

Power-efficient deep learning.

①

① Toy example / finite case.

Obj: introduire de nouveaux critères environnementaux dans l'apprentissage des réseaux de neurones.
vers 5 ans.

- non-parametric statistics / tradeoff complexity / fitting the data
- Bayesian priors.
- Optimal transport and cost function.

1. online learning & exponential weighted averages.

toy example

	t=0	0 1 0 0 1 0 0 1 1 1 0 1	?	1	1	1	1	2	1	$\hat{y}_t = \bar{y}_t(y_1, \dots, y_{t-1})$
N experts	{	0 1 0 0 X X 0 0 1 1 X X 0 0	1	1	$e_{1k}, k=1, N,$ e_{2k}, \dots

solution:

keeper kill.

hypothesis: 1 expert knows the story.

$$\text{err}(T) \leq N-1 \Rightarrow \text{chaque fois au moins un expert au hasard mais après } N-1 \text{ errors il ne reste que "l'expert".}$$

(keep or kill + majority vote) optimal.

$$\text{err}(T) \leq (\log N)$$

w_e = nombre de succès après T errors.

$$w_p \leq \frac{1}{2} w_{p-1} \cdot (N/V).$$

$$w_p \leq 2^{-e} w_0$$

$$w_0 = N$$

$$2^e \leq w_0 = N.$$

$$e \leq \log N.$$

General case: uncountable \mathcal{G} / Deep learning.
 Toy example \rightarrow general case (finite set of experts) \rightarrow uncountable set.

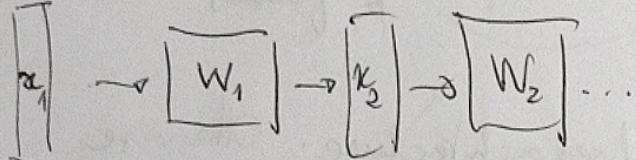
z_1, \dots, z_t ?

$$\mathcal{G} = \{g_\theta, \theta \in \Theta \subseteq \mathbb{R}^P\}$$

\hookrightarrow deep nets, linear regres., etc

$$z_t = (x_t, y_t)$$

$$g_\theta(z_t) = \langle z_t, \theta \rangle$$



$$g_\theta(z_t) = \langle z_t, f_\theta \rangle,$$

$$= \sum_{k=1}^P \theta_k f_k(z_t)$$

$$\Theta = (W_1, \dots, W_L) \subseteq \mathbb{R}^P, \quad \text{span}\{f_1, \dots, f_P\} = \mathcal{G}$$

$P \approx 60 \text{ millions.}$

PAC-Bayesian bound.

$$\sum_{t=1}^T \mathbb{E}_{g \sim \hat{g}_t^*} \ell(g_t^*, z_t) \leq \inf_{\pi \in \Delta(\mathcal{G})} \left\{ \frac{\mathbb{E}_{g \sim \hat{g}_t^*} \sum_{r=1}^T \ell(g_r, z_r)}{\pi(g)} + \frac{k(s, \pi)}{\lambda} \right\},$$

$$|\hat{g}_t - g_{\theta_t}(z_t)|$$

$$(s, \pi) = 0 \Rightarrow g(s) = 0.$$

$$\text{where } k(s, \pi) = \begin{cases} \left| \frac{d\pi}{ds} \right| & s \leq \pi \\ +\infty & \text{if } s > \pi \end{cases} \quad \hat{g}_{t+1} \sim \exp \left(\sum_{u=1}^t \ell(g_u, z_u) \right) \pi$$

Kullback-Leibler divergence -

Proof (sketch). Audibert 2009. Consequence -

power tail.

$$\pi \text{ sparsity prior.} \rightarrow \pi(\theta) = \prod_{j=1}^P \text{student}_2(\theta_j) \sim \frac{1}{|\Theta|} \cdot \frac{\text{const}}{(1 + \frac{w_j}{\sqrt{\lambda}})^4}$$

$$k(\beta_{w_0}, \pi) = \int \beta_{w_0}(w) \ln \frac{\pi(w)}{\pi(w_0)} dw = \|\omega_0\|_2 \ln \left(1 + \frac{\|\omega\|_2}{\|\omega_0\|_2} \right)$$

translate de π

de directe w_0 .

$$\pi_{w_0}(w) = \pi(w - w_0).$$

N naburde coefficients non nulles

③ Generalization - generic regret bounds

(3)

$$g_{t+1} = \arg \min_g \left\{ \sum_{u=1}^t \tilde{l}(g_u, x_u) + \frac{K(g, \pi)}{\lambda} \right\} = \arg \min_g \left\{ \tilde{l}(g, x_t) + \frac{K(g, g_t)}{\lambda} \right\}$$

Convex duality formula

$$p(dg) = \exp(-h(g)) \pi(dg)$$

$$= \arg \min_p \left\{ \mathbb{E}_{g \sim p} h(g) + \frac{K(g, \pi)}{\lambda} \right\}$$

generalization

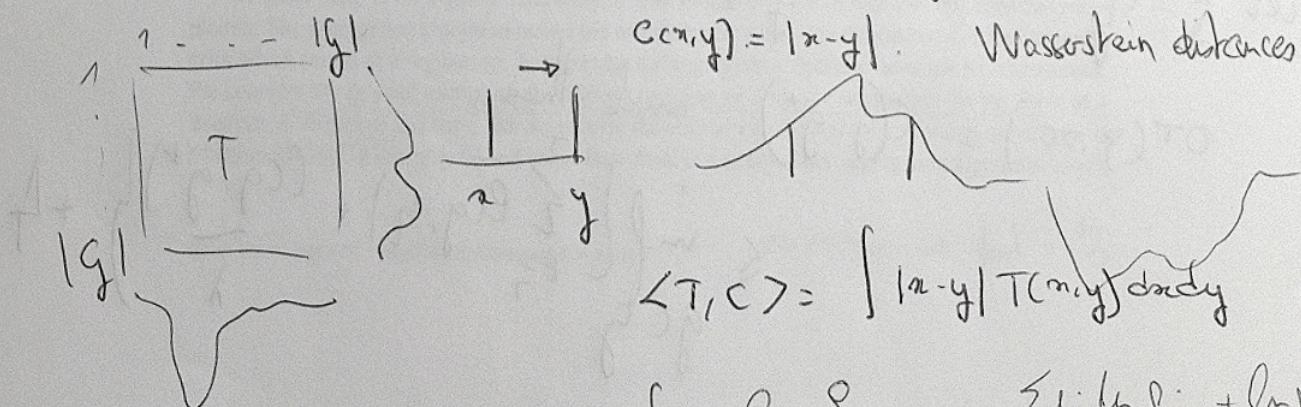
$$\min_p \left\{ \mathbb{E}_{g \sim p} h(g) + \frac{D(g, \pi)}{\lambda} \right\}, \text{ where } D(g, \pi) = \begin{array}{l} \text{Bregman} \\ \text{divergence} \\ \text{optimal transport} \end{array}$$

$$D_\phi(g, \pi) = \phi(g) - \phi(\pi) - \nabla \phi(\pi)(g - \pi).$$

$$D(g, \pi) = \min_{T \in \Delta(g, \pi)} \langle T, C \rangle \quad \text{where } C: \mathcal{G} \times \mathcal{G} \rightarrow \mathbb{R} \quad \text{cost function}$$

to move a unit mass from g to π .

$$\Delta(g, \pi) = \{T: \mathcal{G} \times \mathcal{G} \rightarrow \mathbb{R}_+, \sum_{g'} T(g, g') = \delta(g), \sum_g T(g, g') = \pi(g')\}.$$



$$H(g, \pi) - H(g) = K(g, \pi) = \int g \ln \frac{g}{1/N} = \sum_i f_i \ln f_i + \ln N$$

$$= \ln N - H(g).$$