Online optimization and machine learning: Applications to resource allocation problems in wireless networks

E. Veronica Belmega^{1,2}

¹Université Gustave Eiffel, CNRS, LIGM ²ETIS, CY Cergy-Paris University, ENSEA, CNRS

January 6th, 2024



Collaborators and students

P. Mertikopoulos (CNRS, LIG)
I. Fijalkow (ENSEA, ETIS)
A. Savard (IMT Nord Europe)
R. Negrel (ESIEE Paris, LIGM)
M. Debbah (Khalifa Univ. of Sci. and Technol.)

A. Marcastel (PhD)H. El Hassani (PhD)I. Chafaa (PhD)O. Bilenne (PostDoc)

Supported by

Chair IoT Orange with Univ. Cergy Pontoise foundation, ANR ORACLESS, ANR ELIOT, GdR ISIS, ENSEA, ...

Based on



E.V. Belmega, P. Mertikopoulos, and R. Negrel, "Online convex optimization in wireless networks and beyond: The feedback - performance trade-off", *invited paper at RAWNET intl. workshop in conjunction with WiOpt*, Sep. 2022.

I. Chafaa, R. Negrel, E.V. Belmega, and M. Debbah, "Self-supervised deep learning for mmWave beam steering exploiting sub-6 GHz channels", *IEEE Trans. on Wireless Commun.*, Mar. 2022.

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024



I. Online resource optimization policies: feedback vs. performance

- Preliminaries
- First order (gradient) feedback
- Zeroth-order (scalar) feedback
- 1-bit feedback

II. Deep learning vs. online policies

• mmWave beamforming from sub-6GHz channels

I. Online resource optimization policies: feedback vs. performance

- Preliminaries
- First order (gradient) feedback
- Zeroth-order (scalar) feedback
- 1-bit feedback

II. Deep learning vs. online policies

• mmWave beamforming from sub-6GHz channels

Distributed and dynamic wireless networks

B5G, 6G, IoT

[Saad'19, Tataria'21]

- Characteristics and requirements: dense, heterogeneous, autonomous, decentralized, energy-efficient, ...
- Challenges

Highly *mobile* radios & networks Unpredictable and arbitrary connectivity patterns Limited network/channel knowledge (*potentially outdated*) Limited processing capabilities

Obj: Design **energy-efficient** resource allocation policies coping with **arbitrary** network dynamics and **scarce** feedback

Toy example



Rk: Arbitrary time-varying objective, unknown at the transmission instant

E. V. Belmega

Télécom Paris - Dép. ICE



 How to optimize an unknown L(**p**, ω_t) objective at the decision time t?



Online iterative process

For t = 1 to T

- Choose policy $\mathbf{p}(t)$
- Incur loss $L(\mathbf{p}(t), \omega_t)$
- Receive feedback
- Update: $\mathbf{p}(t+1) \leftarrow \mathbf{p}(t)$ based on feedback

- Goal: develop efficient online processes based on strictly causal feedback
- Link with machine learning: multi-armed bandits (MABs) reinforcement learning
- Assumptions: on the objective function $L(\mathbf{p}, \omega_t)$ w.r.t. \mathbf{p} but **not** on the underlying dynamics of ω_t

· Ideal benchmark

$$\forall t, L_t(\mathbf{p}(t)) - \min_{\mathbf{q} \in \mathcal{P} \atop \mathbf{dynamic optimal loss}} L_t(\mathbf{q})$$

Too ambitious under the worse-case assumption of an arbitrarily time-varying network!

• Less ambitious benchmark: regret [Hannan'57]



• **Rk:** both benchmarks require knowledge of the future!

Objective: design online policies $\mathbf{p}(t)$ that lead to **no regret**:

 $\limsup_{T \to \infty} \frac{1}{T} \operatorname{Reg}_T \leq 0.$

No-regret policies

 $\begin{array}{ll} \mbox{Static convex optimization} & L_t(\mathbf{p}) = L(\mathbf{p}), \forall t \\ \mathbf{p}(t) \mbox{ no-regret policy} \Rightarrow \overline{\mathbf{p}}(T) = \frac{1}{T} \sum_{t=1}^{T} \mathbf{p}(t) \mbox{ converges to the optimal solution} \end{array}$

Stochastic convex optimization $L_t(\mathbf{p}) = L(\mathbf{p}, \omega_t), \forall t$ $\mathbf{p}(t)$ no-regret policy $\Rightarrow \overline{\mathbf{p}}(T) = \frac{1}{T} \sum_{t=1}^{T} \mathbf{p}(t)$ converges to the optimal solution

Non-cooperative games

[Viossat'13, Mertikopoulos'19]

A no-regret policy (cumulative) converges to the Nash equilibrium in two-player (discrete or convex) **zero-sum** games, potential (discrete or convex) games, ...



Adversarial (non stochastic) MABs

- S arms of slot machines
- Variable: probability distribution $\mathbf{p}(t)$

$$\mathbf{p}(t) \in \mathcal{S} \triangleq \left\{ \mathbf{p} \in \mathbb{R}^{S}_{+} : \sum_{s=1}^{S} p_{s} = 1 \right\}$$

- **Objective:** maximize the linear expected gains $U_t(\mathbf{p}(t)) = \mathbf{r}(t)^T \mathbf{p}(t)$ r(t) - vector of arms' rewards at time t
- **Optimal policy**: exponential/multiplicative weights [Auer'02]

Our online power allocation problem

- S channels or frequency subcarriers
- Variable: power allocation vector $\mathbf{p}(t)$

$$\mathbf{p}(t) \in \mathcal{P} \triangleq \left\{ \mathbf{p} \in \mathbb{R}^{S}_{+} : \sum_{s=1}^{S} p_{s} \leq P_{\max} \right\}$$

- **Objective:** minimize the convex loss $L_t(\mathbf{p}(t))$ tradeoff power consumption vs. rate
- Hint: adapted version of exponential weights



I. Online resource optimization policies: feedback vs. performance

- Preliminaries
- First order (gradient) feedback
- Zeroth-order (scalar) feedback
- 1-bit feedback

II. Deep learning vs. online policies

• mmWave beamforming from sub-6GHz channels

Gradient-based online policies

Low complexity, parallel computing

Online Gradient Descent (OGD)

- Gradient feedback: $\mathbf{v}(t) = -\nabla L_t(\mathbf{p}(t))$
- Update: Euclidean projection of $\mathbf{p}(t) + \gamma \mathbf{v}(t)$

$$\mathbf{p}(t+1) \triangleq \arg \min_{\mathbf{q} \in \mathcal{P}} \frac{1}{2} \|\mathbf{p}(t) + \gamma \mathbf{v}(t) - \mathbf{q}\|^2$$
$$= \arg \min_{\mathbf{q} \in \mathcal{P}} \left\{ \gamma \mathbf{v}(t)^T (\mathbf{p}(t) - \mathbf{q}) + \frac{1}{2} \|\mathbf{p}(t) - \mathbf{q}\|^2 \right\}$$

Based on the mirror descent [Nemirovski'83]

Online Mirror Descent (OMD)

[Shalev-Shwartz'07]

- Gradient feedback: $\mathbf{v}(t) = -\nabla L_t(\mathbf{p}(t))$
- Update: Mirror-mapping on $\mathbf{p}(t) + \gamma \mathbf{v}(t)$

$$\mathbf{p}(t+1) \triangleq \arg\min_{\mathbf{q}\in\mathcal{P}} \left\{ \gamma \mathbf{v}(t)^T (\mathbf{p}(t) - \mathbf{q}) + D_h(\mathbf{q}, \mathbf{p}(t)) \right\}$$

 $D_h(\mathbf{q}, \mathbf{p}) \triangleq h(\mathbf{q}) - h(\mathbf{p}) - \nabla h(\mathbf{p})^T (\mathbf{q} - \mathbf{p})$ the Bregman divergence of regularizer $h(\cdot)$

E. V. Belmega

Jan. 6th, 2024



[Zinkevich'03]

Online Mirror Descent (OMD)

- Particular cases
 - OGD: Euclidean regularizer $h(\mathbf{p}) = \frac{1}{2} \|\mathbf{p}\|^2$, $D_h(\mathbf{q}, \mathbf{p}) = \frac{1}{2} \|\mathbf{q} \mathbf{p}\|^2$
 - Exponential weights for MAB:

entropic regularizer
$$h(\mathbf{p}) = \sum_{s=1}^{S} p_s \log p_s$$
, Kullback-Leibler $D_h(\mathbf{q}, \mathbf{p}) = \sum_{s=1}^{S} q_s \log \frac{q_s}{p_s}$

• MAB (simplex) regret performance: scalability

• OGD:
$$Reg_T = \mathcal{O}\left(\sqrt{T S}\right)$$
 [Zinkevich'03]

• Exponential weights:
$$Reg_T = \mathcal{O}\left(\sqrt{T\log S}\right)$$
 [Auer'02]

OMD allows to design algorithms that are tailored to the feasible set compared with OGD.

E. V. Belmega

Télécom Paris - Dép. ICE

OXL based on gradient feedback

Main idea: we exploit the similarity of \mathcal{P} with the simplex and the entropic regularizer

Online exponential learning (OXL)

- Gradient feedback: $\mathbf{v}(t) = -\nabla L_t(\mathbf{p}(t))$
- Update:

$$\begin{aligned} \mathbf{y}(t+1) &= \mathbf{y}(t) + \gamma \mathbf{v}(t) & \text{cumulative gradient score} \\ p_s(t+1) &= P_{\max} \frac{\exp(y_s(t+1))}{1 + \sum_{i=1}^{S} \exp(y_i(t+1))}, \ \forall s & \text{exponential mapping} \end{aligned}$$

- · Similar to the exponential weights for MAB
- Closed-form update as opposed to OGD, which requires an Euclidean projection
- We proved: no-regret property, $Reg_T \leq P_{\max} \sqrt{2V \log(1+S)T}$, $V = \max \|v(t)\|^2$

Extends to the imperfect gradient feedback case: $\tilde{\mathbf{v}}(t) = \mathbf{v}(t) + \mathbf{error}$

A. Marcastel, E. V. Belmega, P. Mertikopoulos, and I. Fijalkow, "Online power optimization in feedback-limited, dynamic and unpredictable IoT networks", IEEE Trans. on Signal Processing, 2019.

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024

I. Online resource optimization policies: feedback vs. performance

- Preliminaries
- First order (gradient) feedback
- Zeroth-order (scalar) feedback
- 1-bit feedback

II. Deep learning vs. online policies

• mmWave beamforming from sub-6GHz channels

Can the feedback be reduced to a scalar?

- Received feedback: value of incurred loss $L_t(\mathbf{p})$ as in MABs with bandit feedback
- Stochastic gradient approximation: v
 ˜(t) = S/δ L_t(**p** + δ**u**) **u**, **u** - uniformly drawn over the unit S-dimensional Sphere
- To estimate $\nabla L_t(\mathbf{p}(t))$, transmit at $\tilde{\mathbf{p}}(t) = \mathbf{p}(t) + \delta \mathbf{u} \longrightarrow \tilde{\mathbf{p}}(t)$ may fall outside \mathcal{P} !
- Our solution: shrink the feasible set s.t. $\forall \mathbf{p}(t) \in \mathcal{P}_{\delta} \subset \mathcal{P} \Rightarrow \tilde{\mathbf{p}}(t) \in \mathcal{P}$



A. Marcastel, E. V. Belmega, P. Mertikopoulos, and I. Fijalkow, "Online power optimization in feedback-limited, dynamic and unpredictable IoT networks", IEEE Trans. on Signal Processing, 2019.

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024

16/30

[Spall'97, Flaxman'05]

Our modified OXL₀

OXL₀ algorithm

- Transmit at $\tilde{\mathbf{p}}(t) = \mathbf{p}(t) + \delta \mathbf{u}$
- Receive scalar feedback: $L_t(\tilde{\mathbf{p}}(t))$
- Gradient estimation: $\tilde{\mathbf{v}} = -\frac{S}{\delta} L_t(\tilde{\mathbf{p}}(t)) \mathbf{u}$
- Update:

$$\mathbf{y}(t+1) = \mathbf{y}(t) + \gamma \tilde{\mathbf{v}}(t)$$
 cumulative score

$$p_s(t+1) = \delta + P_{\max}(1 - C_{\delta}) \frac{\exp(y_s(t+1))}{1 + \sum_{i=1}^{S} \exp(y_i(t+1))}, \quad \forall s \quad \text{modified exponential mapping}$$

 $C_{\delta} = \frac{\delta}{P_{\max}} \left(S + \sqrt{S}\right)$

- We proved: no-regret property, $EReg_T = \mathcal{O}(T^{3/4})$
- Zeroth-order feedback greatly impacts the decay rate of the average regret

Tradeoff: regret decay rate vs. amount of feedback!

A. Marcastel, E. V. Belmega, P. Mertikopoulos, and I. Fijalkow, "Online power optimization in feedback-limited, dynamic and unpredictable loT networks", IEEE Trans. on Signal Processing, 2019.

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024

Impact of reducing the feedback

OXL vs. OXL₀



 OXL_0 and S - problem dimension

Zeroth-order information reduces the decay rate of the average regret. The decay rate becomes worse by increasing S - the problem dimensionality.

COST-HATA wireless channels, M = 10 transmitting devices, $S \in \{1, 2, 4\}$ bands, N = 1 receiver, $P_{\max} \in [0.5, 2]$, $R_{\min} \in [0.5, 3]$ bps/Hz, $\lambda \in [0.5, 10]$, average regret over 200 random draws of **u**

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024

Zeroth-order feedback with callbacks

- Issues: slow decay rate \$\mathcal{O}(T^{3/4})\$, poor scaling with problem dimensionality
 - ---> MXL0 performs poorly in MIMO networks
- Motivation: exponential weights yield $\mathcal{O}(\sqrt{T})$ regret in bandit feedback MABs
- Our idea: exploit the current feedback R_t jointly with the past R_{t-1}
 - \longrightarrow improved two-point gradient estimator
- For static convex objectives, we can recover the $\mathcal{O}(\sqrt{T})$ regret as with perfect gradient feedback

K-user MIMO multiple access channel

$$R(\mathbf{Q}) = \log \det \left(\mathbf{I} + \sum_{k=1}^{K} \mathbf{H}_{k} \mathbf{Q}_{k} \mathbf{H}_{k}^{\dagger} \right)$$

$$\mathbf{Q}_k \succeq 0, \quad \operatorname{Tr}(\mathbf{Q}_k) \leq P_{\max}$$



Arbitrary static channels, K = 20 users, 3 transmit and 16 receive antennas

O. Bilenne, P. Mertikopoulos, and E.V. Belmega, "Fast Gradient-Free Optimization in Distributed Multi-User MIMO Systems", IEEE Trans. on Signal Processing, 2020.

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024

I. Online resource optimization policies: feedback vs. performance

- Preliminaries
- First order (gradient) feedback
- Zeroth-order (scalar) feedback
- 1-bit feedback

II. Deep learning vs. online policies

• mmWave beamforming from sub-6GHz channels

Can the feedback be reduced to one bit?

- Our idea: exploit MABs with outage (QoS) based rewards and ACK/NACK-type feedback
- Problems: beam alignment in mmWave networks, NOMA with no CSIT/CDIT
- Issue: feasible set needs to be quantized
 - \longrightarrow performance loss



$$\text{GEE}(i, \mathbf{p}) \triangleq \frac{(R_{\min,i} + R_{\min,j}) (1 - \mathbb{P}_{\text{out}}(i, \mathbf{p}))}{p_i + p_j + P_c}$$



H. El Hassani, A. Savard, and E.V. Belmega, "Adaptive NOMA in time-varying wireless networks with no CSIT/CDIT relying on a 1-bit feedback", IEEE Wireless Commun. Lett., 2021.

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024

- Online optimization and no-regret learning: a suitable framework for resource allocation problems in dynamic and unpredictable networks
- Assumptions: Convex (extends to fractional programs) and Lipschitz objectives $L(\mathbf{x}, \omega_t)$ w.r.t. \mathbf{x} , convex feasible sets

\Rightarrow Adaptive online algorithms

Decentralized and reinforcing Cope with **arbitrary network dynamics** Rely on **strictly causal** information Imperfect and **reduced feedback** Regret-based theoretical **guarantees**

- Tradeoff: performance vs. available feedback
- Applies to semi-definite programming: in massive MIMO networks reducing the (matrix) gradient feedback is crucial



- Zeroth-order feedback with callbacks for **online** problems
- One-bit feedback policies: improve the performance, NOMA cope with continuous and discrete variables
- Beyond convex, Lipschitz objectives: more realistic (non convex) energy-efficiency measures
- Compare online algorithms with deep learning ones



I. Online resource optimization policies: feedback vs. performance

- Preliminaries
- First order (gradient) feedback
- Zeroth-order (scalar) feedback
- 1-bit feedback

II. Deep learning vs. online policies

• mmWave beamforming from sub-6GHz channels



Model-based approaches tradeoff

simple (unrealistic) but tractable problems vs. practical but not tractable problems

 \implies machine (deep) learning has the potential to bridge this gap

Deep learning based on neural networks

- + Powerful toolbox: *universal approximators* of complex relationships
- + Pervasive in 6G network design, operation, ...
- Large and relevant training datasets
- High computing/processing capabilities

[Zappone'19, Tataria'21]



mmWave communications

- Exploit high frequency spectrum
- MIMO systems and beam alignment against severe pathloss
- **Challenges:** channel estimation, device mobility and network dynamics
 - \longrightarrow online optimization (MABs) and deep learning [Hashemi'18, Alrabeiah'20]
- Idea: Mapping sub-6GHz channels into mmWave beamforming vectors [Alrabeiah'20]



- Sub-6 GHz channels provide multipath signature, easier to estimate
- $\mathbf{h}^{\mathrm{UL}} \to \mathbf{f}^{\mathrm{DL}}$ highly complex and non linear mapping \longrightarrow deep learning

Télécom Paris - Dép. ICE

Deep learning vs. 1-bit feedback MABs

Network architecture and DeepMIMO

Neural network architecture

- Input: uplink sub-6 GHz channels
- Output: mmWave beamforming vectors
- Fully-connected: structure agnostic

Training

- DeepMIMO dataset: $\left(\{ \mathbf{h}_i[\ell]^{\mathrm{UL}} \}_{\ell=1}^L, \{ \mathbf{h}_i[\ell]^{\mathrm{DL}} \}_{\ell=1}^L \right)_i$
- Custom loss function: $\mathcal{L} = -1/\mathcal{B}\sum_{i=1}^{\mathcal{B}}\mathcal{R}_i$

$$\mathcal{R}_{i} = \frac{1}{L} \sum_{\ell=1}^{L} \log_{2} \left(1 + \frac{P^{\mathrm{DL}}}{L \left(\sigma^{\mathrm{DL}}\right)^{2}} \mid \mathbf{h}_{i}^{\mathrm{DL}\dagger}[\ell] \mathbf{f}_{i} \mid^{2} \right)$$

 \mathbf{f}_i : predicted beamformer, P^{DL} transmit power, \mathcal{B} size of mini-batch, L subcarriers, $(\sigma^{\mathrm{DL}})^2$ noise variance





DeepMIMO setup [Alkhateeb'19]

27/30

I. Chafaa, R. Negrel, E.V. Belmega, and M. Debbah, "Unsupervised deep learning for mmWave beam steering exploiting sub-6 GHz channels", IEEE Trans. on Wireless Commun., 2022.

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024

Deep learning vs. 1-bit feedback MABs

Deep learning or MABs?

MABs

- + No *offline* training, online learning
- + Low complexity $\sim 10^2$ ops/iteration
- I-bit feedback (ACK/NACK)
- Quantized beams → performance loss
- Trial and error → transitory phase

Deep learning

- Offline training on large and relevant datasets, DeepMIMO [Alkhateeb'19]
- Higher complexity $\sim 10^6 10^7$ ops/iteration
- sub-6 GHz channel information
- + Regression (no quantization)
- + Better running performance (no transitory phase)

Depends on the target application, its characteristics and requirements



MABs: EXP3, MEXP3, UCB Deep learning [Alrabeiah'20]: classification

I. Chafaa, R. Negrel, E.V. Belmega, and M. Debbah, "Self-supervised deep learning for mmWave beam steering exploiting sub-6 GHz channels", IEEE Trans. on Wireless Commun., 2022. I. Chafaa, E.V. Belmega, and M. Debbah, "One-bit Feedback exponential learning for beam alignment in mobile mmWave", IEEE Access, Oct. 2020.

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024

- Deep learning approaches capture complex relationships: sub-6GHz uplink channels – downlink mmWave beams
- Extension to multi-cell multi-user networks
- Robustness to data imperfections
- Generalization capability: other wireless settings, problems
- Recurrent networks, reinforcement and online deep learning to capture the temporal network dynamics



AI-enhanced highly mobile and unpredictable IoT networks

- **OBJ-1** Design efficient **online optimization** algorithms requiring **low-cost** and energy-efficient communication feedback
- **OBJ-2** Design online algorithms maximizing **non convex** energy efficiency exploiting non convex online optimization and deep learning

Sustainable wireless communications: low-energy, low-cost and zero added electromagnetic waves

- **OBJ-1** Derive the fundamental **achievable Shannon rates** in multi-user, multi-backscatter/RIS networks
- **OBJ-2** Develop **efficient algorithms** that jointly tune the transmit strategy and the backscattering/RIS strategy

More intel on my webpage: https://sites.google.com/site/evbelmega

E. V. Belmega





O. Bilenne, P. Mertikopoulos, and E.V. Belmega, "Fast Gradient-Free Optimization in Distributed Multi-User MIMO Systems", IEEE Trans. on Signal Processing, Oct. 2020.





- A. Marcastel, E.V. Belmega, P. Mertikopoulos, and I. Fijalkow, "Online power optimization in dynamic IoT networks: The impact of feedback scarcity", *IEEE Trans. on Signal Processing*, Mar. 2019.
- P. Mertikopoulos, and E.V. Belmega, "Learning to be green: robust energy efficiency maximization in dynamic MIMO-OFDM systems", IEEE Journal on Selected Areas in Communication, Apr. 2016.



- H. El Hassani, A. Savard, and E.V. Belmega, "Energy-efficient 1-bit feedback NOMA in wireless networks with no CSIT/CDIT", IEEE SSP Workshop, Jun. 2021.
- H. El Hassani, A. Savard, and E.V. Belmega, "Adaptive NOMA in time-varying wireless networks with no CSIT/CDIT relying on a 1-bit feedback", IEEE Wireless Commun. Lett., Apr. 2021.
- I. Chafaa, R. Negrel, E.V. Belmega, and M. Debbah, "Self-supervised deep learning for mmWave beam steering exploiting sub-6 GHz channels", IEEE Trans. on Wireless Commun., 2022.
- I. Chafaa, E.V. Belmega, and M. Debbah, "One-bit Feedback exponential learning for beam alignment in mobile mmWave", IEEE Access, Oct. 2020.





T. Chen, S. Barbarossa, X. Wang, G. B. Giannakis, and Z.-L. Zhang, "Learning and management for Internet-of-Things: Accounting for adaptivity and scalability", Arxiv preprint arxiv:1810.11613, 2018.



J. Hannan, "Approximation to Bayes risk in repeated play", Contributions to the Theory of Games, Vol. III, fPrinceton University Press, vol. 39, pp. 97–139, 1957.



Y. Viossat and A. Zapechelnyuk, "No-regret dynamics and fictitious play", Journal of Economic Theory, vol. 148, no 2, pp. 825–842, 2013.



M. Zinkevich, "Online convex programming and generalized infinitesimal gradient ascent, ICML'03: Proceedings of the 20th International Conference on Machine Learning, pp. 928–936, 2003.



S. Shalev-Shwartz, "Online learning: Theory, algorithms, and applications," Ph.D. dissertation, Hebrew University of Jerusalem, 2007.



A. S. Nemirovski and D. B. Yudin, "Problem Complexity and Method Efficiency in Optimization", Wiley, New York, NY, 1983.





J. C. Spall, "A one-measurement form of simultaneous perturbation stochastic approximation", Automatica, vol. 33, no. 1, pp. 109–112, 1997.

A. D. Flaxman, A. T. Kalai, and H. B. McMahan, "Online convex optimization in the bandit setting: gradient descent without a gradient", SODA'05: Proceedings of the 16th annual ACM-SIAM symposium on discrete algorithms, 2005, pp. 385–394.

- C. Zhang, P. Patras, and H. Haddadi, "Deep learning in mobile and wireless networking: A survey", Arxiv preprint arxiv:1803.04311, 2018.
- A. Zappone, M. Di Renzo, and M. Debbah, "Wireless networks design in the era of deep learning: Model-based, AI-based, or both?", IEEE Transactions on Communications, vol. 67, no. 10, pp. 7331-7376, 2019.
- E. C. Strinati, S. Barbarossa, J. L. Gonzalez-Jimenez, D. Ktenas, N. Cassiau, and C. Dehos, "6G: The next frontier: From holographic messaging to artificial intelligence using subterahertz and visible light communication", IEEE Vehicular Technology Magazine, vol. 14, no.3, pp.42-50, 2019.
- W. Saad, M. Bennis, and M. Chen, "A vision of 6G wireless systems: Applications, trends, technologies, and open research problems", IEEE Network, vol. 34, no. 3, pp. 134–142, 2019.
- H. Tataria, M. Shafi, A.F. Molisch, M. Dohler, H. Si oland, and F. Tufvesson, "6G wireless systems: Vision, requirements, challenges, insights, and opportunities", Proceedings of the IEEE, vol. 109, no. 7, pp. 1166-1199, 2021.
- M. Alrabeiah and A. Alkhateeb, "Deep learning for mmwave beam and blockage prediction using sub-6 GHz channels", IEEE Trans. Commun., 2020.



- A. Alkhateeb, "DeepMIMO: A generic deep learning dataset for millimeter wave and massive MIMO applications", arXivpreprint arXiv:1902.06435, 2019.
- - X, Song, S, Haghighatshoar, and G, Caire, "A scalable and statistically robust beam alignment technique for mm-Wave systems", IEEE Trans, Wireless Commun., vol. 17, no. 7, pp. 4792-4805, 2018.
 - J.B. Wang, M. Cheng, J. Y. Wang, M. Lin, Y. Wu, H. Zhu, and J. Wang, "Bandit inspired beam searching scheme for mmWave high-speed train communications", arXiv preprint arXiv:1810.06150.2018.
 - M. Hashemi, A. Sabharwal, C.E. Koksal, and N.B. Shroff, "Efficient beam alignment in millimeter wave systems using contextual bandits". IEEE Conference on Computer Communications (INFOCOM), 2018,



• Uplink signal: captured by AP at sub-6 GHz

$$\mathbf{y}^{\mathrm{UL}}[\ell] = \mathbf{h}^{\mathrm{UL}}[\ell] \, x^{\mathrm{UL}}[\ell] + \mathbf{n}^{\mathrm{UL}}[\ell]$$

 $\mathbf{h}^{\mathrm{UL}}[\ell] \in \mathbb{C}^{M \times 1}; \ \mathbb{E}[\mid x^{\mathrm{UL}}[\ell] \mid^2] = P^{\mathrm{UL}}/L; \ \mathbf{n}^{\mathrm{UL}}[\ell] \sim \mathcal{N}(0, (\sigma^{\mathrm{UL}})^2)$

• Downlink signal: received by the user at the mmWave band

$$y^{\mathrm{DL}}[\ell] = \mathbf{h}^{\mathrm{DL}\dagger}[\ell] \mathbf{f} x^{\mathrm{DL}}[\ell] + n^{\mathrm{DL}}[\ell]$$

 $\mathbf{h}^{\mathrm{DL}}[\ell] \in \mathbb{C}^{N \times 1}; \ \mathbf{f} \in \mathbb{C}^{N \times 1}, \|\mathbf{f}\|^2 = 1; \ \mathbb{E}[|\ x^{\mathrm{DL}}[\ell] \ |^2] = P^{\mathrm{DL}}/L; \ n^{\mathrm{DL}}[\ell] \sim \mathcal{N}(0, (\sigma^{\mathrm{DL}})^2)$

• Dataset samples: $({\mathbf{h}_i[\ell]^{\text{UL}}}_{\ell=1}^L, {\mathbf{h}_i[\ell]^{\text{DL}}}_{\ell=1}^L)_i$

L number of OFDM subcarriers, i sample index

- Sub-6 GHz channels capture wireless characteristics, easy to estimate
- Dataset split at each link: 80% for training (15% validation set) and 20% for testing



• System parameters:

- Learning parameters: 100 epochs, ADAM optimizer (10^{-4}) , $\mathcal{B} = 256$
- ▶ Dataset parameters: 4 AP, outdoor scenario: $3.5 \text{ GHz} \rightarrow 28 \text{ GHz}$
- Samples in each link: 108 419, 108 781, 63 350 and 117 650
- ▶ $P^{DL} = 34$ dBm, M = 4, N = 64 and L = 32, mmWave bandwidth 0.5 GHz
- Comparison benchmarks:
 - Centralized learning: same neural network, full access to all data samples
 - Perfect downlink CSI: full and instantaneous knowledge of downlink channels
 - Individual learning: fully distributed, no cooperation between APs
MABs

+		No offline	training,	online	learning	
---	--	------------	-----------	--------	----------	--



+ 1-bit feedback (ACK/NACK)

- Quantized beams → performance loss
- Trial and error → transitory phase

Deep learning

- Offline training on large and relevant datasets, *DeepMIMO* [Alkhateeb'19]
- Higher complexity $\sim 10^6 10^7$ ops/iteration
- sub-6 GHz channel information
- + Regression (no quantization)
- + Better running performance (no transitory phase)



MABs: EXP3, MEXP3, UCB

Deep learning [Alrabeiah'20]: classification

Depends on the target application, its characteristics and requirements

I. Chafaa, R. Negrel, E.V. Belmega, and M. Debbah, "Self-supervised deep learning for mmWave beam steering exploiting sub-6 GHz channels", IEEE Trans. on Wireless Commun., 2022. I. Chafaa, E.V. Belmega, and M. Debbah, "One-bit Feedback exponential learning for beam alignment in mobile mmWave", IEEE Access, Oct. 2020.

E. V. Belmega

Télécom Paris - Dép. ICE



Deep learning or MABs?

MABs

+ No offline training, online learning

+ Low complexity $\sim 10^2$ ops/iteration

+ 1-bit feedback (ACK/NACK)

- Quantized beams → performance loss
- Trial and error → transitory phase

Deep learning

- Offline training on large and relevant datasets, *DeepMIMO* [Alkhateeb'19]
- Higher complexity $\sim 10^6 10^7$ ops/iteration
- sub-6 GHz channel information
- + Regression (no quantization)
- + Better running performance (no transitory phase)



MABs: EXP3, MEXP3, UCB

Deep learning [Alrabeiah'20]: classification

Mobility model: average distance between the AP and user increases in time

Depends on the target application, its characteristics and requirements

I. Chafaa, R. Negrel, E.V. Belmega, and M. Debbah, "Self-supervised deep learning for mmWave beam steering exploiting sub-6 GHz channels", IEEE Trans. on Wireless Commun., 2022. I. Chafaa, E.V. Belmega, and M. Debbah, "One-bit Feedback exponential learning for beam alignment in mobile mmWave", IEEE Access, Oct. 2020.

E. V. Belmega

Télécom Paris - Dép. ICE





• FL approach is efficient in case of local data scarcity



- · Uplink channels are contaminated with different levels of noise
- Variance of uplink channels: 10^{-9}
- SFL outperforms CL and individual learning at high noise regime



• N = 64 and L = 32

• When the input size M increases, the optimal size decreases

E. V. Belmega

Télécom Paris - Dép. ICE





• M = 4 and L = 32.

• When the output size N increases, the optimal size increases as well

E. V. Belmega

Télécom Paris - Dép. ICE

MABs

- UCB attains its ε -optimal performance (ε -regret) after $1/\varepsilon$ iterations (up to logarithmic factors)
- EXP3 and MEXP3 attains it after $1/\varepsilon^2$ iterations
- One iteration of these algorithms scales linearly with the codebook size: $\mathcal{O}(A)$

Deep learning

• Complexity of $\mathcal{O}(LMS + S^2 + SN)$ at each iteration

Entropic regularizer $h(\cdot)$: $D_h(\mathbf{y}, \mathbf{x}) = \sum_j y_j \log \frac{y_j}{x_j}$ (Kullback-Leibler divergence)

Exponential weights (EW / HEDGE / Multiplicative weights)

- Full information (gradient feedback): $\mathbf{v}(t) \equiv \mathbf{r}(t) = \nabla U_t(\mathbf{p}(t))$ (reward vector)
- Recursive update (multiplicative):

$$p_s(t+1) = \frac{p_s(t) \exp(\gamma r_s(t))}{S}, \forall s$$
$$\sum_{i=1}^{S} p_i(t) \exp(\gamma r_i(t))$$

• Equivalent to:

$$(t+1) = \mathbf{y}(t) + \gamma \mathbf{r}(t)$$

cumulative reward score

$$p_s(t+1) = rac{\exp(y_s(t+1))}{S}, \ orall s$$
 exponential mapping $\sum_{i=1}^{S} \exp(y_i(t+1))$

• Regret bound: $Reg_T \leq \sqrt{2T\log S}$ with $\gamma^* = \sqrt{2\log S/T}$

 \mathbf{y}

E. V. Belmega

Jan. 6th, 2024

[Auer'02]

• Convex conjugate of the regularizer $h(\mathbf{p})$ - strongly convex:

$$\begin{aligned} h^*(\mathbf{y}) &= \max_{\mathbf{p} \in \mathcal{S}} \mathbf{y}^T \mathbf{p} - h(\mathbf{p}) \\ &= \log\left(\sum_i \exp(y_i)\right) \end{aligned}$$

• Its gradient is equivalent with the update of EW algorithm $\mathbf{p}(t+1) = \nabla h^*(\mathbf{y}(t+1))$

$$\frac{\partial h^*(\mathbf{y})}{\partial y_i} = \frac{\exp(y_i)}{\sum_k \exp(y_k)}$$

• The update is also equivalent with

$$\mathbf{p}(t+1) = \arg \max_{\mathbf{p} \in S} \mathbf{y}(t)^T \mathbf{p} - h(\mathbf{p})$$

E. V. Belmega

Télécom Paris - Dép. ICE



30/30

Feasible set:
$$\mathcal{P} = \left\{ \mathbf{p} \in \mathbb{R}_{+}^{S} : \sum_{s=1}^{S} p_{s} \leq P_{\max} \right\}$$

Our regularizer $h(\mathbf{p}) = \sum_{i} p_{i} \log p_{i} + (P_{\max} - \sum_{i} p_{i}) \log(P_{\max} - \sum_{i} p_{i})$
 $D_{h}(\mathbf{q}, \mathbf{p}) = \sum_{j} q_{j} \log \frac{q_{j}}{p_{j}} + (P_{\max} - \sum_{k} q_{k}) \log \frac{P_{\max} - \sum_{k} q_{k}}{P_{\max} - \sum_{k} p_{k}}$

Kullback-Leibler type of divergence with additional variables: $q_{s+1} = P_{\max} - \sum_{k=1}^{S} q_k, p_{s+1} = P_{\max} - \sum_{k=1}^{S} p_k$ (the unused power of zero reward)

OMD algorithm

- Gradient feedback: $\mathbf{v}(t) = -\nabla L_t(\mathbf{p}(t))$
- Undate:

$$p_{s}(t+1) = \frac{P_{\max} p_{s}(t) \exp(\gamma v_{s}(t))}{\sum_{i=1}^{S} p_{i}(t) \exp(\gamma v_{i}(t)) + P_{\max} - \sum_{k} p_{k}(t)}, \quad \forall s$$
• Equivalent to OXL:

$$\mathbf{y}(t+1) = \mathbf{y}(t) + \gamma \mathbf{v}(t) \qquad \text{cumulative gradient score}$$

$$p_{s}(t+1) = P_{\max} \frac{\exp(y_{s}(t+1))}{1 + \sum_{i=1}^{S} \exp(y_{i}(t+1))}, \quad \forall s \quad \text{exponential mapping}$$
E. V. Belmega

$$\mathbf{Télécom Paris - Dép. ICE} \qquad Jan. 6th, 2024$$

• Convex conjugate of $h(\mathbf{p})$

$$h^{*}(\mathbf{y}) = \max_{\mathbf{p} \in \mathcal{P}} \mathbf{y}^{T} \mathbf{p} - h(\mathbf{p})$$
$$= P_{\max} \log \left(1 + \sum_{i} \exp(y_{i}) \right) - P_{\max} \log P_{\max}$$

• $h(\mathbf{p})$ is $1/P_{\max}$ - strongly convex wrt $\|\cdot\|_1 \Longrightarrow h^*(\mathbf{y})$ is P_{\max} strongly smooth wrt $\|\cdot\|_{\infty}$

• Its gradient is equivalent with the update of OXL algorithm: $\mathbf{p}(t+1) = \nabla h^*(\mathbf{y}(t+1))$

$$\frac{\partial h^*(\mathbf{y})}{\partial y_i} = P_{\max} \frac{\exp(y_i)}{1 + \sum_k \exp(y_k)}$$

• The update is also equivalent with

$$\mathbf{p}(t+1) = \arg \max_{\mathbf{p} \in \mathcal{P}} \ \mathbf{y}(t)^T \mathbf{p} - h(\mathbf{p})$$

E. V. Belmega

Télécom Paris - Dép. ICE



• First-order convexity condition of $L_t(\cdot)$ (linearization), $\mathbf{v}(t) = \nabla L_t(\mathbf{p}(t))$:

$$Reg_{\mathbf{q}}(T) \leq -\sum_{t=1}^{T} \langle \mathbf{v}(t) | \mathbf{p}(t) - \mathbf{q} \rangle$$

• Cumulative score: $\mathbf{y}(t+1) = \mathbf{y}(t) + \gamma \mathbf{v}(t)$ and $\mathbf{y}(1) = 0$

$$Reg_{\mathbf{q}}(T) \leq -\sum_{t=1}^{T} \left< \mathbf{v}(t) | \mathbf{p}(t) \right> + \frac{1}{\gamma} \left< \mathbf{y}(T+1) | \mathbf{q} \right>$$

• Using $\mathbf{p}(t) = \nabla h^*(\mathbf{y}(t))$ and P_{\max} -strong smoothness of $h^*(\mathbf{y})$

$$h^{*}(\mathbf{y}(t+1)) \leq h^{*}(\mathbf{y}(t)) + \gamma \langle \mathbf{v}(t) | \nabla h^{*}(\mathbf{y}(t)) \rangle + \frac{\gamma^{2}}{2} P_{\max} \| \mathbf{v}(t) \|_{\infty}^{2},$$

$$\operatorname{Reg}_{\mathbf{q}}(T) \leq \frac{1}{\gamma} \left[h^*(0) - h^*(\mathbf{y}(T+1)) \right] + \frac{\gamma}{2} \operatorname{Pmax} \sum_{t=1}^T \|\mathbf{v}(t)\|_{\infty}^2 + \frac{1}{\gamma} \left\langle \mathbf{y}(T+1) | \mathbf{q} \right\rangle.$$

• Fenchel inequality: $h^*(\mathbf{y}(T+1)) + h(\mathbf{q}) \ge \langle \mathbf{y}(T+1) | \mathbf{q} \rangle$

$$\operatorname{Reg}_{\mathbf{q}}(T) \leq \frac{1}{\gamma} \left[h(\mathbf{q}) + P_{\max} \log(1+S) - P_{\max} \log(P_{\max}) \right] + \frac{\gamma}{2} P_{\max} V T.$$

• From $h(\mathbf{q}) \leq P_{\max} \log(P_{\max})$ $Reg(T) \leq \frac{P_{\max} \log(1+S)}{\gamma} + \frac{\gamma}{2} P_{\max} VT$ E. V. Belmega | Télécom Paris - Dép. ICE | Jan. 6th, 2024 | 30/30

Fixed step-size γ



 $V = \max \|\mathbf{v}(t)\|^2$

• If T is known: optimal step-size $\gamma^* = \sqrt{\frac{2\log(1+S)}{TV}}$

$$Reg_T \le P_{\max} \sqrt{2V \log(1+S)T}$$

• If T is not-known: window doubling trick

Size 1 2 4 T $T_m = 2^m$ Window 0 1 2 $m = \lceil \log_2 T \rceil$

$$Reg_T \le \frac{2}{\sqrt{2}-1} P_{\max} \sqrt{2V \log(1+S)T}$$

E. V. Belmega

Télécom Paris - Dép. ICE



Having outdated feedback impacts the waterfilling performance at high channel variability.

Single device time-varying setup: $M = N = 1, S = 4, g(t + 1) = \alpha g(t) + (1 - \alpha)z(t)$ with $z(t) \sim \mathcal{N}(0, \sigma_z^2)$ (average over 90 draws), $R_{\max} = 2, R_{\min} = 3, \lambda = 1$

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024

30/30

• Received feedback: unbiased gradient estimation, $\tilde{\mathbf{v}}(t) = -\nabla L_t(\mathbf{p}(t)) + \text{error}$

$\mathbb{E}\left[\tilde{\mathbf{v}}(t) \mid \mathcal{F}_{t-1}\right]$	=	$-\nabla L_t(\mathbf{p}(t))$	no systematic errors
$\mathbb{E}\left[\ \tilde{\mathbf{v}}(t)\ ^2 \mid \mathcal{F}_{t-1}\right]$	\leq	\widetilde{V}	bounded mean square

 \mathcal{F}_t - history of play up to t

• OXL with
$$\mathbf{v}(t) \leftarrow \tilde{\mathbf{v}}(t)$$
 and an optimal step-size $\gamma^* = \sqrt{\frac{2 \log(1+S)}{T \widetilde{V}}}$
 $\mathbb{E}\left[Reg_T\right] \leq P_{\max} \sqrt{2\widetilde{V} \log(1+S) T}$

Rk: The expected regret is not highly impacted by unbiased errors in the gradient estimation.

• First-order convexity condition (linearization):

$$\mathbb{E}[Reg_{\mathbf{q}}(T)] \leq \mathbb{E}\left[\sum_{t=1}^{T} \langle \nabla L_t(\mathbf{p}(t)) | \mathbf{p}(t) - \mathbf{q} \rangle\right]$$

• Using $\nabla L_t(\mathbf{p}(t)) = -\mathbb{E}[\tilde{\mathbf{v}}(t)|\tilde{\mathbf{v}}(t-1),...,\tilde{\mathbf{v}}(1)]$ and the law of total expectation

$$\mathbb{E}\left[\sum_{t=1}^{T} \langle \nabla L_t(\mathbf{p}(t)) | \mathbf{p}(t) - \mathbf{q} \rangle\right] = -\mathbb{E}\left[\sum_{t=1}^{T} \langle \tilde{\mathbf{v}}(t) | \mathbf{p}(t) - \mathbf{q} \rangle\right]$$

• Using $\mathbb{E}\left[\|\tilde{\mathbf{v}}(t)\|_2^2\right] \leq \tilde{V}$

$$\mathbb{E}[Reg(T)] \le \frac{P_{\max}\log(1+S)}{\mu} + \frac{\mu}{2}P_{\max}T\tilde{V}$$

E. V. Belmega

Télécom Paris - Dép. ICE

Stochastic gradient approximation: $\tilde{\mathbf{v}}(t) = \frac{S}{\delta} L_t(\mathbf{p} + \delta \mathbf{u}) \mathbf{u}$,

[Spall'97, Flaxman'05]

by sampling the loss function not at \mathbf{p} , but at a nearby point \mathbf{u} - uniformly drawn over the unit S-dimensional Sphere

Properties:

$$\begin{split} \mathbb{E}\left[\mathbf{\tilde{v}}\right] &= \nabla \tilde{L}_t(\mathbf{p}) & \text{link with the gradient of } \tilde{L}_t(\mathbf{p}) \\ \\ \tilde{L}_t(\mathbf{p}) &\triangleq \mathbb{E}\left[L_t(\mathbf{p} + \delta \mathbf{u})\right] & \tilde{L}_t(\mathbf{p}) \text{-smooth approximation of } L_t(\mathbf{p}) \\ \\ |L_t(\mathbf{p}) - \tilde{L}_t(\mathbf{p})| &\leq L\delta \end{split}$$

L - the Lipschitz constant of $L_t(\mathbf{p})$

• For fixed γ and δ

$$\mathbb{E}\left[Reg_{T}\right] \leq \frac{P_{\max}\,\log(1+S)}{2\gamma} + \gamma\,\mathbf{T}\,S^{2}\left(\frac{B}{\delta} + L\right)^{2} + L\,\mathbf{T}\,\delta\,\left(3 + P_{\max}\left(S + 2\sqrt{S}\right)\right)$$

 $B = \max L_t(\mathbf{p})$

• Choosing
$$\delta^* = \frac{P_{\max}}{(S+\sqrt{S})T^{1/4}}$$
 and $\gamma^* = \sqrt{\frac{P_{\max}\log(1+S)}{2T}} \left[S\left(\frac{B}{\delta^*} + L\right)\right]^{-1}$, the expected regret is sublinear

$$\mathbb{E}\left[Reg_T\right] = \mathcal{O}(T^{3/4}S^3\sqrt{\log(1+S)})$$

 $\ensuremath{\mathbf{Rk}}\xspace$ The zero-th order feedback greatly impacts the decay rate of the average regret.

 \rightarrow **Tradeoff:** regret decay rate vs. amount of required feedback!

E. V. Belmega

Télécom Paris - Dép. ICE

New regularizer:
$$h(\mathbf{p}) = \sum_{i} (p_{i} - \delta) \log(p_{i} - \delta) + (P_{\max} - \sqrt{S}\delta - \sum_{i} p_{i}) \log(P_{\max} - \sqrt{S}\delta - \sum_{i} p_{i})$$
$$D_{h}(\mathbf{q}, \mathbf{p}) = \sum_{j} (q_{j} - \delta) \log \frac{q_{j} - \delta}{p_{j} - \delta} + (P_{\max} - \sqrt{S}\delta - \sum_{k} q_{k}) \log \frac{P_{\max} - \sqrt{S}\delta - \sum_{k} q_{k}}{P_{\max} - \sqrt{S}\delta - \sum_{k} p_{k}}$$

Kullback-Leibler type of divergence adapted to the shrunk set

OMD algorithm

- Gradient feedback: $\mathbf{v}(t) = -\nabla L_t(\mathbf{p}(t))$
- Update:

$$p_{s}(t+1) = \delta + \frac{P_{\max}(1-C_{\delta}) \left(p_{s}(t)-\delta\right) \exp(\gamma v_{s}(t))}{\sum_{i=1}^{S} \left(p_{i}(t)-\delta\right) \exp(\gamma v_{i}(t)) + P_{\max} - \sqrt{S}\delta - \sum_{k} p_{k}(t)}, \quad \forall s$$

• Equivalent to OXL₀:

$$\begin{aligned} \mathbf{y}(t+1) &= \mathbf{y}(t) + \gamma \tilde{\mathbf{v}}(t) & \text{cumulative score} \\ p_s(t+1) &= \delta + P_{\max}(1-C_{\delta}) \frac{\exp(y_s(t+1))}{1+\sum_{i=1}^{S} \exp(y_i(t+1))}, \quad \forall s \quad \text{modified exponential mapping of } \end{aligned}$$

adapted to the shrunk set \mathcal{P}_{δ}

 $C_{\delta} = \frac{\delta}{P_{\max}}(S + \sqrt{S})$ E. V. Belmega | **Télécom Paris - Dép. ICE** | Jan. 6th, 2024 | 30/30 • Convex conjugate of $h(\mathbf{p})$:

$$h^{*}(\mathbf{y}) = \max_{\mathbf{p} \in \mathcal{P}_{\delta}} \mathbf{y}^{T} \mathbf{p} - h(\mathbf{p})$$
$$= \delta \sum_{i} y_{i} + P_{\max}(1 - C_{\delta}) \log \left(1 + \sum_{i} \exp(y_{i})\right)$$
$$-P_{\max}(1 - C_{\delta}) \log(P_{\max}(1 - C_{\delta}))$$

• Its gradient is equivalent with the update of OXL algorithm: $\mathbf{p}(t+1) = \nabla h^*(\mathbf{y}(t))$

$$\frac{\partial h^*(\mathbf{y})}{\partial y_i} = \delta + P_{\max}(1 - C_{\delta}) \frac{\exp(y_i)}{1 + \sum_k \exp(y_k)}$$

• The update is also equivalent with

$$\mathbf{p}(t+1) = \arg \max_{\mathbf{p} \in \mathcal{P}_{\delta}} \mathbf{y}(t)^T \mathbf{p} - h(\mathbf{p})$$

E. V. Belmega

Télécom Paris - Dép. ICE

• Compare $L_t(\mathbf{p}_{\delta}(t) + \delta \mathbf{u}), L_t(\mathbf{q})$ to $L_t(\mathbf{p}_{\delta}(t)), L_t(\mathbf{q}_{\delta})$ using the L-Lipschitz property (move to the shrunk set)

$$\mathbb{E}[\operatorname{Reg}_{\mathbf{q}}(T)] \leq \mathbb{E}\left[\sum_{t=1}^{T} L_t(\mathbf{p}_{\delta}(t)) - L_t(\mathbf{q}_{\delta})\right] + LT\delta\left(1 + P_{\max}\left(S + 2\sqrt{S}\right)\right)$$

• Compare $L_t(\mathbf{p}_{\delta}(t)), L_t(\mathbf{q}_{\delta})$ to smooth approx. $\tilde{L}_t(\mathbf{p}_{\delta}(t)), \tilde{L}_t(\mathbf{q}_{\delta})$ for linking the regret with $\tilde{v}(t)$

$$\mathbb{E}[\operatorname{Reg}_{\mathbf{q}}(T)] \leq \mathbb{E}\left[\sum_{t=1}^{T} \tilde{L}_{t}(\mathbf{p}_{\delta}(t)) - \tilde{L}_{t}(\mathbf{q}_{\delta})\right] + LT\delta\left(3 + P_{\max}\left(S + 2\sqrt{S}\right)\right)$$

- First-order convexity condition of $\tilde{L}_t(\mathbf{p})$: $\mathbb{E}\left[\sum_{t=1}^T \tilde{L}_t(\mathbf{p}_{\delta}(t)) \tilde{L}_t(\mathbf{q}_{\delta})\right] \leq \mathbb{E}\left[\sum_{t=1}^T \langle \nabla \tilde{L}_t(\mathbf{p}_{\delta}(t)) | \mathbf{p}_{\delta}(t) \mathbf{q}_{\delta} \rangle \right]$
- Using $\nabla \tilde{L}_t(\mathbf{p}_{\delta}(t)) = \mathbb{E}[\tilde{\mathbf{v}}(t)|\mathbf{u}(1), ..., \mathbf{u}(t-1)]$ and the total law of expectation: $\mathbb{E}\left[\sum_{t=1}^T \langle \nabla \tilde{L}_t(\mathbf{p}_{\delta}(t)) | \mathbf{p}_{\delta}(t) - \mathbf{q}_{\delta} \rangle\right] \leq -\mathbb{E}\left[\sum_{t=1}^T \langle \tilde{\mathbf{v}}(t) | \mathbf{p}_{\delta}(t) - \mathbf{q}_{\delta} \rangle\right]$

• Using the convex conjugate of $h(\mathbf{p}_{\delta})$, the links with the update and $||\tilde{v}(t)||_{\infty}^2 \leq S^2 (B/\delta + L)^2$

$$\mathbb{E}[\operatorname{Reg}(T)] \leq \frac{P_{\max}\log(1+S)}{2\gamma} + \gamma TS^2 \left(\frac{B}{\delta} + L\right)^2 + LT\delta \left(3 + P_{\max}\left(S + 2\sqrt{S}\right)\right).$$

where $B = \max_{t,\mathbf{p}} L_t(\mathbf{p})$

E. V. Belmega

• Based on gradient feedback, we proposed

matrix gradient descent and exponential learning for rate/energy-efficiency maximization

- Exponential mapping preserves the positivity of the input covariance matrices
- Our online policies perform close to the **ideal benchmark:** dynamic optimal one

COST-HATA wireless channels, mobile users [3 - 130] km/h, K=15 users, 4 \times 8 MIMO, S=8 bands



P. Mertikopoulos, and E.V. Belmega, "Learning to be green: robust energy efficiency maximization in dynamic MIMO-OFDM systems", IEEE JSAC, vol. 34, no. 4, pp. 743 – 757, Apr. 2016.

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024

30/30

- Multi-antenna devices: M_k transmit antennas at user $k \in \{1, \ldots, K\}$, N receive antennas
- Achieved Shannon sum rate

$$R(\mathbf{Q}) = \log \det \left(\mathbf{I} + \sum_{k=1}^{K} \mathbf{H}_k \mathbf{Q}_k \mathbf{H}_k^{\dagger} \right)$$

 $\mathbf{H}_k \in \mathbb{C}^{N \times M_k}$ - channel matrix between transmitter k and receiver, $\mathbf{Q} = (\mathbf{Q}_1, \ldots, \mathbf{Q}_K), M = \max_k M_k$

- Distributed optimization: each transmitter k tunes its input covariance matrix
 - $\mathbf{Q}_k \in \mathbb{C}^{M_k imes M_k}$ s.t.

$$\mathbf{Q}_k \succeq 0, \quad \operatorname{Tr}(\mathbf{Q}_k) = P_{\max,k}$$

P. Mertikopoulos, and E.V. Belmega, "Learning to be green: robust energy efficiency maximization in dynamic MIMO-OFDM systems", IEEE JSAC, vol. 34, no. 4, pp. 743 – 757, Apr. 2016.

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024

30/30

- Assumptions: static channels, S = 1
- Rate maximization algorithms based on gradient feedback:

MXL - *matrix exponential learning* [*MertikopoulosB'16*] *IWF* - *iterative waterfilling* [*Yu'04*]

$$\nabla_k R(\mathbf{Q}) = \mathbf{H}_k^{\dagger} \left[\mathbf{I} + \sum_{\ell=1}^K \mathbf{H}_\ell \mathbf{Q}_\ell \mathbf{H}_\ell^{\dagger} \right]^{-1} \mathbf{H}_k$$

• Massive MIMO systems: N >> M, K >> 1

Issue: too much feedback $\mathcal{O}(\min\{KM^2, N^2\})$ at each iteration!

MXL algorithm

[MertikopoulosB'16]

- Transmit at $\mathbf{Q}(t)$
- Receive matrix gradient feedback:
 V(t) = ∇_k R(Q(t))
- Update:

$$\mathbf{Y}(t+1) = \mathbf{Y}(t) + \gamma \mathbf{V}(t)$$

$$\mathbf{Q}(t+1) = \Lambda(\mathbf{Y}_t), \ \Lambda_k = P_{\max,k} \ \frac{\exp(\mathbf{Y}_k)}{\operatorname{tr}(\exp(\mathbf{Y}_k))}$$

- MXL with utility-based gradient estimator: feedback $R(t) = R(\mathbf{Q}(t))$, slow convergence $\mathcal{O}(T^{-1/4})$, poor scaling with problem dimensionality (MXL0)
- Idea: exploit the current sum rate feedback R(t) jointly with the previous one $R(t-1) = R(\mathbf{Q}(t-1))$ (MXL0⁺)

 \longrightarrow improved two-point gradient estimator with no additional feedback



Arbitrary static channels, K = 20 users, 3 transmit and 16 receive antennas

Remark: for static convex and Lipschitz objectives, $MXL0^+$ recovers the convergence rate $O(1/\sqrt{T})$ as MXL with perfect gradient feedback !

O. Bilenne, P. Mertikopoulos, and E.V. Belmega, "Fast Gradient-Free Optimization in Distributed Multi-User MIMO Systems", IEEE Trans. on Signal Processing, Oct. 2020.

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th. 2024

30/30

- First-order derivative: $\ell'(x) = \lim_{\delta \to 0} \frac{\ell(x+\delta) \ell(x-\delta)}{2\delta}$
- A two query estimator for small δ : $\hat{v}_2 = \frac{\ell(x+\delta) \ell(x-\delta)}{2\delta}$
- One query random estimator: $u \sim \text{Bernoulli}(1/2)$

$$\hat{v} = \frac{\ell(x+\delta u)u}{\delta}, \qquad \mathbb{E}[\hat{v}] = \hat{v}_2$$

Bias: $|\mathbb{E}[\hat{v}] - v| = \mathcal{O}(\delta)$, variance: $|\hat{v}| = \mathcal{O}(1/\delta)$

• One query with callback ρ :

$$\hat{v}_{\rho} = rac{\left(\ell(x+\delta u)-\rho\right)u}{\delta}, \qquad \rho = \ell(x)$$

Bias: $|\mathbb{E}[\hat{v}_{\rho}] - v| = \mathcal{O}(\delta)$, variance: $|\hat{v}_{\rho}| \leq L$ L - Lipschitz constant of ℓ

MXL0+

• Ensure the query point falls in the feasible set

$$\hat{\mathbf{Q}}_k = \mathbf{Q}_k + rac{\delta}{r_k} (\mathbf{C}_k - \mathbf{Q}_k) + \delta \mathbf{Z}_k \in \mathcal{Q}_k$$

K users, M_k transmit antennas, N receive antennas, $M = \max_k M_k$ $\mathbf{C}_k = \mathbf{I}/M_k, r_k = 1/\sqrt{M_k(M_k-1)}, \delta < r_k$ \mathbf{Z}_k uniformly drawn on the unit sphere, $d_k = M_k^2 - 1$ dimension of \mathcal{Q}_k

- One-point estimator: $\hat{\mathbf{V}}_k(\mathbf{Q}) = \frac{d_k}{\delta} R(\hat{\mathbf{Q}}) \mathbf{Z}_k$
- One-point estimator with callback: $\tilde{\mathbf{V}}_{k,t}(\mathbf{Q}) = \frac{d_k}{\delta} \left[R(\hat{\mathbf{Q}}_t) R(\hat{\mathbf{Q}}_{t-1}) \right] \mathbf{Z}_{k,t}$
- Remark: our algorithm does not require two queries at each iteration t!

• Convergence:
$$\mathcal{O}(\sqrt{K^4 M^6 \log M/T})$$
 with $\gamma \propto 1/\sqrt{T}, \delta \propto 1/\sqrt{T}$



Qu: Can we design an *adaptive NOMA* scheme via MABs relying on **little feedback** and **no CSIT/CDIT**?

Main challenge: Unknown channels that vary in time

- Simple case: K = 2 (or user pairing)
- Decoding order is a *discrete* control variable $i \in \{1, 2\}$
 - ► User i ∈ {1, 2} performs SIC
 - User $j \neq i$ suffers interference

E. V. Belmega

• Fix the transmit rates $R_{\min,k}$ and measure the outage

$$\mathbb{P}_{\text{out}}(i, \mathbf{p}) = \mathbb{P}[R_i(t) \le R_{\min, i} \cup \min\{R_{j \to j}(t), R_{j \to i}(t)\} \le R_{\min, j}]$$

· Feedback: 1-ACK/NACK bit from each user

[ElHassaniSB'21a, ChafaaBD'20]

30/30

$$R_{i}(t) = \log\left(1 + \frac{|h_{i}(t)|^{2}p_{i}}{\sigma_{i}^{2}}\right), \ R_{j \to j}(t) = \log\left(1 + \frac{|h_{j}(t)|^{2}p_{j}}{|h_{j}(t)|^{2}p_{i} + \sigma_{j}^{2}}\right), \ R_{j \to i}(t) = \log\left(1 + \frac{|h_{i}^{(t)}|^{2}p_{j}}{|h_{i}^{(t)}|^{2}p_{i} + \sigma_{i}^{2}}\right)$$

H. El Hassani, A. Savard, and E. V. Belmega, "Adaptive NOMA in time-varying wireless networks with no CSIT/CDIT relying on a 1-bit feedback", accepted paper IEEE Wireless Commun. Lett., Nov 2020

Télécom Paris - Dép. ICE

· Ideal energy-efficient target

[Zhang'20]

maximize
$$\operatorname{GEE}(i, \mathbf{p}) \triangleq \frac{(R_{\min,i} + R_{\min,j}) (1 - \mathbb{P}_{\operatorname{out}}(i, \mathbf{p}))}{p_i + p_j + P_c}$$

subject to $i \in \{1, 2\}, \quad \mathbf{p} = (p_i, p_j) \in \mathbb{R}^2_+, \ p_i + p_j \leq P_{\max}$

- Quantization of the power allocation
 - User fairness: less power is allocated to the SIC decoding user i

$$\mathbf{p}_{\beta} = (0.25 \ \beta \ P_{\max}, 0.75 \ \beta \ P_{\max})$$
, for discrete $\beta \in \mathcal{B}$

- Uniformly quantized $\mathcal{B} \subset [0, 1]$
- An arm $\mathbf{a} \triangleq (i, \mathbf{p}_{\beta}) \in \mathcal{A}$ decoding and power allocation vector

$$\mathcal{A} = \{1, 2\} \times \{\mathbf{p}_{\beta} = (0.25 \ \beta \ P_{\max}, 0.75 \ \beta \ P_{\max}) \mid \beta \in \mathcal{B}\}$$

Rk: The quantization and the 1/4 - 3/4 user split will imply an optimality loss.

E. V. Belmega

Télécom Paris - Dép. ICE

Adaptive energy-efficient 1-bit feedback NOMA

Initialize: t = 1, $\mathbf{a}(1) = (1, 0.5)$

For t = 1 to T

- Play arm $\mathbf{a}(t) = (i(t), \beta(t))$
 - Inform users of their decoding: i(t)
 - Transmit with $\mathbf{p}_{\beta}(t) = (0.25 \ \beta(t) \ P_{\max}, \ 0.75 \ \beta(t) \ P_{\max})$
- Receive 1-bit ACK/NACK feedback from each user
- Compute reward

$$U_t(\mathbf{a}(t)) = \begin{cases} \frac{R_{\min,1} + R_{\min,2}}{\beta(t)P_{\max} + P_c}, & \text{if } R_i(t, \mathbf{a}(t)) \ge R_{\min,i} \text{ and} \\ & \min\{R_{j \to j}(t, \mathbf{a}(t)), R_{j \to i}(t, \mathbf{a}(t))\} \ge R_{\min,i} \\ 0, & \text{otherwise} \end{cases}$$

• Update policy
$$\mathbf{a}(t+1) \leftarrow \mathbf{a}(t)$$
 via UCB or EXP3

• $t \leftarrow t+1$

- Both UCB and EXP3 have the no regret property
 - \rightarrow adaptive NOMA reaches the best expected reward

$$\max_{\mathbf{a}\in\mathcal{A}} \mathbb{E}[U_t(\mathbf{a})] \equiv \max_{i\in\{1,2\},\beta\in\mathcal{B}} GEE(i,\mathbf{p}_\beta)$$

- Regret decay rate: $Reg_T^{\text{UCB}} = \mathcal{O}(\log T/T)$ and $Reg_T^{\text{EXP3}} = \mathcal{O}(1/\sqrt{T})$
- Complexity of UCB and EXP3: every iteration scales as $\mathcal{O}(|\mathcal{A}|)$

Rk: There are two fundamental tradeoffs: performance vs. feedback information and performance vs. complexity.

H. El Hassani, A. Savard, and E. V. Belmega, "Adaptive NOMA in time-varying wireless networks with no CSIT/CDIT relying on a 1-bit feedback", submitted paper IEEE SSP, Feb. 2021.

30/30





Our 1-bit feedback adaptive NOMA scheme outperforms quantized OMA with perfect CDIT.

 $\sigma_{k}^{2} = 0.1, R_{\min,1} = 1$ bpcu, $R_{\min,2} = 10$ bpcu, $P_{c} = 1$ W, $P_{\max} = 100$ W, $T = 5000, 10^{3}$ channel (Rayleigh) realizations

E. V. Belmega

Télécom Paris - Dép. ICE



Energy efficiency as a function of $\Gamma_{\min,2}=2^{R\min,2}-1$

Number of iterations to reach a 10% regret level

Our adaptive 1-bit feedback NOMA scheme is sub-optimal compared to the ideal target (with CDIT).

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024



10

Outage minimization problem

Malicious jammer

• Has access to perfect CSI

E. V. Belmega

- Anticipates the arm that will be played via deterministic UCB
- · Interferes such that the system is put systematically in outage



In non-stationary adversarial settings EXP3 reaches no regret while UCB is brought to a fault.

H. El Hassani, A. Savard, and E. V. Belmega, "Adaptive NOMA in time-varying wireless networks with no CSIT/CDIT relying on a 1-bit feedback", accepted paper IEEE Wireless Commun. Lett., Nov. 2020.

Télécom Paris - Dép. ICE



Malicious jammer

• Has access to perfect CSI

E. V. Belmega

- Anticipates the arm that will be played via deterministic UCB
- Interferes such that the system is put systematically in outage



In non-stationary adversarial settings EXP3 reaches no regret while UCB is brought to a fault.

H. El Hassani, A. Savard, and E. V. Belmega, "Adaptive NOMA in time-varying wireless networks with no CSIT/CDIT relying on a 1-bit feedback", accepted paper IEEE Wireless Commun. Lett., Nov 2020

Télécom Paris - Dép. ICE

Jan. 6th, 2024

30/30

- They tradeoff between data exploration and exploitation
 - \rightarrow poorly chosen lead to **positive regret**
- From a theoretical perspective (to reach no regret asimptotically)

$$\eta^* > 2$$
 and $\gamma^* = \sqrt{rac{|\mathcal{A}| \ln |\mathcal{A}|}{(e-1)T}}$

- But these values minimize the upper-bound of the regret
 - \rightarrow not optimal for the actual regret
- They can be further shaped via empirical experiments
Online mirror descent for convex optimization

• **Regret guarantees** $\mathcal{O}(\sqrt{T})$ rely on Lipschitz continuity

$$|\ell(\mathbf{x}) - \ell(\mathbf{y})| \le G ||\mathbf{x} - \mathbf{y}||$$
 or $||\operatorname{grad} \ell(\mathbf{x})|| \le G$

- We extend them to problems with singular gradients of the objective at the boundary of the feasible set
- Idea: generalize the Lipschitz continuity based on local Riemannian norms, $\|\mathbf{z}\|_x = \mathbf{z}^T \mathbf{G}(\mathbf{x}) \mathbf{z}$ with $\mathbf{G}(\mathbf{x}) \succeq 0$

Collaborator: P. Mertikopoulos (CNRS) Student: K. Antonakopoulos (PhD) Supported by: Inria, ANR ORACLESS Publications: 3 confs.

E. V. Belmega

e.g.
$$\ell(\mathbf{x}) = -\log(\mathbf{a}^T \mathbf{x}), \|\text{grad }\ell\|^2 = \frac{\sum_i a_i^2}{(\mathbf{a}^T \mathbf{x})^2}$$

$$\mathbf{G}(\mathbf{x}) = \mathbf{I}/\left(\sum_i x_i\right)^2 \Rightarrow \|\text{grad }\ell\|_x^2 \le \frac{\sum_i a_i^2}{(\min_j a_j)^2}$$

Jan. 6th. 2024

30/30

K. Antonakopoulos, E.V. Belmega, and P. Mertikopoulos, "Online and stochastic optimization beyond Lipschitz continuity: A Riemannian approach", ICLR, spotlight talk, 2020.

Télécom Paris - Dép. ICE

30/30

- · Applications: tomography, astronomical data
- Ill-conditioned reconstruction problems: $\hat{\mathbf{u}} = \mathbf{H}\mathbf{x} + \hat{\mathbf{z}} \quad \dim \hat{\mathbf{u}} \ll \dim \mathbf{x}$ H - (Toeplitz, circulant, convoluting) optical system, $\hat{\mathbf{z}}$ - Poisson noise $\hat{u}_i \sim \text{Poisson}(\mathbf{H}\mathbf{x})_i$
- **Objective:** minimize Kullback-Leibler divergence $\ell(\mathbf{x}, \hat{\mathbf{u}}) = \sum_{i} \left[\hat{u}_i \log(\hat{u}_i / (\mathbf{H}\mathbf{x})_i) + (\mathbf{H}\mathbf{x})_i - \hat{u}_i \right]$
- Example: test image contaminated with Poisson noise









RMD reconstruction



 $\begin{array}{l} \dim \, \hat{u} = \dim \, \hat{x}/2, \dim \, \hat{x} \simeq 10^5 \\ \mathrm{LR-Luy-Richardson, mirror descent with entropic regularizer, [Bertero-2009] \\ \mathrm{CMD-mirror descent with local Riemann norm \\ \mathrm{RMD-mirror descent with local Riemann norm } \end{array}$

K. Antonakopoulos, E.V. Belmega, and P. Mertikopoulos, "Online and stochastic optimization beyond Lipschitz continuity: A Riemannian approach", ICLR, spotlight talk, 2020.



Télécom Paris - Dép. ICE

Jan. 6th. 2024

Online metric learning for multimedia indexing



$$d_X(\mathbf{p}, \mathbf{q}) = \|\mathbf{X}^{1/2}\mathbf{p} - \mathbf{X}^{1/2}\mathbf{q}\|^2$$

p, q - input data, e.g., images

• **Objective:** learn the best discriminative linear data transformation matrix **X**

 $\mathbf{X} \succeq 0, \quad \operatorname{Tr}(\mathbf{X}) \le c$

- · Training examples are obtained online, on-the-fly
- · Our matrix exponential learning is competitive in a toy setup
- Other examples: online matrix completion, universal linear filtering, online dictionary learning



MDML - mirror descent with Frobenius norm regularizer [Kunapuli-2012]

P. Mertikopoulos, E. V. Belmega, R. Negrel, and L. Sanguinetti, "Distributed stochastic optimization via matrix exponential learning", IEEE Trans. on Signal Processing, 2017.

E. V. Belmega

Télécom Paris - Dép. ICE

Jan. 6th, 2024

30/30

- R. M. Lee, M. J. Assante, and T. Conway, "Analysis of the cyber attack on the Ukrainian power grid. Defense use case", Electricity Information Sharing and Analysis Center, Mar. 2016.
- M. Ozay, I. Esnaola, F. T. Yarman Vural, S. R. Kulkarni, and H. V. Poor, "Machine learning methods for attack detection in the smart grid", *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 27, no. 8, pp. 1773–1786, 2016.
- A. Altman, A. Bercovici-Boden, and M. Tennenholtz, "Learning in one-shot strategic form games", In European Conf. on Machine Learning, Sep. 2006.
- Y. Wu, A. Khisti, C. Xiao, G. Caire, K.K. Wong, and X. Gao, "A survey of physical layer security techniques for 5G wireless networks and challenges ahead", IEEE J. Sel. Areas Commun., vol. 36, no 4, pp. 679–695, 2018.
- H. Xing, K.K. Wong, A. Nallanathan, and R. Zhang, "Wireless powered cooperative jamming for secrecy multi-AF relaying networks", IEEE Trans. Wireless Commun., vol. 15, no. 12, pp. 7971–7984, 2016.
- H. Zhang, J. Zhang, and K. Long, "Energy efficiency optimization for NOMA UAV network with imperfect CSI", *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 12, pp. 2798–2809, 2020.
- G. Kunapuli and J. Shavlik, "Mirror descent for metric learning: A unified approach", Joint European Conference on Machine Learning and Knowledge Discovery in Databases. Springer, 2012, pp. 859–874.
- W. Yu, W. Rhee, S. Boyd, and J. M. Cioffi, "Iterative water-filling for Gaussian vector multiple-access channels", IEEE Trans. Inf. Theory, vol. 50, no. 1, pp. 145–152, 2004.

- M. Bertero, P. Boccacci, G. Desidera, and G. Vicidomini, "Image deblurring with Poisson data: from cells to galaxies", Inverse Problems, 2009.
- N. He, Z. Harchaoui, Y. Wang, and L. Song, "Fast and simple optimization for poisson likelihood models", arXiv preprint arXiv:1608.01264, 2016.

E. V. Belmega

Télécom Paris - Dép. ICE