# Unified Approach to Optimization Techniques in Shannon Theory

Peter Harremoës and Flemming Topsøe[*]

October 12, 2001

## 1 Theory

### 1.1 The model

To focus on the main ideas and avoid technical problems we will assume that we work with alphabets $\mathbb{A}, \mathbb{B}$ and $\mathbb{E}$ of finite size. Consider an information system $(\mathbb{A}, \mathbb{B}, \mathbb{E}, \mathcal{M}, \Gamma, c)$. Let $\mathcal{M} : \mathbb{A} \to \mathbb{M}_1^+ (\mathbb{B} \times \mathbb{E})$ be a Markov kernel. Let $\Gamma$ be a subset of $M_1^+ (\mathbb{A})$. Let $c : \Gamma \to \mathbb{R}$ be a convex function. A coding strategy $\kappa$ is a map from $\mathbb{B} \times \mathbb{E}$ to $[0; \infty]$ which satisfy a variant of Krafft's inequality

$$\sum_{b \in \mathbb{B}} \exp \left( -\kappa \left( b \mid e \right) \right) \leq 1 \text{ for all } e \in \mathbb{E} \, .$$

The interpretation is as follows. For each $e \in \mathbb{E}$ we code the elements in $\mathbb{B}$ in such way that the code length of letter $b \in \mathbb{B}$ given $e \in \mathbb{E}$ is $\kappa \left( b \mid e \right)$. Let a 2-persons 0-sum game between Alice and Eve be defined by that Alice plays a distribution $P$ from $\Gamma$ and Eve plays a coding strategy $\kappa$. The object function should be

$$\langle \kappa, P \rangle_c = \sum_{e \in \mathbb{E}} \mathcal{M} \left( P \right) \left( e \right) \langle \kappa \left( \cdot \mid e \right), \mathcal{M} \left( P \right) \left( \cdot \mid e \right) \rangle - c \left( P \right) \, .$$

The cost function $c$ is introduced such that the model will cover as many examples as possible. At this state we may think of $c \left( P \right)$ as a cost for Alice of sending $P$ through the channel given by $\mathcal{M}$. Alice wants the object function to be big and Eve wants it to be small. One interpretation of the model is that Alice sends a secret message $P$ and an eavesdropper Eve wants to extract as much information as possible. Alice wants to prevent Eve from extract information.

Alice is only allowed to send $P \in \Gamma$, because otherwise Bob will not be able to recover the encrypted information. Define

$$H_c\left(P\right) = \inf_{\kappa} \left\langle \kappa, P \right\rangle_c$$
$$R\left(\kappa\right) = \sup_{P \in \Gamma} \left\langle \kappa, P \right\rangle_c$$

and

$$H_c\left(\Gamma\right) = \sup_{P \in \Gamma} H_c\left(P\right)$$
$$R_{\min} = \inf_{\kappa} R\left(\kappa\right)$$

With these definitions one automatically has

$$H_c\left(\Gamma\right) \leq R_{\min} .$$

The information divergence $D$ shall always mean the mean of the conditional divergence on the output side given the letter in $\mathbb{E}$. The divergence from a distribution to a code/ciding strategy shall denote the divergence from the the distribution to the distribution/kernel corresponind to the code/coding strategy.

## 1.2 Results

**Proposition 1** $\left\langle \kappa, P \right\rangle_c = \left\langle \kappa_P, P \right\rangle_c + D\left(P \parallel \kappa\right)$ *and therefore* $H_c\left(P\right) = \left\langle \kappa_P, P \right\rangle_c$ .

**Proposition 2** *The function* $P \to H\left(P\right)$ *is continuous.*

    **Proof.** First remark that

$$H\left(P\right) = \left\langle \kappa_P, P \right\rangle$$
$$= \sum_{e \in \mathbb{E}} \mathcal{M}\left(P\right)\left(e\right) \cdot H\left(\left(\mathcal{M}P\right)\left(\cdot \mid e\right)\right)$$

    Assume $P_\lambda \to P$. Then $\mathcal{M}\left(P_\lambda\right)\left(e\right) \to \mathcal{M}\left(P\right)\left(e\right)$. If $\mathcal{M}\left(P\right)\left(e\right) \neq 0$ then $\left(\mathcal{M}P_\lambda\right)\left(\cdot \mid e\right) \to \left(\mathcal{M}P_\lambda\right)\left(\cdot \mid e\right)$ and therefore

$$H\left(\left(\mathcal{M}P_\lambda\right)\left(\cdot \mid e\right)\right) \to H\left(\left(\mathcal{M}P\right)\left(\cdot \mid e\right)\right) .$$

If $\mathcal{M}\left(P_\lambda\right)\left(e\right) \to \mathcal{M}\left(P\right)\left(e\right)$ and $\mathcal{M}\left(P\right)\left(e\right) = 0$ then

$$\left| \mathcal{M}\left(P_\lambda\right)\left(e\right) \cdot H\left(\left(\mathcal{M}P_\lambda\right)\left(\cdot \mid e\right)\right) \right| \leq \mathcal{M}\left(P_\lambda\right)\left(e\right) \cdot \log |B|$$
$$\to 0$$

    ■

**Lemma 3** *Assume that* $c$ *is continuous, and that* $\Gamma$ *is convex and compact. Then there exists an input distribution* $P_{opt}$ *such that*

$$H_c\left(P\right) + D\left(P \parallel P_{opt}\right) \leq H_c\left(\Gamma\right) .$$

2

**Proof.** The map $P \to H_c(P)$ is continuous and by compactness of $\Gamma$ there exists an input distribution $P_{opt}$ such that

$$H_c(P_{opt}) = H_c(\Gamma)$$

Let $P \in \Gamma$ be an arbitrary input distribution. By convexity we have

$$
\begin{aligned}
&H_c(\Gamma) \\
&\geq H_c((1-\alpha)P_{opt} + \alpha P) \\
&\geq (1-\alpha)H_c(P_{opt}) + \alpha H_c(P) \\
&+ (1-\alpha)D(P_{opt} \| (1-\alpha)P_{opt} + \alpha P) + \alpha D(P \| (1-\alpha)P_{opt} + \alpha P) \\
&\geq (1-\alpha)H_c(\Gamma) + \alpha H_c(P) + \alpha D(P \| (1-\alpha)P_{opt} + \alpha P)
\end{aligned}
$$

and therefore

$$0 \geq (-\alpha)H_c(\Gamma) + \alpha H_c(P) + \alpha D(P_{opt} \| (1-\alpha)P_{opt} + \alpha P)$$
$$H_c(\Gamma) \geq H_c(P) + D(P \| (1-\alpha)P_{opt} + \alpha P)$$

Let $\alpha \to 0$ and use lower semi continuity to obtain

$$H_c(P) + D(P \| P_{opt}) \leq H_c(\Gamma)$$
$$\left\langle \kappa_{P_{opt}}, P \right\rangle_c \leq H_c(\Gamma)$$

∎

**Theorem 4** *Assume that $c$ is continuous and that $\Gamma$ is convex and compact. Then $H_c(\Gamma) < \infty$, and the value of the game exists and equals $H_c(\Gamma)$. Further Alice has an optimal input distribution $P_{opt}$ and Eve has an optimal coding strategy $\kappa_{opt}$. The trivial inequality*

$$H_c(P) \leq H_c(\Gamma) \leq R(\kappa)$$

*can be improved to*

$$H_c(P) + D(P \| \kappa_{opt}) \leq H_c(\Gamma) \leq R(\kappa) - D(P_{opt} \| \kappa)$$

**Proof.** The equation $H_c(P) + D(P \| P_{opt}) = \left\langle \kappa_{P_{opt}}, P \right\rangle_c$ shows that

$$\left\langle \kappa_{P_{opt}}, P \right\rangle_c \leq H_c(\Gamma) \tag{1}$$

This shows that $\kappa_{P_{opt}}$ is an optimal strategy for Eve.

Further we have

$$
\begin{aligned}
\langle \kappa, P_{opt} \rangle_c &= D(P_{opt} \| \kappa) + H_c(P_{opt}) \\
&\geq H_c(\Gamma)
\end{aligned}
$$

To demonstrate the last inequality we write

$$R(\kappa) \geq \langle \kappa, P_{opt} \rangle_c$$

∎

We remark that the inequality (1) is the condition for $\kappa_{opt}$ to be a Nash equilibrium coding strategy.

**Condition 5 (Kuhn-Tucker conditions)** *If $\Gamma$ is compact, then a necessary and sufficient condition for $\kappa$ to be optimal is that there exists input distributions $P_1, P_2, \cdots, P_n \in \Gamma$ such that $\kappa$ is induced by a convex combination of $P_1, P_2, \cdots, P_n$ and such that $R(\kappa) \geq \langle \kappa, P_i \rangle_c$ for all $i$.*

For the value of the game to exist neither continuity of $c$ nor compactness of $\Gamma$ is necessary, and one can get

**Theorem 6** *Assume that $\Gamma$ is convex and that $H_c(\Gamma) < \infty$. Then the value of the game exists and equals $H_c(\Gamma)$. Further Eve has an optimal coding strategy $\kappa_{opt}$.*

# 2  Examples

## 2.1  Maximum Entropy Principle

First we consider the *code length game* where the cost function is zero and there are no side conditions. In this game $\mathbb{A} = \mathbb{B}$, and $\mathcal{M}, \mathbb{E}$ and $c$ are trivial. We arrive at a game where $H_c(\Gamma)$ is the maximum entropy (if it exists). This game was considered in great detail in [3].

## 2.2  Minimal free energy

In thermodynamics one defines Helmholz' free energy $A = U - TS$ where $U$ is the inner energy, $T$ is the absolute temperature and $S$ is the thermodynamic entropy. In a physical or chemical system where the volume and temperature are kept fixed Helmholz' free energy will tend to a minimum. The thermodynamic entropy $S$ is related to the information theoretic entropy $H$ by the relation $S = nkH$ where $n$ is the number of molecules and $k$ is Bolzmann's constant. Now, consider a system where the molecules can be in states with energies $E_1$, $E_2,...,E_n$ , and let $p_i$ denote the probability that a molecule is in state $E_i$ . Then we have to minimize

$$n \sum p_i E_i - nkTH(p_1, p_2, ..., p_n)$$

or, equivalently maximize

$$H(p_1, p_2, ..., p_n) - \sum p_i \cdot \frac{E_i}{kT} \ .$$

This corresponds to a game with $\mathbb{A} = \mathbb{B}$, and $\mathcal{M}, \mathbb{E}$ are trivial, and the cost function $c$ is given by $c(p) = \sum p_i \cdot \frac{E_i}{kT}$ . Thus the cost for the system to use state $i$ is proportional to the energy $E_i$ of the state. The Kuhn-Tucker conditions shows that an optimal code is obtained if $\kappa(i) = \frac{E_i}{kT} + \gamma$ where $\gamma$ is a suitable chosen constant. This is in agreement with the solution found in textbooks [5], where the distribution with minimal free energy is shown to have point probabilities proportional to $\exp\left(-\frac{E_i}{kT}\right)$.

## 2.3 Minimum Information Principle

Again we consider a game with no side conditions. In this game $\mathbb{A} = \mathbb{B}$, and $\mathcal{M}$, $\mathbb{E}$ are trivial. A reference distribution $Q$ in $M_+^1(\mathbb{A})$ is given and the cost function is given by $c(P) = \left\langle \log \frac{1}{Q(i)}, P \right\rangle$. The the object function is $\left\langle \kappa(i) - \log \frac{1}{Q(i)}, P \right\rangle$ and we get *the relative game* described in [6]. If $\Gamma$ is convex and compact the optimal strategy for Alice is the information projection of $Q$ on $\Gamma$, and the optimal strategy of Eve is the corresponding code. This amounts to the minimum information principle. The Kuhn-Tucker conditions leads naturally to the exponential families well-known from statistics.

## 2.4 Discrete Memoryless Channel

To model a discrete memoryless channel without side information, put $\Gamma = M_1^+(\mathbb{A})$ and $\mathbb{E}$ trivial. $M$ is the Markov kernel defining the channel as usual. The cost function is given by $c(P) = H(\mathcal{M}(P))$. Then the equilibrium of the game is the Gallager-Ryabko Theorem which states that maximal transmission rate equals minimal redundancy. The conditions 5 are the well-known Kuhn-Tucker conditions for DMC's [2, p. 191].

If $\mathbb{E}$ is non trivial this corresponds to having a channel with side conditions. The model actually also covers rate distortion theory as described in [1]

## 2.5 Exact Prediction

Some problems in prediction theory can be solved exactly. Here we shall only consider the simplest prediction problem where we have to guess the next letter in a sequence of letters from an alphabet $\mathbb{B} = \{a, b\}$. We consider a sequence of length 2, and assume that first and second letter are independent and identically distributed. We have to specify conditional probabilities $P(a \mid a)$, $P(b \mid a)$, $P(a \mid b)$ and $P(b \mid b)$, and we will measure the performance of the predictor $P$ by the supremum of

$$D(Q \parallel P) = \sum Q(i,j) \log \frac{Q(i \mid j)}{P(i \mid J)}$$

where the supremum is taken over all i.i.d. over $\mathbb{B}^2$. To use our general model put $\mathbb{E} = \mathbb{B}$, $\mathbb{A} = [0; 1]$ and $\Gamma = M_+^1(\mathbb{A})$. For $q \in [0; 1]$ $\mathcal{M}$ is given by

$$\mathcal{M}(a,a) = q^2, \mathcal{M}(a,b) = q(1-q), \mathcal{M}(b,a) = (1-q)q, \mathcal{M}(a,a) = (1-q)^2$$

If $Q$ is given by a number $q = Q(a)$ then

$$
\begin{aligned}
& D\left(Q \parallel P\right) \\
& = q^2 \log \frac{q}{P\left(a \mid a\right)} + q\left(1-q\right) \log \frac{q}{P\left(a \mid b\right)} \\
& + \left(1-q\right)q \log \frac{1-q}{P\left(b \mid a\right)} + \left(1-q\right)^2 \log \frac{1-q}{P\left(b \mid b\right)} \\
& = \left\langle \log \frac{1}{P\left(i \mid j\right)}, Q \right\rangle - H\left(q, 1-q\right) = \left\langle \log \frac{1}{P\left(i \mid j\right)}, Q \right\rangle + c\left(q\right)
\end{aligned}
$$

where $c\left(q\right) = H\left(q, 1-q\right)$. We state that the predictor given by

$$
P\left(a \mid a\right) = \frac{4}{5},\ P\left(b \mid a\right) = \frac{1}{5},\ P\left(a \mid b\right) = \frac{1}{5},\ P\left(b \mid b\right) = \frac{4}{5}
$$

is optimal. For this predictor

$$
\begin{aligned}
D\left(Q \parallel P\right) &= \left(q^2 + \left(1-q\right)^2\right) \log \frac{1}{4} + \log 5 - H\left(q, 1-q\right) \\
&\leq \log \frac{5}{4}
\end{aligned}
$$

with equality for $q \in \left\{0, \frac{1}{2}, 1\right\}$. For the inequality see [4]. Using the Kuhn-Tucker conditions we just have to show that $P$ is induced by a mixture of probability distributions $Q$ with $q = 0$, $q = \frac{1}{2}$ and $q = 1$.

# References

[1] M. Chiang and T. Cover. Unified duality between channel capacity and rate distortion with state information. In *Proceedings 2001 IEEE Int. Symp. on Inform. Theory*, page 301, 2001.

[2] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. Wiley, New York, 1991.

[3] P. Harremoës and F. Topsøe. Maximum entropy fundamentals. *Entropy*, (3):191 – 226, Sept. 2001. ONLINE http://www.mdi.org/entropy/.

[4] Peter Harremoës and Flemming Topsøe. Information diagrams: Entropy, index of coincidence and probability of error. In *Proceedings IEEE International Symposium on Information Theory*, page 43. IEEE, June 2001.

[5] P. T. Landsberg. *Thermodynamics and Statistical Mechanics*. Dover, New York, 1990.

[6] Flemming Topsøe. Information theoretic optimization technics. *Kybernetika*, 15(1), 1979.