = \frac{m}{m+1}$ , representing the unformativeness bound. It does not vanish in  $p = \tilde{p} = 0$  as its constant coefficient is  $\alpha(\mathbf{v}') - \alpha(\mathbf{v}) \neq 0$ . Consequently,  $Q$  is not the zero polynomial and has at most  $n - 1$  further zeros on  $p \in (0, 1] \setminus \{\frac{m}{m+1}\}$ .  $\square$

*Proof of Proposition 7.1.* We adapt the proof of Jäger et al. (2011) to account for noise. It suffices to show continuity and boundedness of  $L$  as this implies continuity of  $\Lambda$  in the weak topology. The rest of the assertions follow well-known lines (Heifetz et al. (2007), Hofbauer et al. (2009)) as well as Bhatia & Szegö (2002) for the last statement.

Let  $(\pi_k)_k$  be a sequence of communication strategies converging uniformly to  $\pi$ , i.e., for all  $\rho' > 0$  there is an  $M$  such that for all  $t \in T$  we simultaneously find  $d(\pi_k(t), \pi(t)) < \rho'$  for  $k > M$ . As  $d$  has only values in  $\{0, \dots, n\}$ , this is equivalent to  $\pi_k \equiv \pi$  for all  $k > N_0$  for some  $N_0$ . Let  $\rho > 0$  be arbitrary. As  $T$  is compact and  $|\cdot|$  as well as  $\ell$  are continuous, there is  $\delta > 0$  such that  $|\ell(a) - \ell(b)| < \rho$  if  $\|a - b\| < \delta$ . Furthermore, let  $(\alpha_k)_k$  be a sequence converging to  $\alpha$  uniformly on  $T^{\#W} \subseteq \mathbb{R}^{\#W}$ . Then there is  $N_1 \geq N_0$  such that for all  $t \in T$  and  $\mathbf{w} \in W$  we have

$\|t - \alpha_k(\mathbf{w}) - (t - \alpha(\mathbf{w}))\| = \|\alpha_k(\mathbf{w}) - \alpha(\mathbf{w})\| < \delta$  for all  $k > N_1$ . Hence, for all  $k > N_1$  we find

$$|L(\pi_k, \alpha_k) - L(\pi, \alpha)| \tag{59}$$

$$\leq \int_T \left| \sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \pi_k(t)) \ell(\|t - \alpha_k(\mathbf{w})\|) - \sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \pi(t)) \ell(\|t - \alpha(\mathbf{w})\|) \right| \mu_0(dt) \tag{60}$$

$$= \int_T \left| \sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \pi(t)) (\ell(\|t - \alpha_k(\mathbf{w})\|) - \ell(\|t - \alpha(\mathbf{w})\|)) \right| \mu_0(dt) \tag{61}$$

$$\leq \int_T \sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \pi(t)) |\ell(\|t - \alpha_k(\mathbf{w})\|) - \ell(\|t - \alpha(\mathbf{w})\|)| \mu_0(dt) \tag{62}$$

$$\leq \int_T \sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \pi(t)) \cdot \rho \mu_0(dt) = \rho. \tag{63}$$

Finally, boundedness of  $L$  follows from compactness of  $T$  and continuity of  $\ell$ , since  $\bar{\ell} := \sup_{t \in T} |\ell(\|t\|)| < \infty$ . For any  $\pi, \alpha$

$$|L(\pi, \alpha)| \leq \int_T \sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \pi(t) \cdot) |\ell(\|t - \alpha(\mathbf{w})\|)| \mu_0(dt) \tag{64}$$

$$\leq \int_T \sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \pi(t)) \cdot \bar{\ell} \mu_0(dt) = \bar{\ell} < \infty. \tag{65}$$

□

### 9.3 Calculations of Section 6.2

We prove the assertions and formulae of Section 6.2 by the following lengthy calculations which are structured as follows. We begin by calculating the minimizing interpretations and the expected loss. Afterwards, fixing any of the calculated interpretation maps we show that the given tessellation is indeed an optimal one, even uniquely up to null sets. Denote by  $\mathbb{E}_C$  the expectation operator of the measure  $\mu_0$  restricted to the cell  $C$ . We begin by noting that for any considered cell  $\mu_0(C_v) = \frac{1}{4}$  and thus for any  $\mathbf{w}$

$$\lambda^\pi(\mathbf{w}) = \int_T \varepsilon(\mathbf{w} | \pi(t)) dt = \sum_v \varepsilon(\mathbf{w} | v) \cdot \mu_0(C_v) = \frac{1}{4}. \tag{66}$$

#### (i) Interpretations and expected Loss

Denote by  $\square(AA)$  and  $\Delta(AA)$  the resp. square or triangular cell in Figure 4. The following calculations are down for the resp. languages.

1 a) We derive the center of gravity of each cell, say the one for  $AB$ .

$$\mathbb{E}_{\square(AB)}[t] := \mathbb{E}_{\mu_0, \square(AB)}[t] = \mu_0(\square(AB))^{-1} \cdot \int_{\square(AB)} t \, dt \quad (67)$$

$$= 4 \cdot \int_{-\frac{1}{2}}^0 \int_{-\frac{1}{2}}^0 (t_1, t_2) \, dt_1 dt_2 = \left(-\frac{1}{4}, -\frac{1}{4}\right). \quad (68)$$

Similarly, or by using symmetry arguments, we obtain the expected values for  $AA, BA, BB$  which are, resp.  $(-\frac{1}{4}, \frac{1}{4}), (\frac{1}{4}, \frac{1}{4}), (\frac{1}{4}, -\frac{1}{4})$ . Let us now calculate, e.g., the optimal action  $\hat{\alpha}(AA)$ .

$$\hat{\alpha}(AA) = \mathbb{E}_{\mu_{AA}^\pi}[t] = \lambda^\pi(AA)^{-1} \cdot \int_T \varepsilon(AA | \pi(t)) \cdot t \, \mu_0(dt) \quad (69)$$

$$= 4 \cdot \sum_w \varepsilon(AA | w) \cdot \int_{\square(w)} t \, dt \quad (70)$$

$$= \sum_w \varepsilon(AA | w) \cdot \mathbb{E}_{\square(w)} \quad (71)$$

$$= (1-p)^2 \cdot \left(-\frac{1}{4}, \frac{1}{4}\right) + p(1-p) \cdot \left(\left(-\frac{1}{4}, -\frac{1}{4}\right) + \left(\frac{1}{4}, \frac{1}{4}\right)\right) + p^2 \cdot \left(\frac{1}{4}, -\frac{1}{4}\right) \quad (72)$$

$$= \frac{1}{4} \cdot (-1 + 2p, 1 - 2p). \quad (73)$$

Analogously, by symmetry arguments or using Lemma 6.1 we get  $\hat{\alpha}(AB) = \frac{1}{4} \cdot (-1 + 2p, -1 + 2p)$ ,  $\hat{\alpha}(BA) = \frac{1}{4} \cdot (1 - 2p, 1 - 2p)$ ,  $\hat{\alpha}(BB) = \frac{1}{4} \cdot (1 - 2p, -1 + 2p)$ . For each word  $w$  we see  $\|\alpha(w) - \alpha_{\text{pool}}\|_2 \searrow 0$  for  $p \rightarrow \frac{1}{2}$ , where  $\alpha_{\text{pool}} = (0, 0)$  is the center of the whole state space.

We are now set to calculate the expected loss and start by observing that each interpretation has the same norm:

$$\|\hat{\alpha}(w)\|_2^2 = \left\| \frac{1}{4} \cdot (1 - 2p, 1 - 2p) \right\|_2^2 = \frac{1}{8} \cdot (1 - 2p)^2. \quad (74)$$

Having calculated  $\mathbb{E}_T[\|t\|_2^2] = \frac{1}{6}$ , we use (20) to obtain the expected loss

$$L(\pi_{1,a}, \hat{\alpha}) = \frac{1}{6} - \sum_w \frac{1}{4} \cdot \frac{1}{8} \cdot (1 - 2p)^2 = \frac{1}{6} - \frac{1}{8} \cdot (1 - 2p)^2. \quad (75)$$

One clearly sees that the expected loss is monotonically increasing in  $p \in [0, \frac{1}{2}]$

- b) Using the calculations from (a) we can directly compute the optimal interpretations, only keeping in mind that the centers of gravity are switched for  $BA$  and  $BB$ . We obtain  $\hat{\alpha}(AA) = \frac{1}{4} \cdot (-1 + 2p, 1 - 4p + 4p^2)$ ,  $\hat{\alpha}(AB) = \frac{1}{4} \cdot (-1 + 2p, -1 + 4p - 4p^2)$ ,  $\hat{\alpha}(BA) = \frac{1}{4} \cdot (1 - 2p, -1 + 4p - 4p^2)$ ,  $\hat{\alpha}(BB) = \frac{1}{4} \cdot (1 - 2p, 1 - 4p + 4p^2)$ . Thus, for any word  $w$  we have

$$\|\hat{\alpha}(w)\|_2^2 = \frac{1}{16} \cdot ((1 - 2p)^2 + (1 - 2p)^4), \quad (76)$$

resulting in an expected loss of

$$L(\pi_{1,b}, \hat{\alpha}) = \frac{1}{6} - \frac{1}{16} \cdot ((1 - 2p)^2 + (1 - 2p)^4) \quad (77)$$

$$= \frac{1}{6} - \frac{1}{8} \cdot (1 - 2p)^2 \cdot (1 - 2p + 2p^2). \quad (78)$$

We observe for  $0 < p < \frac{1}{2}$

$$L(\pi_{1,a}, \hat{\alpha}) < L(\pi_{1,b}, \hat{\alpha}), \quad (79)$$

thus, the language putting distant words farther away from one another achieves a lower expected loss.

- 2 a) The expected values of each colored area can be determined to be  $\mathbb{E}_{\Delta(AA)}[t] = (0, \frac{1}{3})$ ,  $\mathbb{E}_{\Delta(AB)}[t] = (-\frac{1}{3}, 0)$ ,  $\mathbb{E}_{\Delta(BA)}[t] = (\frac{1}{3}, 0)$ ,  $\mathbb{E}_{\Delta(BB)}[t] = (0, -\frac{1}{3})$ . Optimal actions can be computed to be  $\hat{\alpha}(AA) = (0, -\frac{1}{3} + \frac{2}{3}p)$ ,  $\hat{\alpha}(AB) = (-\frac{1}{3} + \frac{2}{3}p, 0)$ ,  $\hat{\alpha}(BA) = (\frac{1}{3} - \frac{2}{3}p, 0)$ ,  $\hat{\alpha}(BB) = (0, \frac{1}{3} + \frac{2}{3}p)$ .

We thus get

$$\|\alpha(w)\|_2^2 = \|(0, \frac{1}{3} - \frac{2}{3}p)\|_2^2 = \frac{1}{9} \cdot (1 - 2p)^2. \quad (80)$$

The resulting expected loss is

$$L(\pi_{2,a}, \hat{\alpha}) = \frac{1}{6} - \frac{1}{9} \cdot (1 - 2p)^2, \quad (81)$$

which is strictly higher than  $L(\pi_{1,a}, \hat{\alpha})$  for any  $p \in [0, \frac{1}{2})$ .

- b) Optimal actions can be calculated to be  $\hat{\alpha}(AA) = \frac{1}{3} \cdot (-p + 2p^2, 1 - 3p + 2p^2)$ ,  $\hat{\alpha}(AB) = \frac{1}{3} \cdot (-1 + 3p - 2p^2, p - 2p^2)$ ,  $\hat{\alpha}(BA) = \frac{1}{3} \cdot (p - 2p^2, -1 + 3p - 2p^2)$ ,  $\hat{\alpha}(BB) = \frac{1}{3} \cdot (1 - 3p + 2p^2, -p + 2p^2)$ .

We thus get for any word  $w$

$$\|\hat{\alpha}(w)\|_2^2 = \frac{1}{9} (1 - 2p)^2 (1 - 2p + 2p^2) \quad (82)$$

and hence

$$L(\pi_{2,b}, \hat{\alpha}) = \frac{1}{6} - \frac{1}{9} (1 - 2p)^2 (1 - 2p + 2p^2), \quad (83)$$

which is worse than  $L(\pi_{2,a}, \hat{\alpha})$  for  $0 < p < \frac{1}{2}$ .

**(ii) Optimal cell structure**

To start with, we simplify the expressions from Proposition 6.3 and Proposition 6.4 for  $W = \{A, B\}^2$ . To this end, fix w.l.o.g. the word  $AA$  and derive conditions on a fixed  $t \in T$  for  $AA$  to be the optimal word.

(i) In state  $t$  the sender prefers to send  $AA$  over  $BB$  if and only if

$$\sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \mathbf{v}) \|t - \alpha(AA)\|_2^2 < \sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \mathbf{v}) \|t - \alpha(BB)\|_2^2 \quad (84)$$

$$\iff \|t - \alpha(AA)\|_2 < \|t - \alpha(BB)\|_2. \quad (85)$$

(ii) In state  $t$  the sender prefers  $AA$  over  $AB$  (the case  $BA$  is analogous) if and only if

$$\sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \mathbf{v}) \|t - \alpha(AA)\|_2^2 < \sum_{\mathbf{w}} \varepsilon(\mathbf{w} | \mathbf{v}) \|t - \alpha(AB)\|_2^2 \quad (86)$$

$$\iff \|t - \alpha(AA)\|_2^2 - \|t - \alpha(AB)\|_2^2 < \tilde{p} (\|t - \alpha(BB)\|_2^2 - \|t - \alpha(BA)\|_2^2) \quad (87)$$

$$\iff 2 \langle t, \alpha(AB) - \alpha(AA) + \tilde{p}(\alpha(BB) - \alpha(BA)) \rangle + \|\alpha(AA)\|_2^2 - \|\alpha(AB)\|_2^2 + \tilde{p}(\|\alpha(BA)\|_2^2 - \|\alpha(BB)\|_2^2) < 0. \quad (88)$$

Whereas in (i) we clearly see that the set of states for which the sender is indifferent between sending  $AA$  and  $BB$  lie on the perpendicular bisector of  $\alpha(AA)$  and  $\alpha(BB)$  if the interpretations do not agree, it is not so obvious in case (ii). What we can say for sure is, that, as long as  $\alpha(AB) - \alpha(AA) + \tilde{p}(\alpha(BB) - \alpha(BA))$  is not the zero vector, the set of indifferent states is again a null set as it is the intersection of a line and  $T$ .

To drop some notation, we just write  $AA$  instead of  $\alpha(AA)$  from Table 1 when talking about points in  $T$ . Consider the variants (a) and (b) respectively and let  $t \in \square(AA)$  (resp.  $t \in \Delta(AA)$ ) be in the interior.

(a) Observe that

$$\|t - AA\|_2 < \|t - AB\|_2, \|t - BA\|_2 < \|t - BB\|_2. \quad (89)$$

Obviously, sending  $AA$  is preferred to  $BB$  as  $\|t - AA\|_2 < \|t - BB\|_2$ .

Realizing that

$$\|t - AA\|_2^2 - \|t - AB\|_2^2 < 0 < \tilde{p} \cdot (\|t - BB\|_2^2 - \|t - BA\|_2^2), \quad (90)$$

reveals that sending  $AA$  is preferred to  $AB$  (and analogously  $BA$ ). Thus,  $AA$  is the unique best word to be send.

(b) As before, preferring  $AA$  to  $BB$  is clear from  $\|t - AA\|_2 < \|t - BB\|_2$ . Since

$$0 \leq \|t - AA\|_2 < \|t - AB\|_2, \|t - BB\|_2 < \|t - BA\|_2, \quad (91)$$

we find

$$\|t - BA\|_2^2 - \|t - AA\|_2^2 > \left| \|t - AB\|_2^2 - \|t - BB\|_2^2 \right| \quad (92)$$

$$> \tilde{p} \cdot \left| \|t - AB\|_2^2 - \|t - BB\|_2^2 \right| \quad (93)$$

$$\geq \tilde{p} \cdot (\|t - AB\|_2^2 - \|t - BB\|_2^2), \quad (94)$$

proving that  $AA$  is preferred to  $BA$ . Eventually, using  $AA = -BA$ ,  $AB = -BB$  and that  $\|\alpha(w)\|$  is constant, we find

$$2 \cdot \langle t, -AA + AB + \tilde{p}(BB - BA) \rangle + \|AA\|_2^2 - \|AB\|_2^2 + \tilde{p}(\|BB\|_2^2 - \|BA\|_2^2) \quad (95)$$

$$= 4 \cdot \left\langle t, \underbrace{\frac{AB+BA}{2}}_{=:P} \right\rangle. \quad (96)$$

The expression (96) is smaller than zero in both cases for  $t \in \Delta(AA)$ :

1.  $t_1 < 0, t_2 > 0$  and  $P_1 = 0, P_2 < 0$ .
2.  $t = (y, z)$  with  $z > 0, |y| < z$  and  $P = (-x, x), x > 0$ .

Thus, sending  $AA$  is preferred to  $AB$  as well.

The calculations above show that the borders of the cells consist precisely of the points for which the sender is indifferent between sending the resp. messages.

## References

- Bhatia, N. & Szegö, G. (2002), *Stability Theory of Dynamical Systems*, Classics in Mathematics, Springer Berlin Heidelberg. <https://books.google.de/books?id=wP5dwTS6jg0C>.
- Blume, A. & Board, O. (2013), ‘Language barriers’, *Econometrica* **81**(2), 781–812.
- Blume, A., Board, O. J. & Kawamura, K. (2007), ‘Noisy talk’, *Theoretical Economics* **2**(4), 395–440.
- Cover, T. M. & Thomas, J. A. (2006), *Elements of Information Theory 2nd Edition (Wiley Series in Telecommunications and Signal Processing)*, Wiley-Interscience.

- Crawford, V. P. & Sobel, J. (1982), ‘Strategic information transmission’, *Econometrica* **50**(6), 1431–1451. <http://www.jstor.org/stable/1913390>.
- Cremer, J., Garicano, L. & Prat, A. (2007), ‘Language and the theory of the firm’, *The Quarterly Journal of Economics* **122**(1), 373–407.
- Cressman, R., Hofbauer, J. & Riedel, F. (2006), ‘Stability of the replicator equation for a single species with a multi-dimensional continuous trait space’, *Journal of Theoretical Biology* **239**(2), 273–288. Special Issue in Memory of John Maynard Smith, Available at <https://www.sciencedirect.com/science/article/pii/S0022519305003887>.
- Frege, G. (1892), ‘Über Sinn und Bedeutung’, *Zeitschrift für Philosophie und philosophische Kritik* **100**, 25–50.
- Grice, H. P. (1975), Logic and conversation, in ‘Speech acts’, Brill, pp. 41–58.
- Gärdenfors, P. (2004), *Conceptual spaces: The geometry of thought*, MIT press.
- Hamming, R. W. (1950), ‘Error detecting and error correcting codes’, *The Bell System Technical Journal* **29**(2), 147–160.
- Heifetz, A., Shannon, C. & Spiegel, Y. (2007), ‘What to maximize if you must’, *Journal of Economic Theory* pp. 31–57.
- Hernández, P. & von Stengel, B. (2014), ‘Nash codes for noisy channels’, *Operations Research* **62**(6), 1221–1235.
- Hofbauer, J., Oechssler, J. & Riedel, F. (2009), ‘Brown–von neumann–nash dynamics: The continuous strategy case’, *Games and Economic Behavior* **65**(2), 406–429. Available at <https://www.sciencedirect.com/science/article/pii/S0899825608000651>.
- Hofbauer, J. & Sigmund, K. (1998), *Evolutionary games and population dynamics*, Cambridge university press.
- Jeitschko, T. D. & Normann, H.-T. (2012), ‘Signaling in deterministic and stochastic settings’, *Journal of Economic Behavior & Organization* **82**(1), 39–55.
- Jäger, G. (2007), ‘The evolution of convex categories’, *Linguistics and Philosophy* **30**, 551–564.
- Jäger, G., Metzger, L. P. & Riedel, F. (2011), ‘Voronoi languages: Equilibria in cheap-talk games with high-dimensional types and few signals’, *Games and Economic Behavior* **73**(2), 517 – 537.

- Kamenica, E. & Gentzkow, M. (2011), ‘Bayesian persuasion’, *American Economic Review* **101**(6), 2590–2615. Available at <http://www.aeaweb.org/articles?id=10.1257/aer.101.6.2590>.
- Lewis, D. (1969), *Convention: A philosophical study*, John Wiley & Sons.
- Lipman, B. L. (2009), Why is language vague. Available at <https://sites.bu.edu/blipman/files/2021/10/vague5.pdf>.
- MacKay, D. J. C. (2002), *Information Theory, Inference & Learning Algorithms*, Cambridge University Press, USA.
- Martel, J., Van Wesep, E. D. & Van Wesep, R. (2019), On ratings: A theory of non-strategic information transmission. Working paper, available at .
- Miller, A. (2007), *Philosophy of language*, Taylor & Francis.
- Myerson, R. B. (1991), *Game Theory: Analysis of Conflict*, Harvard University Press. Available at <http://www.jstor.org/stable/j.ctvj522>.
- Nowak, M. A. & Krakauer, D. C. (1999), ‘The evolution of language’, *Proceedings of the National Academy of Sciences* **96**(14), 8028–8033.
- Oechssler, J. & Riedel, F. (2001), ‘Evolutionary dynamics on infinite strategy spaces’, *Economic Theory* **17**(1), 141–162. <https://ideas.repec.org/a/spr/joecth/v17y2001i1p141-162.html>.
- Ortoleva, P. (2012), ‘Modeling the change of paradigm: Non-Bayesian reactions to unexpected news’, *American Economic Review* **102**(6), 2410–2436.
- Pelletier, F. J. (2001), ‘Did Frege believe Frege’s principle?’, *Journal of Logic, Language and information* **10**, 87–114.
- Roth, R. (2006), *Introduction to Coding Theory*, Cambridge University Press, USA.
- Rubinstein, A. (1989), ‘The electronic mail game: Strategic behavior under “almost common knowledge”’, *The American Economic Review* **79**, 385–391.
- Shannon, C. E. (1948), ‘A mathematical theory of communication’, *The Bell system technical journal* **27**(3), 379–423.
- Smith, J. M. & Price, G. R. (1973), ‘The logic of animal conflict’, *Nature* **246**(5427), 15–18.



- Sobel, J. (2015), Broad terms and organizational codes. Unpublished paper, Department of Economics, University of California, San Diego.[1138], Available at <https://ipl.econ.duke.edu/seminars/system/files/seminars/1400.pdf>.
- Sorensen, R. (2023), Vagueness, *in* E. N. Zalta & U. Nodelman, eds, ‘The Stanford Encyclopedia of Philosophy’, Winter 2023 edn, Metaphysics Research Lab, Stanford University.
- Spence, M. (1978), Job market signaling, *in* ‘Uncertainty in economics’, Elsevier, pp. 281–306.
- Szabó, Z. G. (2022), Compositionality, *in* E. N. Zalta & U. Nodelman, eds, ‘The Stanford Encyclopedia of Philosophy’, Fall 2022 edn, Metaphysics Research Lab, Stanford University.