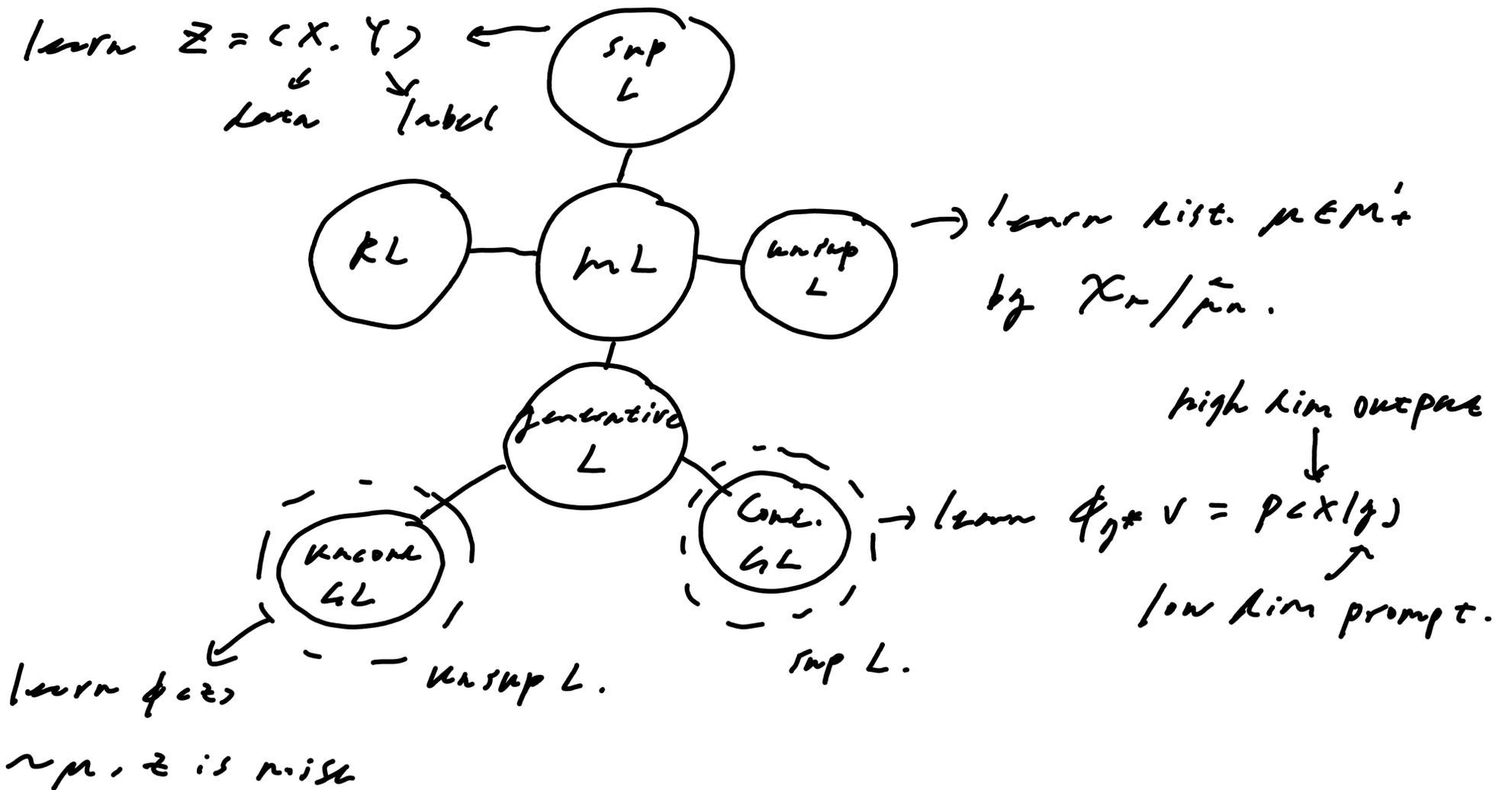


Supervised Learning



Remark: Supervised learning is easier than unsupervised one since the dim. of label is low in common but data of unsupervised has high dim.

In supervised learning, we observe $Z_j = (Y_j, X_j) = (L, A, \Phi) \rightarrow (k^s \times k^L, B_{k^s \times k^L})$ in sample $X_n = (Z_1, \dots, Z_n)$, Z_k i.i.d. where Y_j is label/output and X_j is covariate/input.

Our goal is not to learn distribution of Z_j

rather to learn the dist. of $Y|X$, i.e.
propose a dist. $\hat{\mu}_{Y|X}$ from $\mathcal{X}_n = \{z_1, \dots, z_n\}$
to get close to $\mu_{Y|X} = P(Y|X)$.

ex. Y is r.v. of products customer will buy
 X is shopping history ... info. of customer
 \Rightarrow interested in $P(Y = \text{"soap"} | X) = ?$

(1) Regular conditional probn.:

Lemma: For $\tilde{A} \in \mathcal{A}$. $\Rightarrow E(\cdot | \tilde{A}) = \text{Proj.} : L^2(\mathcal{R}, \mathcal{A}, P; \mathbb{R}^k) \rightarrow L^2(\mathcal{R}, \tilde{A}, P; \mathbb{R}^k)$. ortho. proj.

This fixes $E(\cdot | \tilde{A})$. P -a.s.

RMK: Recall i) $\text{Proj } X = \arg \min_{Y \in L^2(\tilde{A})} E\|X - Y\|^2$.

ii) $\text{Proj} \circ \text{Proj} = \text{Proj}$. iii) $\text{Proj} = i_{L^2(\tilde{A})}$

iv) $E\|X - \text{Proj } Y\|^2 = E\| \text{Proj } X - Y \|^2$

(By ortho. decomposition of X, Y)

Lemma: If $Y \in \sigma(X)$, $Y: \Omega \rightarrow \mathbb{R}^k$ r.v. $X: \Omega \rightarrow \mathbb{R}^k$ r.v. Then: $\exists g: (\mathbb{R}^k, B_{\mathbb{R}^k}) \rightarrow (\mathbb{R}^k, B_{\mathbb{R}^k})$
measurable. st. $Y = g(X)$. g is unique

\mathbb{P}_x -a.s. i.e. $\mathcal{Y}' = \mathcal{Y}$. \mathbb{P}_x -a.s. $\Rightarrow g(x) = \mathcal{Y}$. \mathbb{P} -a.s.

Def: i) $\{P_{Y|X}(A)\}_{A \in \mathcal{A}}$ is regular conditional probability (r.c.p.) if it's induced by a

Markov kernel $k(x, B) = P_{Y|X=x}(B)$. \mathbb{P}_x -a.s. $x \in \mathcal{X}$

a) $B \mapsto k(x, B)$ is p.m. for $\forall x \in \mathcal{X}$.

b) $x \mapsto k(x, B)$ is $\mathcal{B}_{\mathcal{X}}$ -measurable. $\forall B \in \mathcal{A}$.

ii) $P_{Y|X=x}(B)$ admits a regular version if

$\bar{P}_{Y|X}(\cdot)$ is r.c.p. and $P_{Y|X} = \bar{P}_{Y|X}$ a.s.

Remark: Note if let $P_{Y|X}(A) = \mathbb{E}[\mathbb{1}_{\{Y \in A\}} | \mathcal{G}(X)]$.

then: $P_{Y|X}(\cdot)$ also only

satisfies property of p.m. \mathbb{P} -a.s.

But the null set N will depend

on $\{\mathcal{A}_n\}$ we choose if it's not r.c.p.

\Rightarrow it may not be p.m. for \mathbb{P}_x -a.s. x .

Supervised learning is to learn Markov kernel.

Ex: $f(y|x) = \frac{1}{(\sqrt{2\pi}\sigma)^2} e^{-\frac{1}{2}\left(\frac{y-x}{\sigma}\right)^2} k_y \Rightarrow$

$k(x, B) = \int_B f(y|x) k_y$ is a Markov kernel.

for linear regression

Thm. $\mathcal{Z} = \langle Y, X \rangle = \langle \mathcal{X}, \mathcal{A}, \mathbb{P} \rangle \rightarrow \langle \mathcal{X}^s \times \mathcal{X}^t, \mathcal{B}_{\mathcal{X}^s \times \mathcal{X}^t} \rangle$ s.t.
 $Y \in L^2$ -r.v. \Rightarrow There exists a regular
 version $P_{Y|X}(A)$.

Lemma. The dist. of $\mathcal{Z} = \langle Y, X \rangle = P_{\mathcal{Z}}$ is uniquely
 determined by its marginal P_X and r.c.p.
 $P_{Y|X=X}(\cdot)$. i.e. we have: $\forall B \in \mathcal{B}_{\mathcal{X}^s \times \mathcal{X}^t}$.

$$P_{\mathcal{Z}}(B) = \int_{\mathcal{X}^t} \left(\int_{\mathcal{X}^s} \mathbb{I}_B(\eta, x) \kappa P_{Y|X=X}(\eta) \right) \kappa P_X(x).$$

Lemma. If $\mathcal{Z} = \langle X, Y \rangle$ has density $f_{\mathcal{Z}}(x, \eta)$. let

$$f_X(x) = \int_{\mathcal{X}^t} f_{\mathcal{Z}}(x, \eta) \kappa \eta \text{ and } \int_{\mathcal{X}^t}$$

$$f(\eta|x) = \begin{cases} f_{\mathcal{Z}}(x, \eta) / f_X(x) & \text{if } f_X(x) > 0 \\ \mathbb{I}_{\{0,1\}^s}(\eta) & \text{otherwise.} \end{cases}$$

Then: $P_{Y|X=X}(B) = \int_{\mathcal{X}^t} \mathbb{I}_B(\eta) f(\eta|x) \kappa \eta$ is
 a r.c.p. of Y over $X=x$.

(2) Divergence:

Next we want to introduce a divergence to
 measure success of learning from $\hat{P}_{Y|X, x_n}$
 to true Markov kernel $P_{Y|X}$.

Def: For two Markov kernel $\mu_{|x=x}, V_{|x=x}$ and divergence $k(\cdot, \cdot)$ on $\mathcal{M}_1^+(\mathbb{R}^d)$, $\Sigma \subset \mathbb{R}^d$
 $\mapsto k(\mu_{|x=x} \parallel V_{|x=x}) \in (\mathbb{R}^d, \mathcal{B}_{\mathbb{R}^d})$ is measurable for $\forall \mu_{|x}, V_{|x} \in \mathcal{K}_1^+(\mathbb{R}^d, \mathcal{B}_{\mathbb{R}^d})$ space of Markov kernels. Set:

$$D_p(\mu_{|x} \parallel V_{|x}) = D_{p, x, \mu}(\mu_{|x} \parallel V_{|x})$$

$$= \mathbb{E}_x(k(\mu_{|x} \parallel V_{|x}))^p)^{\frac{1}{p}}$$

$$= \left(\int_{\mathbb{R}^d} k(\mu_{|x=x} \parallel V_{|x=x})^p \mathcal{P}_x(dx) \right)^{\frac{1}{p}}, 1 \leq p < \infty.$$

$$D_\infty(\mu_{|x} \parallel V_{|x}) = \text{ess sup}_{x \in \mathbb{R}^d} k(\mu_{|x=x} \parallel V_{|x=x})$$

Remark: i) $p \uparrow$, then sensitivity of $k(\mu_{|x=x} \parallel V_{|x=x})$ at each pt. $x \uparrow$.

ii) $D_{p, TV}$ is weaker than $D_{p, KL}$.

CD will mostly inherit the strength of k . But it's not from Pinsker's)

iii) If $x \sim P_x$ but we let $x \sim \tilde{P}_x$ knowing inference. Suppose that

$$D_{1, x \sim P_x, \mu}(\mu_{|x} \parallel \hat{\mu}_{|x, n}) \rightarrow 0, \text{ a.s.}$$

$$\text{Note } D_{1, x \sim \tilde{P}_x, \mu}(\mu_{|x} \parallel \hat{\mu}_{|x, n}) =$$

$$\mathbb{E}_{x \sim P_x} \left(\mathcal{L}(\mu_{1x} \| \tilde{\mu}_{1x,n}, \frac{AP_x}{AP_x}) \right) \leq$$

$$\left\| \frac{AP_x}{AP_x} \right\|_n D_{1,x \sim P_x, \mathcal{L}}(\mu_{1x} \| \tilde{\mu}_{1x,n}) \xrightarrow{a.s.} 0 \text{ if } \left\| \frac{AP_x}{AP_x} \right\|_n < \infty$$

So the real space can still be covered in training

Def. Denote $V_{1x} P_x(B) \stackrel{\Delta}{=} \int \int I_B(y, x) dV_{1x=x}(y) dP_x(x)$

For $V_{1x}, \mu_{1x} \in \mathcal{K}_1(\mathbb{R}^L, \mathcal{B}_{\mathbb{R}^L})$. We have:

$$D_{1,x, KL}(\mu_{1x} \| V_{1x}) = \mathcal{L}_{KL}(\mu \| V). \text{ where } \mu_{1x} = \mu_{1x} P_x, V = V_{1x} P_x.$$

Pr.: In supervised learning we always know μ and V .

Pf.: i) We first prove: $\frac{\mathcal{L}_V}{\mathcal{L}_\mu}(y, x) \stackrel{a.s.}{=} \frac{\mathcal{L}_{V_{1x=x}}}{\mathcal{L}_{\mu_{1x=x}}}(y)$.

$$a) \mu \ll V \Rightarrow \mu_{1x=x} \ll V_{1x=x}, P_x \text{ a.s. } x.$$

otherwise, $\exists B_1, \text{ s.t. } P_x(B_1) > 0$ satisfy:

$$\exists B_2(x), \text{ s.t. } \mu_{1x=x}(B_2(x)) > 0 = V_{1x=x}(B_2(x))$$

$$\text{Set } B = \{(y, x) : x \in B_1, y \in B_2(x)\}$$

$$\Rightarrow \mu(B) = \int \mu_{1x=x}(B_2(x)) I_{B_1}(x) dP_x > 0$$

$$\text{But } V(B) = \int V_{1x=x}(B_2(x)) I_{B_1}(x) dP_x = 0.$$

\Rightarrow contradiction.

$$\begin{aligned}
 b) \text{ And we see } & \int I_B(\eta, x) \frac{\mu_{V|x=x}(\eta) \mu(\eta, x)}{\mu_{\mu|x=x}} \\
 &= \int \left(\int I_B(\eta, x) \frac{\mu_{V|x=x}(\eta) \mu_{\mu|x=x}(\eta)}{\mu_{\mu|x=x}} \right) \mu_{P_x}(x) \\
 &= \int \left(\int I_B(\mu_{V|x=x}(\eta)) \mu_{P_x}(x) \right) = V(B)
 \end{aligned}$$

$$\int_0 : \mu_{V|x=x} \ll \nu_{V|x=x}, \mu \text{-a.s. } x \Rightarrow \mu \ll \nu.$$

$$\text{and } \frac{\mu}{\mu_{\mu|x=x}}(\eta, x) = \frac{\mu_{V|x=x}(\eta)}{\mu_{\mu|x=x}}. \mu \text{-a.s.}$$

$$\begin{aligned}
 2) \text{ KL}(\mu || \nu) & \stackrel{i)}{=} - \int \log \left(\frac{\mu_{V|x=x}(\eta)}{\mu_{\mu|x=x}}(\eta) \right) \mu(\eta, x) \\
 &= \int \left(- \int \log \left(\frac{\mu_{V|x=x}(\eta)}{\mu_{\mu|x=x}}(\eta) \right) \mu_{\mu|x=x}(\eta) \right) \mu_{P_x}(x) \\
 &= D_{l.x.KL}(\mu_{V|x=x} || \nu_{V|x=x}).
 \end{aligned}$$

(3) Framework:

Denote $\mathcal{J} \subseteq \mathcal{K}_1^+(\mathcal{K}^t, \mathcal{B}_{\mathcal{K}^t})$ is set of target

Markov kernels. $\mathcal{X}_n = (z_1, \dots, z_n)$ is samples

and $\tilde{\mu}_{n|x} \stackrel{\Delta}{=} \tilde{\mu}_{n|x} \circ \mathcal{X}_n$. For D div. We'll

require: $\mathcal{X} \in \mathcal{K}^{(s+t)n} \mapsto D(\mu_{V|x} || \hat{\mu}_{n,V|x}(x)) \in \mathcal{R}^+$

is $\mathcal{B}_{\mathcal{K}^{(s+t)n}}$ -measurable.

Def: $\mathcal{K}_n \subseteq \mathcal{K}_1^+(\mathcal{K}^t, \mathcal{B}_{\mathcal{K}^t})$ hypothesis space of

Markov kernel. $z_j = (Y_j, X_j) \stackrel{i.i.d.}{\sim} \mu = \mu_{V|x} \mu_{P_x}$

Let $D = D_{p,x,d}$ div. on \mathcal{X}_i^+ .

i) Empirical risk function $\hat{\mathcal{L}}_n = \hat{\mathcal{L}}_n(v_{|X}, \mathcal{X}_n)$

satisfies:

a) $\mathbb{R}^{(k+1)n} \ni \mathcal{X}_n \mapsto \hat{\mathcal{L}}_n(v_{|X}, \mathcal{X}_n)$ is measurable

for $\forall v_{|X} \in \mathcal{H}_n$

b) $\exists h(\cdot)$ on $\mathcal{M}_i^+(\mathcal{C}_n) \subset \mathbb{R}^+$ s.t.

$\mathcal{C}_n \ni \hat{\mathcal{L}}_n(v_{|X}, \mathcal{X}_n) + h_n(\mu) \xrightarrow{pr} D_{p,x,d}(\mu_{|X} || v_{|X})$

$\forall \mathcal{C}_n$. s.t. $v_{|X} \in \mathcal{H}_n$.

ii) $\hat{\mathcal{L}}_n$ is unbiased if $\mathbb{E}_{\mathbb{Z}}(\mathcal{C}_n \ni \hat{\mathcal{L}}_n(v_{|X}, \mathcal{X}_n) + h_n(\mu)) = D_{p,x,d}(\mu_{|X} || v_{|X})$.

iii) $\hat{\mu}_{n|X}$ is supervised ERM-learning for $\hat{\mathcal{L}}_n$ if $\hat{\mu}_{n|X} \in \operatorname{argmin}_{v_{|X} \in \mathcal{H}_n} \hat{\mathcal{L}}_n(v_{|X}, \mathcal{X}_n)$.

Note we still have error decomposition:

$$0 \in D_{p,x,d}(\mu_{|X} || \hat{\mu}_{n|X}) \leq \varepsilon_{n,\text{mod}}(\mu_{|X}) + \varepsilon_{n,\text{learn}} + 2\varepsilon_{n,\text{samp}}$$

$$\varepsilon_{n,\text{mod}}(\mu) = \inf_{v_{|X} \in \mathcal{H}_n} D_{p,x,d}(\mu_{|X} || v_{|X});$$

$$\varepsilon_{n,\text{learn}} = c_n \left(\hat{\mathcal{L}}_n(\hat{\mu}_{n|X}, \mathcal{X}_n) - \inf_{v_{|X} \in \mathcal{H}_n} \hat{\mathcal{L}}_n(v_{|X}, \mathcal{X}_n) \right);$$

$$\varepsilon_{n,\text{samp}} = \sup_{v_{|X} \in \mathcal{H}_n} \left| D(\mu_{|X} || v_{|X}) - (c_n \hat{\mathcal{L}}_n(v_{|X}, \mathcal{X}_n) + h_n(\mu)) \right|.$$

Cor. For $\mathcal{J} \subset \mathcal{K}$, $|\mathcal{K}| < \infty$. Then: \mathcal{J} is PAC-learnable.

Cor. \mathcal{K} is opt. $\mathcal{J} \subset \mathcal{K}$. For $\hat{\mathcal{I}}_n(V_{1:n}, \gamma_n) = \sum_{j=1}^n \ell(z_j | V_{1:n})$ with $\ell(z | V_{1:n}) \leq k \ell(z) \in \mathcal{L}'(\mathbb{R}^{s+k}, \mu)$, $\forall \mu = \mu_{1:n} P_x$, $\mu_{1:n} \in \mathcal{K}$. If $V_{1:n} \in \mathcal{K} \mapsto \ell(z | V_{1:n})$ is $D_{P_x}(\cdot, \|\cdot\|)$ -cont.

Then \mathcal{J} is learned by ERM learner $\hat{\mu}_n$.

Thm. $\mathcal{J}, \mathcal{K} = \mathcal{K}_1^T(\mathbb{R}^d, \mathcal{B}_{\mathbb{R}^d}) \cap \{ \mu \mid \mu_{1:n}(y) = f_{V_{1:n}}(\eta(x)) \}$

If $\ell(z | V_{1:n}) \stackrel{\Delta}{=} -\log f_{V_{1:n}}(\eta(x)) \in \mathcal{L}'(\mathbb{R}^{s+k}, \mu)$

$\forall \mu = \mu_{1:n} P_x$, $\forall \mu_{1:n} \in \mathcal{J}$, $\forall V_{1:n} \in \mathcal{K}$. And let:

$\hat{\mathcal{I}}_n(V_{1:n}, \gamma_n) = \sum_{j=1}^n \ell(z_j | V_{1:n})$. Then:

i) $\hat{\mathcal{I}}_n$ is unbiased ERMF with $\mathcal{L}_n = n^{-1}$ and

$$h_n(\mu) = h(\mu) = \int_{\mathbb{R}^{s+k}} \log f_{\mu_{1:n}=x}(\eta(x)) \mu_{1:n=x}(\eta) \mu P_x(x)$$

ii) If $X \sim f_X(x) dx$, $f_V(z) \stackrel{\Delta}{=} f_{V_{1:n}}(\eta(x)) f_X(x)$

$$\bar{\mathcal{J}} \stackrel{\Delta}{=} \{ \mu = \mu_{1:n} P_x \mid \mu_{1:n} \in \mathcal{K} \}$$

$$\bar{\mathcal{K}} \stackrel{\Delta}{=} \{ V = V_{1:n} P_x \mid V_{1:n} \in \mathcal{J} \}. \text{ Then:}$$

$\hat{\mu}_n$ is ERM learner w.r.t. $\mathcal{L}(\cdot, \|\cdot\|)$. $\hat{\mu}_n \in$

$$\arg \min_{\nu \in \bar{\mathcal{K}}} - \sum_{j=1}^n \log f_{\nu}(z_j) \Leftrightarrow \hat{\mu}_{n,1:n} \text{ r.c.p. of}$$

$\hat{\mu}_n$ w.r.t $\sigma(x)$ is ERM learner for \hat{I}_n

Pf: i) $\hat{I}_n = \mathbb{E}_z (L_n(\hat{I}_n(V_{1X}, \mathcal{X}_n) + h_n(\mu))) \stackrel{\text{Lem.}}{=} - \int \log(f_{V_{1X}}(y|x) / f_{\mu_{1X}}(y|x)) d\mu_{X=X}(y) dP_X(x)$
 $= \int A_{f,L}(\mu_{1X=X} \parallel V_{1X=X}) dP_X(x) = D_{1X, f, L}(\mu_{1X} \parallel V_{1X})$

ii) Note $f_v(z) = f_{V_{1X}}(y|x) f_X(x)$
 $= \hat{I}_n(V_{1X}, \mathcal{X}_n) - \hat{I}_n(\log f_X(x_j))$
 \Rightarrow The 2nd term doesn't depend on V_{1X}
with last Lem. of (2). and def of $\hat{\mu}_n$. We obtain the corresponding.

Remark: We can see supervised learning is as unsupervised learning with part of list.

(4) Common Principles for models:

For $(P_{y|x})$ Markov kernel. Next, we want to know what kind type of Y :

a) continuous b) discrete: \mathbb{N} , $\{1, 2, \dots, q\}$, $\{0, 1\}$...

$P_{y|x, \theta}$ should be constructed to give all cond. proba. and label space of Y .

e.g. i) (Binary labels)

$B(p) = \text{Bernoulli dist. on } \{0,1\}$. i.e.

$$B(p)(1) = p = 1 - B(p)(0).$$

Note all dist. on $\{0,1\}$ can be repr
in $\{B(p(x)) \mid p(x) \text{ is measurable}\}$ all Markov
kernels over \mathbb{R}^k and label space = $\{0,1\}$.

ii) (Conti. dist.: Classical regression)

$$\eta \in \mathbb{R}. P(\eta | x) = (\sqrt{2\pi}\sigma)^{-1} \exp\left(-\frac{1}{2}\left(\frac{\eta - m(x)}{\sigma}\right)^2\right)$$

RMK: If we use it to model all conti.

dist. \Rightarrow error will be from:

- a) Model error (may not be gaussian)
- b) Deviation of $m(x)$

Steps:

1) identify the data structure of labels:
i.e. conti. or discrete?

2) choose a parametric family of dist. $\{M_\xi:$
 $\xi \in \Xi, \xi \mapsto M_\xi \subset \mathcal{B}\}$ is measurable. $\forall B \in \mathcal{A}$.

3) choose a parametric class of functions $\{m_\theta:$

$\theta \in \mathbb{R}^L \mapsto \exists m_\theta$ is measurable}. Derive the hypothesis space $\mathcal{H} = \{ \mu_{m_\theta} : \theta \in \mathbb{R}^L, \mu_{x=x}(B) := \mu_{m_\theta}(B) \text{ is Markov kernel} \}$.

4) Choose a divergence $\rho(\cdot, \cdot)$ (may KL) and corresponding ERF \tilde{L}_n (may neg. log likelihood)

5) Train with ERM.

① Binary classification: (logistic regression)

label space = $\{0, 1\}$. $\sigma(z) = e^z / (1 + e^z)$.

$m_\theta(x) = \bar{\theta}^T x + \theta_0$. $x \in \mathbb{R}^L$. $\theta = (\theta_0, \bar{\theta}) \in \mathbb{R}^{L+1}$.

(θ_0 is bias and $\bar{\theta}$ is weight.)

Set $\beta(p(x)) \equiv 1 = p(x) = \sigma(m_\theta(x))$ and use

Neg. log likelihood:

$\tilde{L}_n((y_j, x_j)_{j=1}^n) = - \sum_{j=1}^n y_j \log p(x_j) + (1 - y_j) \log(1 - p(x_j))$

\Rightarrow Choose $y_j = \arg \max_{y \in \{0, 1\}} \beta(p(x_j)) \geq y_j$

Or we can use other functions:

$m_\theta(x) = \varphi_{\theta_L}^{(L)} \circ \dots \circ \varphi_{\theta_1}^{(1)}$. [$\varphi_{\theta_k}^{(k)}$] is layers. $\theta =$

$(\theta_1, \dots, \theta_L)$. L is depth. $\theta_j = (\bar{\theta}_j, \theta_j^0)$. s.t.

$\bar{\theta}_j \in \mathbb{R}^{L+1 \times L_j}$. $\theta_j^0 \in \mathbb{R}^{L_j}$.

We set $\phi_{\theta_j}^{(j)}(x^{(j-1)}) = \mathcal{X}(\theta_j^T x^{(j-1)} + \theta_j^0) \in \mathbb{R}^{k_j}$

where $\mathcal{X}: \mathbb{R}^{k_j} \rightarrow \mathbb{R}^1$, $\mathcal{X}(z) = (\mathcal{X}(z_1) \dots \mathcal{X}(z_{k_j}))^A$

for $z \in \mathbb{R}^{k_j}$, $\forall j$. activation function

e.g. $\mathcal{X}(z_j) = e^{z_j} / (1 + e^{z_j})$ sigmoid function

and $\text{ReLU}(z) = \max(0, z)$.

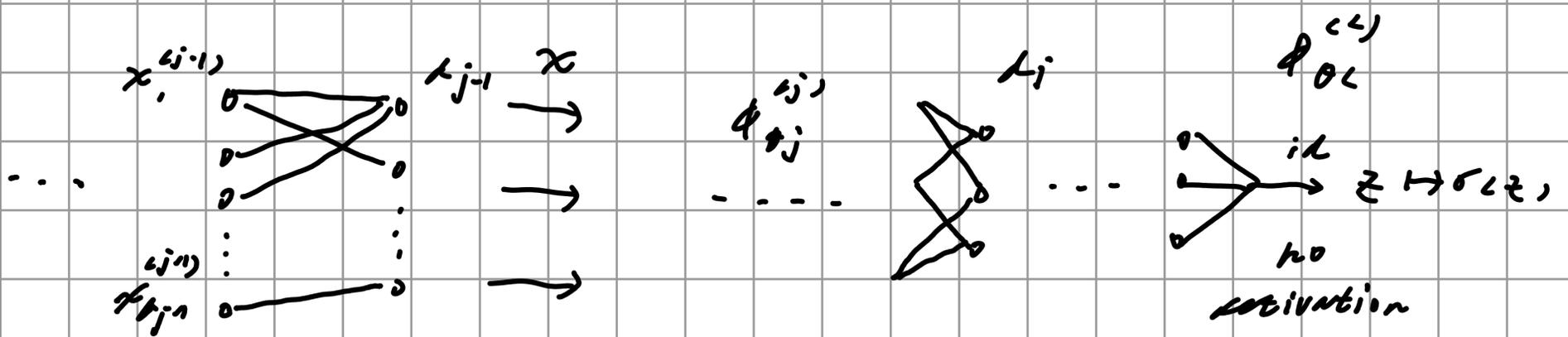
Remark: i) We don't use polynomials generally

because of comp. ineff. and hard to optimize. ($|Z(\phi(x))|$ is large)

ii) $\max\{, \}$ isn't differentiable. So:

optimal can't be trained to get

Fully connected Neural Network (FCNN):



Multi-class Classification:

$q \in \mathcal{L} = \{1, 2, \dots, \ell\}$. Set $\text{softmax}(z)_q =$

$$e^{z_q} / \sum_{y \in \mathcal{L}} e^{z_y}, \quad z \in \mathbb{R}^{\ell}$$

RMK: Note $\beta + c \mathbb{1}$ still corresponds the same
 $(P(\eta|\beta))$. So it's not identifiable
 Hence we generally force $\beta_1 = 0$.

Next, we do a logistic regression: $\theta = (\bar{\theta}, \theta_0)$

$\bar{\theta} \in \text{Mat}_{k \times d-1}(\mathbb{R})$, $\theta_0 \in \mathbb{R}^{d-1}$. $m_{\theta}(x) = \bar{\theta}^T x + \theta_0$.

$\mathcal{M}_{X \times \mathcal{Y}}^{\theta}(g) = P_{\theta}(g|x) := \text{softmax}(0, m_{\theta}(x))_g$

So: $\sum_g P_{\theta}(g|x) = 1 \Rightarrow$ it's dist. on \mathcal{L} .

And we get Markov kernel $\mathcal{M}_X(\cdot)$.

RMK: All dist. on \mathcal{L} can be represented
 by one β .

$$\begin{aligned} \tilde{L}_n(\langle \eta_j, x_j \rangle) &= - \sum_{j=1}^n \log(\text{softmax}(0, \bar{\theta}^T x_j + \theta_0)_{\eta_j}) \\ &= - \sum_{j=1}^n \sum_{i=1}^k \delta_{i, \eta_j} \log \Pi. \end{aligned}$$

③ Count regression:

$P_{\text{oi}}(y|\eta) := e^{-\eta} \eta^y / y!$. $y \in \mathcal{L} = \mathbb{N}$. See:

$$m_{\theta}(x) = \begin{cases} \bar{\theta}^T x + \theta_0 & \text{or} \\ \eta_{\theta}(x) \cdot (\text{FCNN}) \end{cases}$$

mult Markov kernel $\mu_{1|x=x}^{\theta}(y) = P_{\theta}(y|x) =$
 $e^{-m_{\theta}(x)} m_{\theta}(x)^y / y!$

$$\Rightarrow \hat{\mathcal{L}}((y_j, x_j), \theta) = \sum_j -m_{\theta}(x_j) + y_j \log m_{\theta}(x_j) - \log y_j!$$

④ Regression:

Markov kernel is $(\sqrt{2\pi}\sigma)^{-1} e^{-\frac{1}{2}(\frac{y-m_{\theta}(x)}{\sigma})^2} / \sigma$

$$\begin{aligned} \hat{\mathcal{L}}_n((y_j, x_j), \theta, \sigma) &= -\sum -\frac{1}{2} \left(\frac{y_j - m_{\theta}(x_j)}{\sigma} \right)^2 - \log \sqrt{2\pi}\sigma \\ &= (2\sigma^2)^{-1} \sum (y_j - m_{\theta}(x_j))^2 + n \log \sqrt{2\pi}\sigma \\ &= \frac{n}{2\sigma^2} \mathcal{L}_n^{MSE}(\theta) + n \log \sqrt{2\pi}\sigma. \end{aligned}$$

Remark: For $m_{\theta}(x) = \theta x + \theta_0$. We call it ordinary

Least square (OLS) linear regression:

if $\theta_0 = 0$, then $\theta^* = (X^T X)^{-1} X^T Y$ is the

MLE if $|X^T X| \neq 0$. Where $X = (x_1, \dots, x_n)$, $Y =$

(y_1, \dots, y_n) . Also: $\hat{\sigma}_x^2 = \mathcal{L}_n^{MSE}(\theta^*)$ is optimal.