LOW-COMPLEXITY DOWNLINK BEAMFORMING FOR MAXIMUM SUM CAPACITY

Goran Dimić

Dept. of ECE, Univ. of Minnesota, Minneapolis MN 55455, U.S.A. E-mail: goran@ece.umn.edu

ABSTRACT

The problem of simultaneous multiuser downlink beamforming has recently attracted significant interest in both the Information Theory and Signal Processing communities. The idea is to employ a transmit antenna array to create multiple 'beams' directed towards the individual users, and the aim is to increase throughput, measured by sum capacity. Optimal solutions to this problem require convex optimization and so-called Dirty Paper (DP) precoding for known interference, which are prohibitively complex for actual online implementation at the base station. Motivated by recent results by Viswanathan et al and Caire and Shamai, we propose a computationally simple user selection method coupled with zero-forcing beamforming. Our results indicate that the proposed method attains a significant fraction of sum capacity, and thus offers an attractive alternative to DP-based schemes.

1. INTRODUCTION

Depending on whether or not Channel State Information (CSI) is available at the transmitter, transmit antenna arrays can be utilized in two basic ways or a combination thereof: space-time coding, and spatial multiplexing. The former can be used without CSI at the transmitter, and allows mitigation and exploitation of fading. The latter requires CSI at the transmitter, but in turn allows for much higher throughput. Until recently, transmit beamforming was mostly considered for voice services in the context of the cellular downlink. With the emergence of 3G and 4G systems, higher emphasis is being placed on packet data, which are more delay-tolerant but require much higher throughput. Hence the recent interest in transmit beamforming strategies for the cellular downlink that aim for attaining the sum capacity of the wireless channel [1, 8, 9, 4, 6, 7, 5].

Nicholas D. Sidiropoulos*

Dept. of ECE, Technical Univ. of Crete, Chania - Crete, 73100, GREECE E-mail: nikos@telecom.tuc.gr

The scenario of interest can be modeled as a non-degraded Gaussian broadcast channel (GBC). Let N be the number of antennas at the transmitter (Base Station (BS) in a cellular context), and consider a cluster of M mobile users, each equipped with a single receive antenna. The channel between each transmit and receive antenna is constant over a certain time interval and known at the BS. The received signal is corrupted by AWGN independent across users. The BS may transmit simultaneously, using multiple transmit beams, to more than one user in the cluster.

Since the receivers cannot cooperate, successful transmission critically depends on the transmitter's ability to simultaneously send independent signals with as small interference between them as possible. Caire and Shamai [1] proposed a multiplexing technique based on coding for known interference, known as "Writing on Dirty Paper" or Costa precoding [2]. In [2], it is proven that in an AWGN channel with additional additive Gaussian interference, which is known at the transmitter in advance (non-causally), it is possible to achieve the same capacity as if there were no interference. Assuming Costa precoding and known channels at the transmitter, Vishwanath et al. [6] and Yu and Cioffi [9] have proposed algorithms that evaluate sum capacity of the GBC along with the associated optimal signal covariance matrix. However, both approaches require convex optimization in (order of) MN variables to find the optimal signal covariance matrix.

The complexity of the proposed optimization algorithms makes them unsuitable for actual implementation at the BS. A reduced-complexity suboptimal solution to sum rate maximization is proposed in [1]. It suggests the use of QR decomposition of the channel matrix combined with dirty paper (DP) coding at the transmitter. The combined approach nulls interference between data streams, and hence, it is named zero-forcing dirty-paper (ZF-DP) precoding. If $N \ge M$, ZF-DP is proven to be asymptotically optimal at both low and high SNR, but suboptimal in general; whereas zero-forcing (ZF) beamforming without DP coding is optimal in the low SNR regime and yields the same slope of throughput versus SNR in decibels as the sum capacity curve at high SNR. If N < M, [1] has shown that random

^{*}Research supported in part by the European Research Office (ERO) of the US Army under Contract N62558-03-C-0012, and in part by the Army Research Laboratory under Cooperative Agreement DADD19-01-2-0011. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of ERO and ARL of the US Army.

selection of $U \leq N$ users incurs throughput loss for both ZF-DP and ZF. Tu and Blum [5] have proposed a selection algorithm that capitalizes on multiuser diversity, thus increasing the throughput of ZF-DP precoding, and significantly narrowing the gap between ZF-DP throughput and capacity.

An important shortcoming of DP coding is that it requires vector coding and a long temporal block length to be well-approximated in practice; furthermore, with current state-of-art, such approximation entails high computational complexity [3, 8, 10]. For this reason, we advocate herein a more pragmatic approach, based on plain ZF beamforming coupled with a new user selection method. Our approach is applicable in the practically important case that the number of users exceeds the number of transmit antennas. Our simulation results indicate that, at moderate and high SNR, the proposed approach has equal slope of throughput versus SNR as the capacity curve, and it achieves a significant fraction of capacity for all SNR.

ZF beamforming without DP coding was also considered by Spencer and Haardt [4], but they did not consider user selection when M > N. Viswanathan et al. [7] have compared the performance of ZF versus ZF-DP, using a simpler user selection scheme that schedules the N users with the highest *individual* SINR. Under this simpler scheme, they reported that ZF is close to ZF-DP in terms of throughput. Our results further qualify [7], showing that the same is true under a more sophisticated user selection strategy that directly aims to optimize sum capacity. Furthermore, we show that with this new user selection strategy ZF comes close to attaining sum capacity.

2. ZERO-FORCING BEAMFORMING AND USER SELECTION STRATEGY

Let $h_{m,n}$ model the quasi-static, flat-fading channel between transmit antenna n and the receive antenna of user m, and denote $\mathbf{h}_m := [h_{m,1} \ h_{m,2} \ \dots \ h_{m,N}]$. Similarly, let $\mathbf{w}_m = [w_{1,m} \ w_{2,m} \ \dots \ w_{N,m}]^T ((\cdot)^T$ denotes transpose) be the beamforming weight vector for user m. Thus the channel matrix, \mathbf{H} , and the beamforming weight matrix, \mathbf{W} , are

$$\mathbf{H} = \begin{bmatrix} \mathbf{h}_1^* \mathbf{h}_2^* \cdots \mathbf{h}_M^* \end{bmatrix}^* \\ \mathbf{W} = \begin{bmatrix} \mathbf{w}_1 \ \mathbf{w}_2 \ \cdots \ \mathbf{w}_M \end{bmatrix},$$
(1)

where $(.)^*$ denotes conjugate-transpose. Collecting the baseband-equivalent outputs, the received signal vector is

$$\mathbf{x} = \mathbf{HWDs} + \mathbf{n} \tag{2}$$

where s is the transmitted signal vector containing uncorrelated unit-power entries,

$$\mathbf{D} = \begin{bmatrix} \sqrt{p_1} & 0 & \cdots & 0 \\ 0 & \sqrt{p_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sqrt{p_M} \end{bmatrix}$$
(3)

accounts for power-loading and n is the noise vector. Note that the elements of x are physically distributed across the M mobile terminals. Multiuser decoding is therefore not feasible, hence each user treats the signals intended for other users as interference. Noise is assumed to be circular complex Gaussian, zero-mean, uncorrelated with variance of each complex entry $\sigma^2 = 1$.

The desired signal power received by user m is given by $|\mathbf{h}_m \mathbf{w}_m|^2 p_m$. The Signal to Interference plus Noise Ratio (SINR) of user m is

$$SINR_m = \frac{|\mathbf{h}_m \mathbf{w}_m|^2 p_m}{\sum_{i \neq m} |\mathbf{h}_m \mathbf{w}_i|^2 p_i + \sigma^2},$$
(4)

The problem of interest can now be formulated as

$$\max_{\mathbf{W}} \sum_{m=1}^{M} log(1 + SINR_m),$$

subject to: $||\mathbf{WD}||_F^2 \le P,$ (5)

where $||.||_F^2$ denotes Frobenius norm and P stands for a bound on average transmitted power.

Attaining capacity requires Gaussian signaling and long codes, yet the logarithmic SINR reward can be motivated from other, more practical perspectives as well: it can be shown that it measures the throughput of QAM-modulated systems over both AWGN and Rayleigh fading channels. The intuition is that SINR improvements eventually yield diminishing throughput returns.

ZF beamforming inverts the channel matrix at the transmitter, so that orthogonal channels between transmitter and receivers are created. It is then possible to encode users individually, as opposed to more complex long-block-vector coding needed to implement DP. Note that ZF at the transmitter does not enhance noise at the receiver. If the number of users, $M \leq N$, and $rank(\mathbf{H}) = M$, then the ZF beamforming matrix is

$$\mathbf{W} = \mathbf{H}^* (\mathbf{H}\mathbf{H}^*)^{-1}, \tag{6}$$

which is the Moore-Penrose pseudoinverse of the channel matrix. However, if M > N it is not possible to use (6) because **HH**^{*} is singular. In that case, one needs to select $n \leq N$ out of M users.

For M > N, the problem is reformulated as follows: Let $U = \{1, 2, \ldots, M\}$, and $S_n = \{s_u \mid s_u \in U\}$, such that $|S_n| = n$. Given $\mathbf{H} \in \mathbb{C}^{M \times N}$, select $n \leq N$, and a set of channels, $\{\mathbf{h}_{s_1}, \ldots, \mathbf{h}_{s_n}\}$, which produce the row-reduced channel matrix

$$\mathbf{H}(S_n) = \begin{bmatrix} \mathbf{h}_{s_1}^* & \mathbf{h}_{s_2}^* & \cdots & \mathbf{h}_{s_n}^* \end{bmatrix}^*$$
(7)

such that the sum rate is the highest achievable:

$$\max_{1 \le n \le N} \max_{S_n} R_{zf}(S_n)$$

subject to
$$\sum_{i \in S_n} \left[\mu - \frac{1}{c_i(S_n)} \right]_+ = P.$$
 (8)

We define,

$$R_{zf}(S_n) := \sum_{i \in S_n} [log_2(\mu c_i(S_n)]_+, \qquad (9)$$

where $[x]_{+} = max\{0, x\},\$

$$c_i(S_n) = \{ [(\mathbf{H}(S_n)\mathbf{H}(S_n)^*)^{-1}]_{i,i} \}^{-1}, \qquad (10)$$

and μ is obtained by solving the water-filling equation in (8). The power-loading then yields

$$p_i = c_i(S_n) \left[\mu - \frac{1}{c_i(S_n)} \right]_+, \quad \forall i \in S_n.$$
(11)

The problem can be conceptually solved by exhaustive search: for each value of n, find all possible n-tuples S_n and select a pair (n, S_n) which yields maximum $R_{zf}(S_n)$. However, such an algorithm has prohibitive complexity.

We propose a reduced-complexity suboptimal algorithm, dubbed Generalized Zero Forcing (GZF), as outlined next.

1. Initialization:

- Set n = 1.
- Find a user, s_1 , such that $s_1 = \arg \max_{u \in U} \mathbf{h}_u \mathbf{h}_u^*$.
- Set $S_1 = \{s_1\}$ and denote the achieved rate $R_{zf}(S_1)_{max}$.

2. While n < N:

•
$$n = n + 1$$
.

• Find a user, s_n , such that

$$s_n = \arg \max_{u \in U \setminus S_{n-1}} R_{zf}(S_{n-1} \cup \{u\}).$$

- Set S_n = S_{n-1} ∪ {s_n} and denote the achieved rate R_{zf}(S_n)_{max}.
- If $R_{zf}(S_n)_{max} \leq R_{zf}(S_{n-1})_{max}$ break and retain solution $(n-1, S_{n-1})$.
- 3. Beamforming: $W = H(S_n)^* (H(S_n)H(S_n)^*)^{-1}$ Power Loading: Water-filling

2.1. Implementation and Complexity

The most complex task is the evaluation of $R_{zf}(S_{n-1} \cup \{u\})$. From (9), it is split into the evaluation of the $c_i(S_{n-1} \cup \{u\})$'s followed by evaluation of μ . An efficient way to evaluate the $c_i(S_{n-1} \cup \{u\})$'s is by using the matrix inversion lemma to invert the matrix $\mathbf{A}(S_{n-1} \cup \{u\}) := \mathbf{H}(S_{n-1} \cup \{u\})\mathbf{H}(S_{n-1} \cup \{u\})^*$. Note that

$$\mathbf{A}(S_{n-1}\cup\{u\}) = \left[\begin{array}{cc} \mathbf{A}(S_{n-1}) & \mathbf{a}_u \\ \mathbf{a}_u^* & a_{u,u} \end{array}\right]$$

where $\mathbf{a}_u = [\mathbf{h}_{s_1}\mathbf{h}_u^*, \mathbf{h}_{s_2}\mathbf{h}_u^*, \dots \mathbf{h}_{s_{n-1}}\mathbf{h}_u^*]^T$ and $a_{u,u} = \mathbf{h}_u\mathbf{h}_u^*$. Noting that $\mathbf{A}(S_{n-1})^* = \mathbf{A}(S_{n-1})$, and writing

$$\mathbf{q} = \mathbf{A}(S_{n-1})^{-1}\mathbf{a}_u,\tag{12}$$

after some algebraic manipulation we obtain

$$\mathbf{A}(S_{n-1} \cup \{u\})^{-1} = \begin{bmatrix} \mathbf{A}(S_{n-1})^{-1} & \mathbf{0}_{n-1} \\ \mathbf{0}_{n-1}^T & \mathbf{0} \end{bmatrix} + (a_{u,u} - \mathbf{a}_u^* \mathbf{q})^{-1} \begin{bmatrix} \mathbf{q}\mathbf{q}^* & -\mathbf{q} \\ -\mathbf{q}^* & 1 \end{bmatrix},$$
(13)

where $\mathbf{0}_{n-1}^T = [0 \ 0 \ \dots \ 0]_{1 \times (n-1)}$. It can be verified that each time *n* is increased $\mathbf{A}(S_{n-1})^{-1}$ and $a_{i,u}$, $i \in S_{n-2}$, are known before the search over $u \in U \setminus S_{n-1}$ starts. Hence, evaluation of $\mathbf{A}(S_{n-1} \cup \{u\})^{-1}$ from (12) and (13) has complexity proportional to $O(n^2)$.

Given a set S_n , we have [1]

$$c_i(S_n) = |\mathbf{h}_{s_i} \mathbf{P}(S_n \setminus \{s_i\})^{\perp}|^2, \qquad (14)$$

where $\mathbf{P}(S_n)^{\perp}$ denotes the projector onto the orthogonal complement of $\Omega(S_n) = span\{\mathbf{h}_{s_l} : s_l \in S_n\}$. It follows that if (8) and (11) yield $p_u = 0$, then $R_{zf}(S_{n-1} \cup \{u\}) < R_{zf}(S_{n-1})$. We discard such u. We also discard u if (8) and (11) yield $p_{s_i} = 0$ for some $s_i \in S_{n-1}$. This is done to keep complexity at bay, for otherwise combinatorial search might effectively emerge. Hence, user u is a candidate for S_n if $p_i > 0$, $\forall i \in S_{n-1} \cup \{u\}$. From the properties of water-filling, this holds if

$$\frac{n}{c_{i_{m\,in}}(S_{n-1}\cup\{u\})} < P + \sum_{i\in S_{n-1}\cup\{u\}} \frac{1}{c_i(S_{n-1}\cup\{u\})},$$
(15)
(15)

where $c_{i_{min}}(S_{n-1} \cup \{u\}) = \min_{i \in S_{n-1} \cup \{u\}} c_i(S_{n-1} \cup \{u\}).$ Then, we have

$$\mu = \frac{1}{n} \left[P + \sum_{i \in S_{n-1} \cup \{u\}} \frac{1}{c_i(S_{n-1} \cup \{u\})} \right].$$
(16)

If (15) is not satisfied, we skip to the next u. The overall complexity of the algorithm is $O(N^3M)$.

We note that the **break** in Step 2 is necessary when GZF is used, but redundant when ZF-DP is used; it is shown in [1, 5] that in the latter case, maximum sum rate can always be achieved with N active users if P > 0 [1]. On the other hand, when ZF alone is used, the optimum number of active users is $n_{opt} \le N$ and decreases as P decreases, so that for $P \rightarrow 0$, the ZF scheme reduces to maximum ratio combining (MRC), $n_{opt} = 1$ [1]. This also holds for the proposed GZF algorithm, which follows from the water-filling equation in (8) and the fact that $[c_1(S_1)]^{-1} = \max_{i \in U} a_{i,i}$.

3. SIMULATION RESULTS

The performance of the proposed algorithm is presented in Fig. 1. The y-axis shows sum capacity and sum rate in bits per channel use. The x-axis shows total power in dB. Noise level of every user is 1. Sum capacity and sum rates are averaged over 100 channels. Channels are complex-valued, drawn from an i.i.d. Rayleigh distribution with unit-variance for each channel entry. Note that GZF exhibits the same slope of rate increase per dB of SNR as the sum capacity curve at moderate and high SNR. Also note that given N, an increase in M narrows the gap between the sum rate, achieved using GZF, and the sum capacity. This is due to multiuser diversity - the more users that contend for transmission, the higher the probability that N of them will be almost orthogonal. This in turn reduces the advantage of DP-coding based schemes over ZF.



Fig. 1. GZF Performance

4. CONCLUSIONS

We have proposed a low-complexity algorithm for downlink transmission in the GBC for the realistic case wherein the number of users is greater than the number of transmit antennas. We have evaluated the throughput performance of the new algorithm via simulations. The results show that ZF beamforming with the proposed user selection method achieves a significant fraction of sum capacity, at a low complexity cost. The simulation results indicate that GZF achieves the same slope of throughput per dB of SNR as the capacity-achieving strategy based on the use of DP coding for known interference cancellation and convex optimization. Due to its simplicity, low complexity, and close to optimal performance, the proposed method offers an attractive alternative to earlier DP-based methods.

5. REFERENCES

- G. Caire and S. Shamai (Shitz), "On the Achievable Throughput of a Multi-Antenna Gaussian Broadcast Channel," in *IEEE Trans. on Info. Theory*, vol. 49, no. 7, July 2003, pp. 1691–1706
- [2] M. H. M. Costa, "Writing on Dirty Paper," *IEEE Trans. on Info. Theory*, vol. IT-29, no. 3, May 1983.
- [3] C.B. Peel, "On Dirty Paper Coding", *Signal Processing Magazine*, May 2003, pp. 112-113.
- [4] Q. Spencer and M. Haardt, "Capacity and Downlink Transmission Algorithms for a Multi-user MIMO Channel," in *Proc. Of the 36th Asilomar Conf. On Sign. Syst. And Comp.*, Pacific Grove, CA, Nov. 2002.
- [5] Z. Tu and R. S. Blum, "Multiuser Diversity for a Dirty Paper Approach," *IEEE Comm. Letters*, vol. 7, no. 8, Aug. 2003, pp. 370–372
- [6] S. Vishwanath, N. Jindal and A. Goldsmith, "Duality, Achievable Rates and Sum-Rate Capacity of Gaussian MIMO Broadcast Channels," *IEEE Trans. on Info. Theory.*, vol. 49, no. 10, Oct. 2003, pp. 2658–2668
- [7] H. Viswanathan, S. Venkatesan and H. Huang, "Downlink Capacity Evaluation of Cellular Networks with Known Interference Cancellation," *IEEE J. on Sel. Areas in Comm.*, vol. 21, no. 5, June 2003, pp. 802–811
- [8] W. Yu and J. M. Cioffi, "Trellis Precoding for the Broadcast Channel," in *Proc. of Globecom 2001*, San Antonio, TX, November 2001.
- [9] W. Yu and J. M. Cioffi, "Sum Capacity of a Gaussian Broadcast Channel," in *Proc. of IEEE Int. Symp. on Inform. Theory*, ISIT 2002, Lausanne, Switzerland, July 2002.
- [10] R. Zamir, S. Shamai (Shitz), and U. Erez, "Nested Linear/Lattice Codes for Structured Multiterminal Binning," *IEEE Trans. on Inform. Theory*, vol. 48, no 6., June 2002, pp. 1250–1276