

# An Optimal Transport View on Generalization

Nemo Fournier

January 13, 2020

# Outline

Framework

Main results

An application

Deep Neural Networks

instance space  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$   
hypothesis space  $\mathcal{W}$   
loss function  $\ell : \mathcal{Z} \times \mathcal{W} \rightarrow \mathbb{R}^+$

learning algorithm  $\mathcal{A} : \mathcal{Z}^n \rightarrow \mathcal{W}$   
underlying distribution  $D$   
training sample  $S_n \sim D^{\otimes n}$

instance space  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$   
 hypothesis space  $\mathcal{W}$   
 loss function  $\ell : \mathcal{Z} \times \mathcal{W} \rightarrow \mathbb{R}^+$

learning algorithm  $\mathcal{A} : \mathcal{Z}^n \rightarrow \mathcal{W}$   
 underlying distribution  $D$   
 training sample  $S_n \sim D^{\otimes n}$

$$\text{risk } R(w) = \mathbb{E}_{z \sim D}[\ell(z, w)]$$

$$\text{empirical risk } R_{S_n}(w) = \mathbb{E}_{z \sim S_n}[\ell(z, h)] = \frac{1}{n} \sum_{i=1}^n \ell(z_i, w)$$

$$\text{generalization error } G(D, P_{W|S_n}) = \mathbb{E}[R(W) - R_{S_n}(W)]$$

$\mu$  and  $\nu$  two measures on  $\mathcal{W}$

**coupling**  $T$  measure on  $\mathcal{W} \times \mathcal{W}$  such that  $\begin{cases} T(X, \mathcal{W}) = \mu(X) \\ T(\mathcal{W}, X) = \nu(X) \end{cases}$

**wasserstein distance**

$$\mathbb{W}_1(\mu, \nu) = \inf_{T \in \Gamma(\mu, \nu)} \mathbb{E}_{(W, W') \sim T} [d_{\mathcal{W}}(W, W')]$$

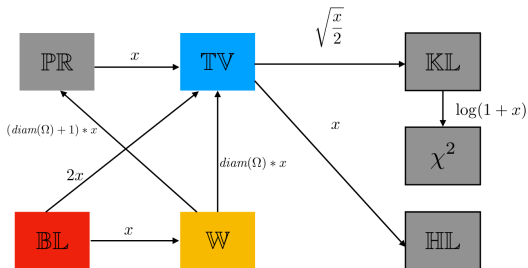
**algorithmic transport cost of algorithm  $\mathcal{A}$**  ( $P_{W|S_n}$ )

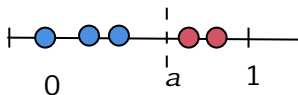
$$\text{Opt}(D, P_{W|S_n}) = \mathbb{E}_{z \sim D} [\mathbb{W}_1(P_W, P_{W|z})]$$

$$A.G.T. Opt(D, P_{W|S_n}) = \mathbb{E}_{z \sim D} [\mathbb{W}_1(P_W, P_{W|z})]$$

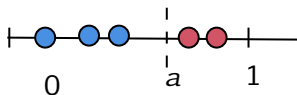
main theorem

$$G(D, P_{W|S_n}) = \mathbb{E}[R(W) - R_{S_n}(W)] \leq K \times Opt(D, P_{W|S_n})$$





$$A: \begin{cases} \mathcal{Z}^n & \rightarrow \\ S_n = \{(x_1, y_1), \dots, (x_n, y_n)\} & \mapsto \end{cases} \begin{cases} \mathcal{W} \\ \max_{1 \leq i \leq n} x_i \text{ if } \{i \mid y_i = 0\} \neq \emptyset \\ \text{s.t. } y_i = 0 \\ 0 \text{ otherwise} \end{cases}$$



$$A: \begin{cases} \mathcal{Z}^n & \rightarrow \mathcal{W} \\ S_n = \{(x_1, y_1), \dots, (x_n, y_n)\} & \mapsto \begin{cases} \max_{1 \leq i \leq n} x_i & \text{if } \{i \mid y_i = 0\} \neq \emptyset \\ \text{s.t. } y_i = 0 \\ 0 & \text{otherwise} \end{cases} \end{cases}$$

$$P_W(w) = (1 - a)^{n-k} + n(w + 1 - a)^n$$

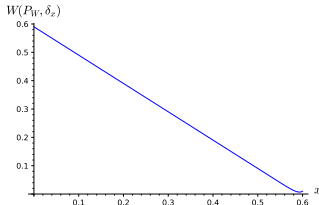
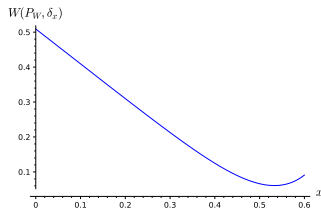
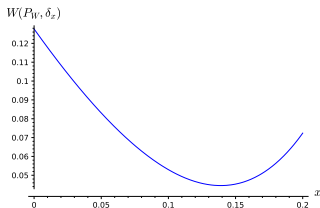
$$P_{W|z} = \delta_x \text{ if } x \leq a$$

$$P_{W|z} = \delta_0 \text{ otherwise}$$

$$\mathbb{W}_1(\mu, \delta_t) = \mathbb{E}_{X \sim \mu}[d(X, t)] \implies \mathbb{W}_1(P_W, \delta_x) = \int_0^a |x - w| P_W(w) dw$$



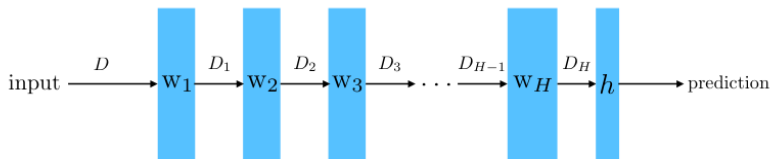
$$\begin{aligned} \mathbb{W}_1(P_W, \delta_x) = & \frac{1}{2(an+a)} \left( a^2((-a+1)^n + 2)n + a^2(3(-a+1)^n + 2) + 2((-a+1)^n n + (-a+1)^n)x^2 \right. \\ & - 4(a^2 - ax - a)(-a+x+1)^n - 2a((-a+1)^n + 1) - 2(a(2(-a+1)^n + 1)n \\ & \left. + a(2(-a+1)^n + 1))x \right) \end{aligned}$$



$$\text{Opt}(D, P_{W_n}) = \int_0^a \mathbb{W}_1(P_W, \delta_x) dx + (1-a) \mathbb{W}_1(P_W, \delta_0)$$

$$\begin{aligned} \text{Opt}(D, P_{W_n}) = & \frac{1}{6(n^2 + 3n + 2)} (2a^2(2(-a+1)^n - 3) - (a^2(4(-a+1)^n + 3) - 3a((-a+1)^n + 2))n^2 \\ & - 3(3a^2 + 3a((-a+1)^n - 2) - 2(-a+1)^n + 2)n) - 6a((-a+1)^n - 2) \end{aligned}$$

$$G(D, P_{W|S_n}) \leq 1 \times \text{Opt}(D, P_{W_n})$$



$$\mathbb{E}\left[R(W) - R_{S_n}(W)\right] \leq \exp\left(-\frac{H}{2} \log \frac{1}{\eta}\right) \sqrt{\frac{K^2 R^2 I(S_n; W)}{2n}}$$

Powerful theoretical tool (average case, link with information theory)

Quite quickly too convoluted to provide concrete bounds