

# Domain generalization (2023)

© Théo Jalabert



$D_{tr}$  training set of  $S$  domains

$$D_{tr} = \{D_1, \dots, D_S\}$$

$$\forall s=1, \dots, S \quad D_s = \{(x_i^s, y_i^s)\}_{i=1}^{n_s}$$

↑ different distributions  
in the domains!  
(but similar)

$\ell(x, y; \theta)$  - loss of model  $\theta$ .

Classical approach: empirical risk  $\min \mathcal{L}_{ERM} = \frac{1}{|D_{tr}|} \sum_{D \in D_{tr}} \mathbb{E}_{(x,y) \sim D} [\ell(x, y; \theta)]$   
it fails to generalize on new domains

## (5) Pitfalls of ERM

- Binary classification
- $x \in \{0, 1\}^4$ ,  $y \in \{0, 1\}$

	$x_1$	$x_2$	$x_3$
	50%	40%	10%
$\mathbb{Q}_1$	(0, 0, 0, 0)	(1, 1, 0, 0)	(1, 0, 0, 0)
$\mathbb{Q}_2$	(0, 0, 0, 0)	(1, 0, 1, 0)	(1, 0, 0, 0)
$\mathbb{Q}_3$	(0, 0, 0, 0)	(1, 0, 0, 1)	(1, 0, 0, 0)
	$P(Y=1 x_1)=0$	$P(Y=1 x_2)=1$	$P(Y=1 x_3)=0.3$

$x = (f_1, f_2, f_3, f_4)$ . Accuracy of the optimal classifier considering only  $f_1$ ?

Optimal classifier  $y(f_1) = \begin{cases} 0 & \text{if } f_1=0 \text{ - obvious} \\ 1 & \text{if } f_1=1 \text{ - higher accuracy than for } y(1)=0 \end{cases}$   
 $\rightarrow y^*(f_1) = f_1$

$$\text{accuracy} = \mathbb{P}(y^*(f_1) = Y) = \mathbb{P}(Y=0|x_1)\mathbb{P}(x_1) + \mathbb{P}(Y=1|x_2)\mathbb{P}(x_2) +$$

$$+ \mathbb{P}(Y=1|x_3) \mathbb{P}(x_3) = 0.93^{0.3}$$

No optimality of ERM principle: if we could take  $x_i$  into account, we would predict  $\begin{cases} \hat{y}^*(f_1, x_3) = 1-f_1 \\ \hat{y}^*(f_1, x_i) = f_1 \end{cases} \quad i=1,2$  and would get better accuracy  $0.9 + 0.7 \cdot 0.1 = 0.97$ .

## (2) Meta-learning - alternative to ERM

$\mathcal{D}_{\text{Tr}} = \mathcal{D}_A \cup \mathcal{D}_B$  Train on  $\mathcal{D}_A$  so it generalizes well on  $\mathcal{D}_B$

$$\mathcal{L}_{\mathcal{D}_A}(\theta) = \mathbb{E}_{\mathcal{D}_A} [\ell(X, Y; \theta)] \quad \mathcal{L}_{\mathcal{D}_B}(\theta) = \mathbb{E}_{\mathcal{D}_B} [\ell(X, Y; \theta)]$$

Meta-loss  $\mathcal{L}_{\text{met-L}}(\theta) = \mathcal{L}_{\mathcal{D}_A}(\theta) + \underbrace{\mathcal{L}_{\mathcal{D}_B}(\theta - \lambda \nabla \mathcal{L}_{\mathcal{D}_A}(\theta))}_{\theta'}$

Gradient step  $\theta'$  should minimize  $\mathcal{L}_{\mathcal{D}_B}$

$$(2.1) \text{ Order 1 Taylor for } \mathcal{L}_{\mathcal{D}_B}(\theta') = \mathcal{L}_{\mathcal{D}_B}(\theta) + \nabla \mathcal{L}_{\mathcal{D}_B}(\theta) \cdot (-\lambda \nabla \mathcal{L}_{\mathcal{D}_A}(\theta)) + \bar{o}(|\theta - \theta'|)$$

$$= \mathcal{L}_{\mathcal{D}_B}(\theta) - \lambda \langle \nabla \mathcal{L}_{\mathcal{D}_B}(\theta), \nabla \mathcal{L}_{\mathcal{D}_A}(\theta) \rangle + \bar{o}(|\theta - \theta'|)$$

$$(2.2) \mathcal{L}_{\text{met-L}}(\theta) \approx \mathcal{L}_{\mathcal{D}_A}(\theta) + \mathcal{L}_{\mathcal{D}_B}(\theta) - \lambda \langle \nabla \mathcal{L}_{\mathcal{D}_A}(\theta), \nabla \mathcal{L}_{\mathcal{D}_B}(\theta) \rangle$$

- (2.3) Interpretation:
- Minimization of loss both on  $\mathcal{D}_A$  and  $\mathcal{D}_B$
  - At the same time maximization of the  $\langle \nabla \mathcal{L}_{\mathcal{D}_A}, \nabla \mathcal{L}_{\mathcal{D}_B} \rangle$  = colinearity of these vectors  $\rightarrow$  simultaneous decrease of loss on  $\mathcal{D}_A$  and  $\mathcal{D}_B \rightarrow$  ability to generalize on domains not included in the train set.

### 3. Gradient alignment

$$G_i := \mathbb{E}_{D_i} \left[ \frac{\partial l(x, y; \theta)}{\partial \theta} \right]$$

$$\text{New loss} = \mathcal{L}_{ERM}(D_{Tr}; \theta) - \frac{2}{N_{D_{Tr}}(N_{D_{Tr}}-1)} \sum_{\substack{i, j \in D_{Tr} \\ \text{nb of sets } i \neq j}} \langle G_i, G_j \rangle$$

(3.1) algo 2 is computationally prohibitive because of the double sum  $\sum$  and the necessity to calculate  $\frac{\partial G}{\partial \theta}$  (numerically or in terms of Hessians)

(3.2) For two domains: explicit expression for updates

$$\text{algo 1: } \theta = \theta - \epsilon \underbrace{\left( \mathbb{E}_{D_{P[1]}} [\partial_\theta l(x, y; \theta)] + \mathbb{E}_{D_{P[2]}} [\partial_\theta l(x, y; \theta - \epsilon g_1)] \right)}_{g_1}$$

$$\begin{aligned} \text{algo 2: } \theta &= \theta - \epsilon \left( \frac{1}{2} \sum_{i=1}^2 \mathbb{E}_\theta \left[ \partial_\theta l(x, y; \theta) \right] - 2\epsilon \partial_\theta \langle G_1, G_2 \rangle \right) \\ &\quad \left. \begin{aligned} \theta \in \mathbb{R}^d && \underbrace{\left( \frac{\partial x_i}{\partial \theta_j} \right)}_{n \times d} G_i && \text{Matrix} \\ \langle X(\theta + \Delta \theta), Y(\theta + \Delta \theta) \rangle &\approx \langle X(\theta) + \underbrace{\nabla X(\theta) \cdot \Delta \theta}_{n \times d}, Y(\theta) + \underbrace{\nabla Y(\theta) \cdot \Delta \theta}_d \rangle = \\ &= \langle X(\theta), Y(\theta) \rangle + \underbrace{(\nabla Y(\theta)^T \cdot \nabla X(\theta) + X(\theta)^T \nabla Y(\theta)) \cdot \Delta \theta}_{\nabla_\theta \langle X(\theta), Y(\theta) \rangle^T} + \tilde{o}(\|\Delta \theta\|) \end{aligned} \right) \end{aligned}$$

$$\text{So, in our case } \nabla_\theta \langle G_1, G_2 \rangle = \mathbb{E}_{D_2} \left[ \nabla_\theta^2 l(x, y; \theta) \right]^T G_1 + \mathbb{E}_{D_1} \left[ \nabla_\theta^2 l(x, y; \theta)^T \right] G_2 \quad \text{Hessian}$$

(3.3) Inner loop of algo 1

$$L_i := L(D_{P[i]}, \hat{\theta}_i), \quad g_i = \nabla L_i(\hat{\theta}_i), \quad \tilde{\theta}_{i+1} = \hat{\theta}_i - \epsilon g_i$$

$$g_{i,1} = \nabla L_i(\theta_1) \quad h_{i,1} = \nabla^2 L_i(\theta_1) \quad \text{at the start of the inner loop}$$

$$(1.3.1) \text{ Show } g_i = \nabla L_i(\theta_1) + \nabla^2 L_i(\theta_1)(\tilde{\theta}_i - \theta_1) + O(\lambda^2) =$$

symmetric

↑ Taylor for  
 $\nabla L_i(\tilde{\theta}_i)$

$$= \lambda \sum_{j=1}^{i-1} g_j - O(\lambda)$$

$$(1.3.2) \quad = g_{i,1} - \lambda h_{i,1} \sum_{j=1}^{i-1} g_j + O(\lambda^2) = g_{i,1} - \lambda h_{i,1} \sum_{j=1}^{i-1} g_{j,1} + O(\lambda^2)$$

↑ Taylor for  $g_j \cdot g_{i,1} = \lambda(\dots) + O(\lambda^2)$

Two domains

$$(1.3.3) \text{ Order: } \Phi_1, \Phi_2 \rightarrow L_1 = L_1(\theta_1) \quad g_1 = \nabla L_1(\theta_1) \quad \tilde{\theta}_2 = \theta_1 - \lambda g_1$$

$$L_2 = L_2(\tilde{\theta}_2) \quad g_2 = \nabla L_2(\tilde{\theta}_2) \quad \tilde{\theta}_3 = \tilde{\theta}_2 - \lambda g_2$$

Outer update

$$\Delta \theta = \theta_1 - \tilde{\theta}_3 = \theta_1 - \tilde{\theta}_2 + \tilde{\theta}_2 - \tilde{\theta}_3 = \lambda(g_1 + g_2) \Theta$$

$$1.3.2. \quad G_1(\theta_1) \quad G_2(\theta_1) \quad \nabla G_2(\theta_1)$$

$$(1.3.4) \quad \Theta \lambda \underbrace{(g_{1,1} + g_{2,1})}_{g_1} - \lambda h_{2,1} \underbrace{g_{1,1}}_{g_2}$$

(1.3.5) What would be the value of  $g_{1,1}, g_{2,1}, h_{2,1}$  if the domains had been considered in order  $\Phi_2, \Phi_1$ ?

$$g_{1,1} = G_2(\theta_1) \quad g_{2,1} = G_1(\theta_1) \quad h_{2,1} = \nabla G_1(\theta_1)$$

iteration indices

$$(1.3.6) \quad E_p[g_1 + g_2] = \frac{1}{2} (G_1(\theta_1) + G_2(\theta_1) - \lambda \nabla G_2(\theta_1) \cdot G_1(\theta_1)) + \frac{1}{2} (G_1(\theta_1) + G_2(\theta_1) -$$

$$+ O(\lambda^2))$$

$$- \lambda \nabla G_1(\theta_1) \cdot G_2(\theta_1) = G_1(\theta_1) + G_2(\theta_1) - \frac{\lambda}{2} (\nabla G_2(\theta_1) \cdot G_1(\theta_1) + \nabla G_1(\theta_1) \cdot G_2(\theta_1))$$

$$\frac{1}{2}(G_1 + G_2) \quad // \text{ see (3.2)} \quad + O(\lambda^2)$$

$$E_p[\Delta \theta] = \lambda E_p[g_1 + g_2] = \lambda \lambda \bar{G} - \frac{\lambda^2}{2} \partial_{\theta_1} (G_1(\theta_1), G_2(\theta_1)) + O(\lambda^2)$$

$$(1.3.7) \text{ For algo 2 we had } \Delta \theta = \bar{G} - 2 \lambda \nabla_{\theta} (G_1(\theta_1), G_2(\theta_1))$$

So, for  $\lambda = 2\sqrt{\kappa}$  it is actually an approximation of order  $O(\kappa)$ .

# Exercise (2022)

© Théo Jalabert

$w$  - r.v., vector of a NN

(1) Objective maximize  $p(w|X, Y) \propto p(Y|X, w)p(w)$

$D = (X, Y)$  training data

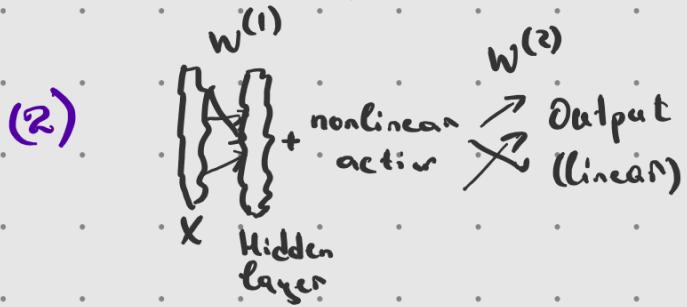
$$\text{loss } L_R(w) = -\ln p(Y|X, w) - \ln p(w) = L(w) + R(w)$$

Prior is Gaussian  $p(w) = \mathcal{N}(w; 0, \sigma^2 I)$

$$(1.1) R(w) = -\ln \frac{1}{(2\pi)^d/2} e^{-\frac{\sigma^2}{2}\|w\|^2} = \underbrace{\frac{d}{2} \ln(2\pi)}_{C} + \underbrace{\frac{\sigma^2}{2}\|w\|^2}_{\lambda w^T w} = C + \lambda w^T w$$

(1.2) Effect of  $R(w)$  = regularization. Prevents  $w$  from being too big  $\rightarrow$  controls the model complexity.

From Bayesian point of view it's a prior assumption that  $\|w\|$  is of order  $\sigma$ .



(2.1) Scale inputs  $x' = ax \rightarrow$  outputs for the first layer

$$w'^{(1)} = \frac{1}{a} w^{(1)} \rightarrow R(w) = \lambda w'^{(1)T} w^{(1)} + \lambda w'^{(2)T} w^{(2)} \Leftrightarrow$$

$$w'^{(2)} \text{ doesn't change} \Leftrightarrow R(w') = \lambda \cdot \frac{1}{a^2} \|w^{(1)}\|^2 + \lambda \|w^{(2)}\|^2$$

Same situation when scaling outputs

(2.2)  $p(w) = p(w^{(1)}) p(w^{(2)}) \quad p(w^{(i)}) = \mathcal{N}(w^{(i)}, 0, \sigma^2 I)$

$$R(w) = -\ln p(w^{(1)}) p(w^{(2)}) = \underbrace{\frac{d_1 + d_2}{2} \ln(2\pi)}_{\text{const}} + \underbrace{\frac{\sigma^2}{2} \|w^{(1)}\|^2}_{\lambda_1} + \underbrace{\frac{\sigma^2}{2} \|w^{(2)}\|^2}_{\lambda_2}$$

© Théo Jalabert 

(2.3) Same problem as in (2.1):

$$w = (w_1, \dots, w_d)$$

(3) Gaussian mixture  $p(w) = \prod_i p(w_i)$   $w_i = \sum_{j=1}^m \tilde{\pi}_j N(w_i; \mu_j, \sigma_j^2)$   
 $\sum \tilde{\pi}_j = 1, \tilde{\pi}_j > 0$

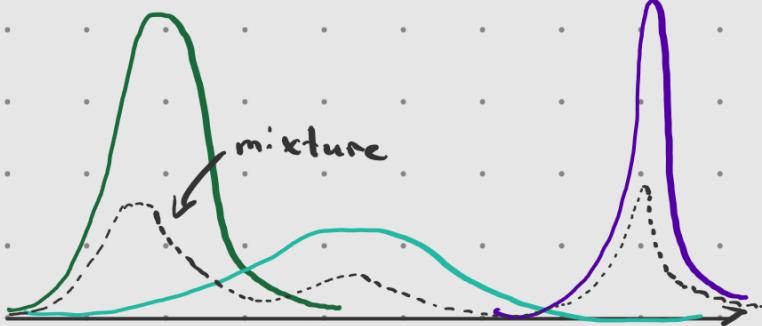
$$R(w) = -\sum_i \ln \left( \sum_{j=1}^m \tilde{\pi}_j N(w_i; \mu_j, \sigma_j^2) \right)$$

$\downarrow$  classical loss

$$L_R(w) = L(w) + \lambda R(w), \lambda \in \mathbb{R}^+$$

$$\tilde{\pi}_i \approx \frac{1}{3}$$

(3.1)



(3.2)  $\mu_j, \sigma_j^2, \tilde{\pi}_j$  fixed  $L_R(w) \rightarrow \min_w$

Effect of the constraint on the weights?

$$L_R(w) = L(w) + \lambda^0 R(w) \rightarrow \min_w$$

can be considered as a conditional minimization of  $L(w)$

under constraints  $R(w) \leq 0$  (in our case  $\sum_i \ln \left( \sum_{j=1}^m \tilde{\pi}_j N(w_i; \mu_j, \sigma_j^2) \right) \geq c$ )

The value of the constraint is greater when all the weights are close to one of the mixture centers  $\mu_i$ .

(3.3.) Training concerns all the parameters.

$$p(j|w) = \pi_j(w) = \frac{\tilde{\pi}_j N(w; \mu_j, \sigma_j^2)}{\sum_{k=1}^m \tilde{\pi}_k N(w; \mu_k, \sigma_k^2)}$$

$w$  is scalar!

$$(3.3.1) \quad \frac{\partial R}{\partial w_i} = -\partial w_i \ln \sum_k \pi_k N(w_i; \mu_k, \sigma_k^2) = \sum_k \pi_k \frac{\partial w_i}{\partial w_i} N(w_i; \mu_k, \sigma_k^2) = \sum_k \pi_k \delta_k(w_i)$$

$$= \left\{ \begin{array}{l} p = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \\ p' = -\frac{(x-\mu)}{\sigma^2} \cdot p \end{array} \right\} = \sum_k \frac{(w_i - \mu_k)}{\sigma_k^2} \delta_k(w_i) \rightarrow$$

$\rightarrow w_i$  tends to approach to  $\mu_k$  with greater posterior proba  $\delta_k(w_i)$  and lower variance  $\sigma_k^2$

$$(3.3.2) \quad \frac{\partial R}{\partial \mu_k} = - \sum_i \partial \mu_k \ln \left( \sum_k \pi_k N(w_i; \mu_k, \sigma_k^2) \right) =$$

$$= - \sum_i \frac{\pi_k \partial \mu_k N(w_i; \mu_k, \sigma_k^2)}{\sum_k \pi_k N(w_i; \mu_k, \sigma_k^2)} = \sum_i \frac{(\mu_k - w_i)}{\sigma_k^2} \delta_k(w_i) \rightarrow$$

$\rightarrow \mu_k$  will approach to the weights  $w_i$  such that  $j^{th}$  component of mixture is probable for them.

$$(3.3.3) \quad \frac{\partial R}{\partial \sigma_k} = - \sum_i \frac{\pi_k \partial \sigma_k N(w_i; \mu_k, \sigma_k^2)}{\sum_k \pi_k N(w_i; \mu_k, \sigma_k^2)} = -\left(\frac{x-\mu_k}{\sigma_k}\right) \varphi\left(\frac{x-\mu_k}{\sigma_k}\right)$$

$$\frac{d}{d\sigma} \left( \frac{1}{\sigma} \varphi\left(\frac{x-\mu_k}{\sigma}\right) \right) = -\frac{1}{\sigma^2} \varphi\left(\frac{x-\mu_k}{\sigma}\right) + \frac{1}{\sigma^3} \varphi'\left(\frac{x-\mu_k}{\sigma}\right) \cdot \left(-\frac{x-\mu_k}{\sigma^2}\right) =$$

$$= -\frac{1}{\sigma^2} P + \frac{(x-\mu_k)^2}{\sigma^3} P = \frac{P}{\sigma^2} \left( \frac{(x-\mu_k)^2}{\sigma^2} - 1 \right)$$

$$\Leftrightarrow \sum_i \frac{1}{\sigma_k^2} \left( 1 - \left( \frac{w_i - \mu_k}{\sigma_k} \right)^2 \right) \delta_k(w_i)$$

$\sigma_k$  will decrease if there are many points very close to  $\mu_k$ .  
 $\sigma_k$  should be positive.

Possible solutions:

- Set  $\sigma_{min} > 0$  and at each step k take  $\sigma_k \vee \sigma_{min}$  in order to keep it positive
- Transform  $\sigma_k = e^{\alpha_k}$ , Optimization w.r.t.  $\alpha_k$

# Convolutional NN (2021)

© Théo Jalabert

Image  $X \in \mathbb{R}^{3 \times 3}$ ,  $W \in \mathbb{R}^{2 \times 2}$  convolution filter,  $Z \in \mathbb{R}^{2 \times 2}$  transformed

image:  $X = \begin{pmatrix} x_1 & x_2 & x_3 \\ x_4 & x_5 & x_6 \\ x_7 & x_8 & x_9 \end{pmatrix}$   $W = \begin{pmatrix} w_1 & w_2 \\ w_3 & w_4 \end{pmatrix}$   $Z = \begin{pmatrix} z_1 & z_2 \\ z_3 & z_4 \end{pmatrix}$

(1)  $z_1 = w_1 x_1 + w_2 x_2 + w_3 x_4 + w_4 x_5$

$$z_2 = w_1 x_2 + w_2 x_3 + w_3 x_5 + w_4 x_6$$

$$z_3 = w_1 x_4 + w_2 x_5 + w_3 x_7 + w_4 x_8$$

$$z_4 = w_1 x_5 + w_2 x_6 + w_3 x_8 + w_4 x_9$$

(2)  $z = \begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{pmatrix} = \text{Flat}(Z)$   $x = \text{Flat}(X)$ ,  $w = \text{Flat}(W)$ ,  $z = wx$

$$\begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{pmatrix} = \underbrace{\begin{pmatrix} w_1 w_2 0 & w_3 w_4 0 & 0 & 0 \\ 0 w_1 w_2 0 & w_3 w_4 0 & 0 & 0 \\ 0 0 0 w_1 w_2 0 & w_3 w_4 0 & 0 & 0 \\ 0 0 0 0 w_1 w_2 0 & w_3 w_4 0 & 0 & 0 \end{pmatrix}}_{w^T} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_9 \end{pmatrix}$$

(3)  $y = w^T z$  transposed convolution

$$w^T z = \begin{pmatrix} w_1 z_1 \\ w_2 z_1 + w_1 z_2 \\ w_2 z_2 \\ w_3 z_1 + w_1 z_3 \\ w_4 z_1 + w_3 z_2 + w_2 z_3 + w_1 z_4 \\ w_4 z_2 + w_2 z_4 \\ w_3 z_3 \\ w_4 z_3 + w_3 z_4 \\ w_4 z_4 \end{pmatrix}$$



(4) How it helps? To build a decoder and increase the shape from  $2 \times 2$  to  $3 \times 3$ .

(5) Padding  $Z$ ,  $\text{Stride} = 1$

$$\begin{matrix} 0 & 0 & 0 & 0 \\ 0 & z_1 & z_2 & 0 \\ 0 & z_3 & z_4 & 0 \\ 0 & 0 & 0 & 0 \end{matrix}$$

$$* \quad \begin{matrix} w_4 & w_3 \\ w_2 & w_1 \end{matrix}$$

$$\begin{matrix} y_1 & y_2 & y_3 \\ y_4 & y_5 & y_6 \\ y_7 & y_8 & y_9 \end{matrix}$$

*W'*

© Théo Jalabert

*Jalabert*

So, transpose convolution in (3) = Padding + convolution with reversed  $W' = (w_4, w_3, w_2, w_1)$

(6) Stride=2

$$\begin{matrix} x_1 & x_2 & x_3 & x_4 \\ x_5 & x_6 & x_7 & x_8 \\ x_9 & x_{10} & x_{11} & x_{12} \\ x_{13} & x_{14} & x_{15} & x_{16} \end{matrix} * \begin{matrix} w_1 & w_2 \\ w_3 & w_4 \end{matrix} = \begin{matrix} z_1 & z_2 \\ z_3 & z_4 \end{matrix}$$

$$(6.1) \quad Z = \left( \begin{matrix} w_1 w_2 & 0 & 0 & w_3 w_4 & 0 & 0 \\ 0 & 0 & w_1 w_2 & 0 & 0 & w_3 w_4 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ w_1 w_2 & 0 & 0 & w_3 w_4 & 0 & 0 \\ 0 & 0 & w_1 w_2 & 0 & 0 & w_3 w_4 \end{matrix} \right) \cdot X$$

(6.2) Decoder

$$\begin{matrix} 0 & 0 & 0 & 0 \\ 0 & z_1 & z_2 & 0 \\ 0 & z_3 & z_4 & 0 \\ 0 & 0 & 0 & 0 \end{matrix} * \begin{matrix} w_1 & w_2 \\ w_3 & w_4 \end{matrix} = \begin{matrix} \text{hatched} & \text{hatched} \\ \text{hatched} & \text{hatched} \end{matrix}$$

NN & conditional density mixture

$$\mathcal{D} = (x_i, y_i)_{i=1}^N, \quad x_i \in \mathbb{R}^n, \quad y_i \in \mathbb{R}$$

$$X \in \mathbb{R}^{n \times N} \quad X = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_N \end{pmatrix} \quad Y = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}$$

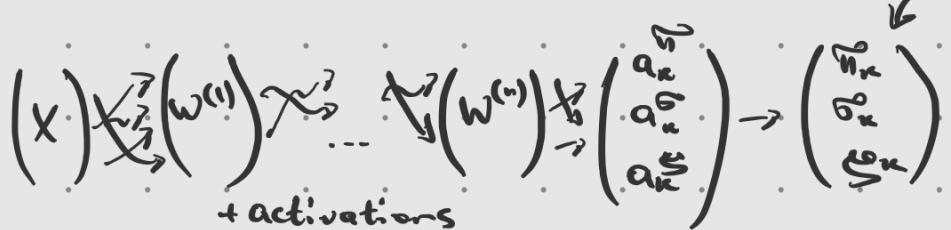
Multinomial  $p(y|x) = \sum_{k=1}^K \pi_k(x) p_k(y|x)$ ,  $p_k(y|x) = \mathcal{N}(y|g_k(x), \Sigma_k(x))$

$\pi_k(x) \in [0,1] \quad \sum \pi_k(x) = 1.$

$g_k(x)$ ,  $\Sigma_k(x)$ ,  $\pi_k(x)$  will be computed by a NN

$\pi_k(x) = \frac{e^{a_k}}{\sum_i e^{a_i}} \quad \Sigma_k(x) = e^{a_k^S} > 0 \quad g_k(x) = a_k^G$

(1) How can be implemented?



(2)  $L(w) = \sum_{i=1}^N L_i(w) \quad L_i(w) = \log \left( \sum_{k=1}^K \pi_k(x_i) p_k(y_i|x_i) \right) / p_k$

Derive  $\partial_{a_i^k} L_i(w)$  for  $k \in \{1, 2, \dots, K\}$

$$\begin{aligned} \partial_{a_i^k} L_i(w) &= \frac{1}{\sum_{k=1}^K \pi_k(x_i) p_k(y_i|x_i)} \cdot \sum_k p_k \left( \underbrace{\frac{-e^{a_k}}{(\sum_i e^{a_i})^2} \cdot e^{a_i} + \frac{e^{a_k}}{\sum_i e^{a_i}} \delta_{ki}}_{\pi_k} \right) \\ &= \frac{\sum_k p_k (\delta_{ki} \cdot \pi_k - \pi_k \pi_i)}{\sum_k \pi_k p_k} \end{aligned}$$

$$\partial_{a_i^k} L_i(w) = \frac{1}{\sum_{k=1}^K \pi_k p_k} \pi_k \cdot \underbrace{\partial_{a_i^k} p_k}_{\text{et les exo de 2022}} \cdot e^{a_i}$$

$$\partial_{a_i^k} b_i(w) = \frac{1}{\sum_{k=1}^K \pi_k p_k} \pi_k \cdot \underbrace{\partial_{a_i^k} p_k}_{\text{ex}}$$

(3) conditional mean  $E[y|x] = ?$  conditional variance  $S^2(x) = ?$

$$E[y|x] = \sum \pi_k(x) E[\mathcal{N}(y|g_k(x), \Sigma_k(x))] = \sum \pi_k(x) g_k(x)$$

$$\mathbb{E}^2(x) = \mathbb{E}\left[\left(y - \mathbb{E}[y|x]\right)^2 | x\right] = \sum_k \pi_k(x) \mathbb{E}\left[\left(\xi_k - \mathbb{E}[\xi_k|x]\right)^2\right] =$$

© Théo Jalabert



$$= \sum_k \pi_k(x) \left( \underbrace{\mathbb{E}\left[\left(\xi_k - \mathbb{E}[\xi_k|x]\right)^2\right]}_{\text{Var}_k(x)} + \left(\mathbb{E}[\xi_k|x] - \mathbb{E}(y|x)\right)^2 \right) = \sum_k \pi_k(x) (\text{Var}_k(x) + (\mathbb{E}[\xi_k|x] - \mathbb{E}[y|x])^2)$$