# 1 Universal Approximation Theorem

$$f_m(x;\theta) = \sum_{k=1}^{m} a_k \sigma(b_k x + c_k), \qquad \sigma(x) = \max(x, 0)$$

$$P_M f = \sum_{k=0}^{M} f(x_k) \, t\left(\frac{x - x_k}{h}\right), \qquad t(x) = \max(1 - |x|, 0)$$

(a) $P_M f$ can be represented as two-layer ReLu network

**Proof.**

$$P_M f(x) = \sum_{k=0}^{M} f(x_k) \max\left(1 - \frac{|x - x_k|}{h}, 0\right)$$

$$= \sum_{k=0}^{M} f(x_k)\left(\max\left(\frac{x - (x_k - h)}{h}, 0\right) + \max\left(\frac{x - (x_k + h)}{h}, 0\right) - 2\max\left(\frac{x - x_k}{h}, 0\right)\right)$$

$$= \sum_{k=0}^{M} f(x_k) \sigma\left(\frac{x - (x_k - h)}{h}\right) + \sum_{k=1}^{M} f(x_k) \sigma\left(\frac{x - (x_k + h)}{h}\right)$$
$$+ \sum_{k=1}^{M}(-2 f(x_k)) \sigma\left(\frac{x - x_k}{h}\right) \qquad \#$$

(b) $f \in C[0,1]$. $\forall \varepsilon > 0$, $\exists f_m$, $\sup_{[0,1]} |f(x) - f_m(x;\theta)| \leq \varepsilon$

**Proof.** Using (a), it suffices to show $\forall \varepsilon > 0$, $\exists M$, $\sup_{[0,1]} |f(x) - P_M f(x)| < \varepsilon$

$f$ is uniform continuous on $[0,1]$. Let $\delta > 0$ be $|x_1 - x_2| \leq \delta \Rightarrow |f(x_1) - f(x_2)| \leq \varepsilon$

Take $M > \frac{1}{\delta}$. $|f(1) - P_M f(1)| = 0$,

Now suppose $x \in [x_k, x_{k+1})$, $k = 0, \cdots M-1$

$$|f(x) - P_M f(x)| = \left| f(x) - f(x_k) \, t\left(\frac{x - x_k}{h}\right) - f(x_{k+1}) \, t\left(\frac{x - x_{k+1}}{h}\right)\right|$$

$$= \left| f(x) - f(x_k)\left(1 - \frac{x - x_k}{h}\right) - f(x_{k+1})\left(1 - \frac{x_{k+1} - x}{h}\right)\right|$$

$$= \left| \left(1 - \frac{x - x_k}{h}\right)(f(x) - f(x_k)) + \left(1 - \frac{x_{k+1} - x}{h}\right)(f(x) - f(x_{k+1}))\right|$$

$$\leq \left(1 - \frac{x - x_k}{h}\right)|f(x) - f(x_k)| + \left(1 - \frac{x_{k+1} - x}{h}\right)|f(x) - f(x_{k+1})|$$

$$\leq \varepsilon\left(2 - \frac{x - x_k + x_{k+1} - x}{h}\right) = \varepsilon \qquad \text{We are done.} \quad \#$$

(c) $f(x) = x^2$ $\varepsilon \in (0,1)$ $\exists\, m \le C\varepsilon^{-\frac{1}{2}}$ s.t. $\sup_{[0,1]} |f(x) - f_m(x;\theta)| \le \varepsilon$

**Proof.** Let $M = \lceil \frac{1}{2\sqrt{\varepsilon}} \rceil$ Consider $f_m(x;\theta) = P_M f$, here $m = 3(M+1) \le 6 + \frac{3}{2\sqrt{\varepsilon}}$

$$\le \frac{10}{\sqrt{\varepsilon}}$$

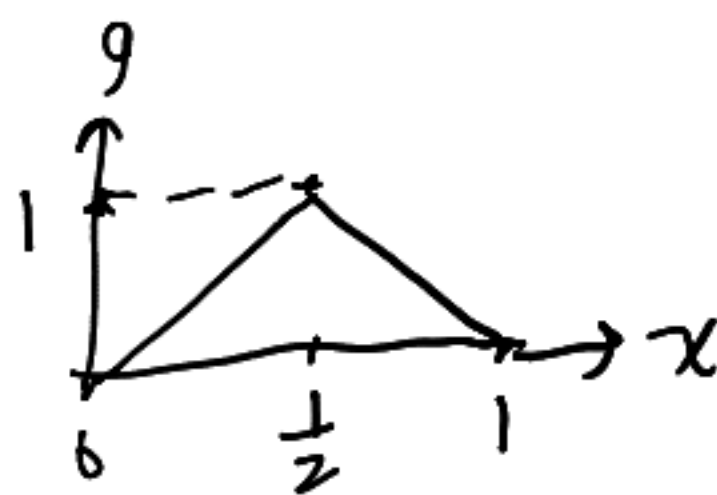Suppose $x \in [x_k, x_{k+1})$ Let $x = x_k + \Delta$, $\Delta < h$

$|f(x) - f_m(x;\theta)| = |f(x) - P_M f(x)|$

$= \left| (1 - \frac{x - x_k}{h})(x^2 - x_k^2) - (1 - \frac{x_{k+1} - x}{h})(x_{k+1}^2 - x^2) \right|$

$= \left| (1 - \frac{\Delta}{h})(\Delta^2 + 2\Delta x_k) - \frac{\Delta}{h}(h-\Delta)(2x_k + h + \Delta) \right|$

$= \left| (1 - \frac{\Delta}{h})(\Delta^2 + 2\Delta x_k - 2\Delta x_k - \Delta h - \Delta^2) \right|$

$= (h - \Delta)\Delta \le \frac{h^2}{4} \le \varepsilon$

$\sup_{[0,1]} |f(x) - f_m(x;\theta)| \le \varepsilon$.  We are done. #

2  Approximate $x^2$ with DNN

$f^*(x) = x^2$ in $[0,1]$  $g(x) = t(2x-1)$  which maps $[0,1]$ to $[0,1]$

$$g_s(x) = g_0 \cdots \circ g(x)$$
$$\underbrace{\phantom{g_0 \cdots \circ g(x)}}_{s}$$



(a) Show that $\sup_{[0,1]} |P_{2^\mu} f^*(x) - f^*(x)| \le \frac{C}{2^{2\mu}}$

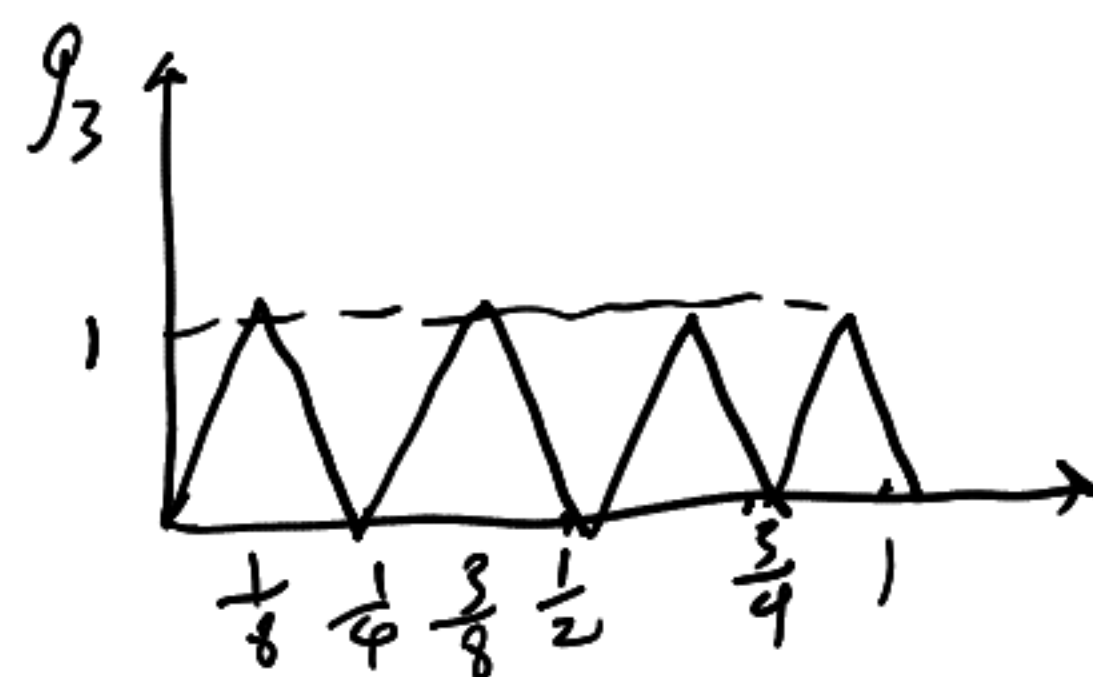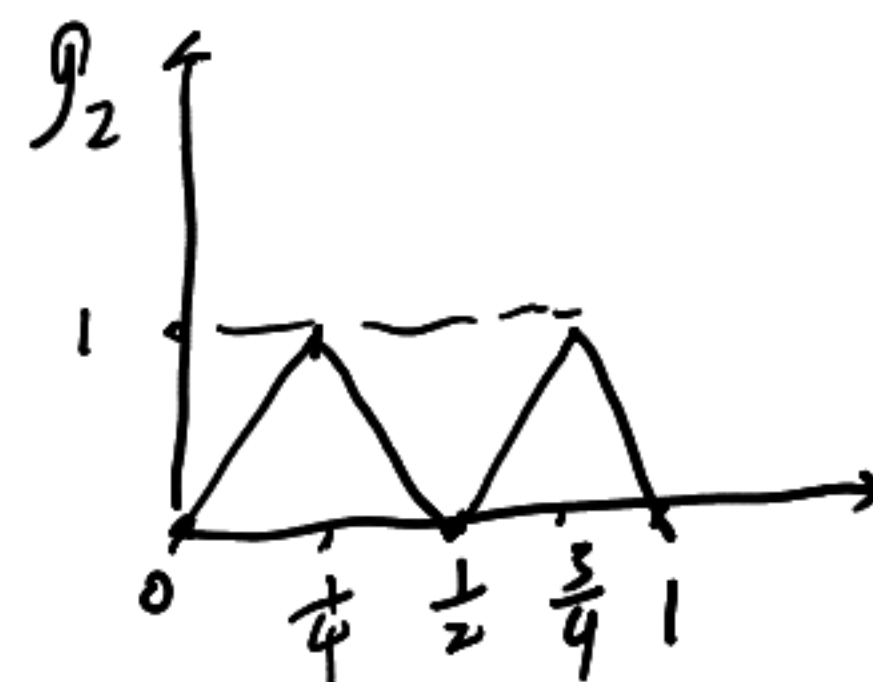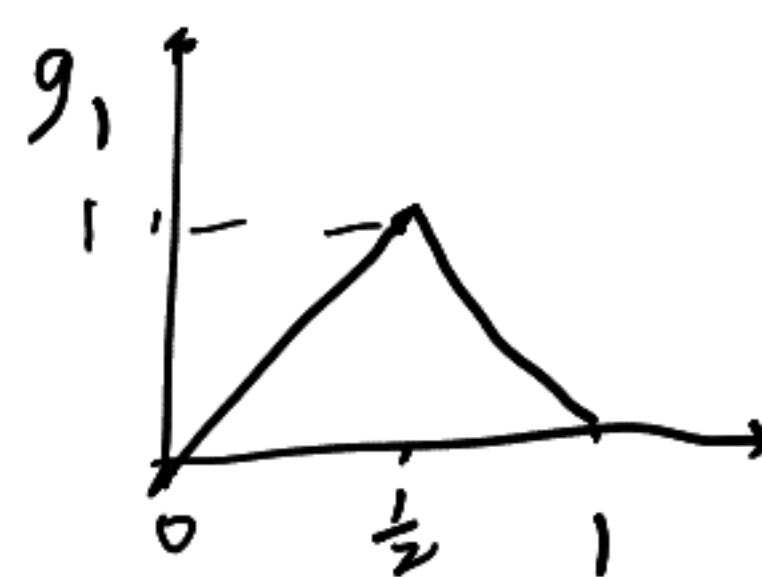**Proof.**  Similar to (c) of last home work,

$|f(x) - P_M f(x)| = \left| (1 - \frac{\Delta}{h})(\Delta^2 + 2\Delta x_k) - \frac{\Delta}{h}(h-\Delta)(2x_k + h + \Delta) \right|$

(where $x = x_k + \Delta$, $0 \le \Delta < h$, $k = 0, 1, \cdots M-1$)

$= \left| (1 - \frac{\Delta}{h})(\Delta^2 + 2\Delta x_k - 2\Delta x_k - \Delta h - \Delta^2) \right|$

$= (h - \Delta)\Delta \le \frac{h^2}{4} = \frac{1}{4M^2} = \frac{1}{4} \cdot \frac{1}{2^{2\mu}}$  #

(b) $l = 2, 3, \cdots$

$$P_{2^{l+1}} f^*(x) - P_{2^l} f^*(x) = \frac{g_l(x)}{2^{2l}}, \quad \forall x \in [0,1].$$

<u>Proof.</u> By induction, it is easy to verify $g_l(x)$ is a piecewise-linear function of step $\frac{1}{2^l}$

with $g_l\left(\frac{2k+1}{2^l}\right) = 1$, $k = 0, \cdots 2^{l-1} - 1$

and $g_l\left(\frac{k}{2^{l-1}}\right) = 0$, $k = 0, \cdots 2^{l-1}$.

From definition, we also have

$P_{2^{l+1}} f^*(x) - P_{2^l} f^*(x)$ is piecewise-linear

function of step $\frac{1}{2^l}$

with $P_{2^{l+1}} f^*\left(\frac{k}{2^{l-1}}\right) - P_{2^l} f^*\left(\frac{k}{2^{l-1}}\right) = \left(\frac{k}{2^{l-1}}\right)^2 - \left(\frac{k}{2^{l-1}}\right)^2 = 0$
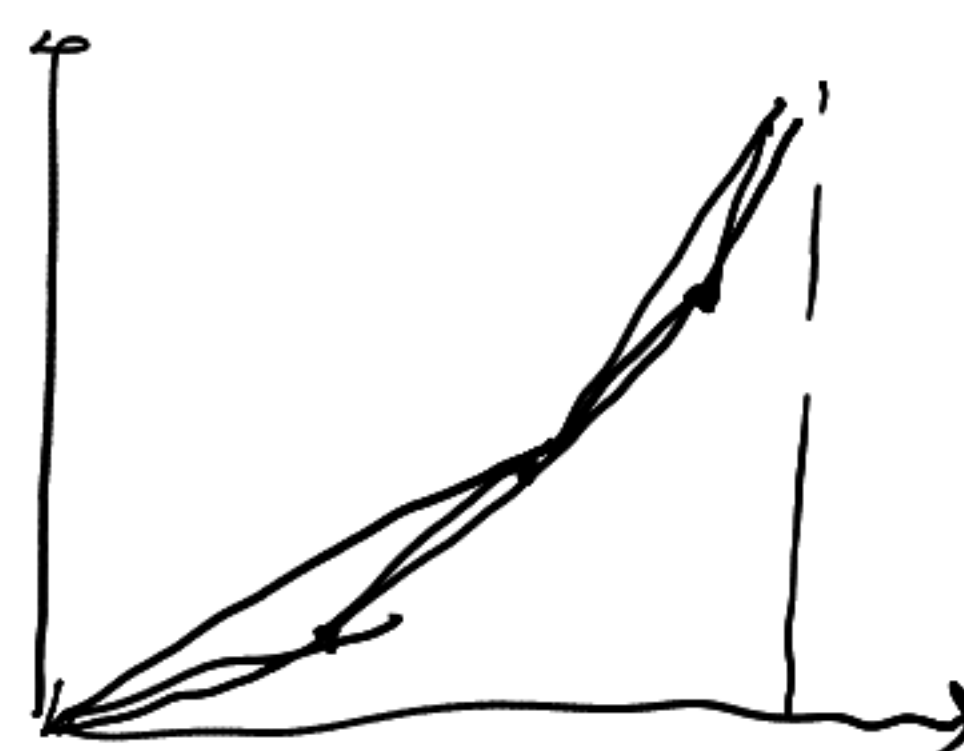
$k = 0, \cdots 2^{l-1}$

$P_{2^{l+1}} f^*\left(\frac{2k+1}{2^l}\right) - P_{2^l} f^*\left(\frac{2k+1}{2^l}\right)$

$= \frac{1}{2}\left(\frac{k}{2^{l-1}}\right)^2 + \frac{1}{2}\left(\frac{k+1}{2^{l-1}}\right)^2 - \left(\frac{2k+1}{2^l}\right)^2$
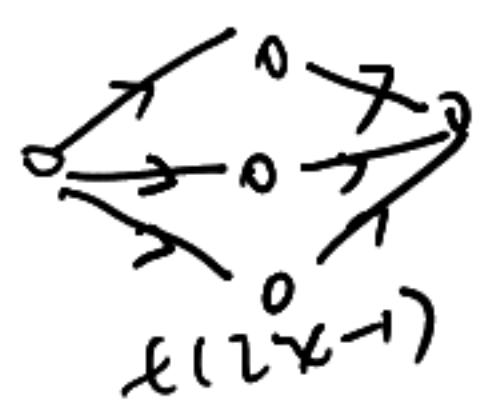
$= \frac{k^2}{2^{2l-1}} + \frac{k^2 + 2k + 1}{2^{2l-1}} - \frac{2k^2 + 2k + \frac{1}{2}}{2^{2l-1}} = \frac{1}{2^{2l}}$, $k = 0, \cdots 2^{l-1} - 1$

We have demonstrated $P_{2^{l+1}} f^*(x) - P_{2^l} f^*(x) = \frac{g_l(x)}{2^{2l}}$, $\forall x \in [0,1]$. #

(c) $P_{2^\ell} f^*$ can be represented as $O(\ell)$-layer NN with $O(1)$ width

Proof. $t(2x-1) = \max(2x-2,0) + \max(2x,0) - 2\max(2x-1,0)$
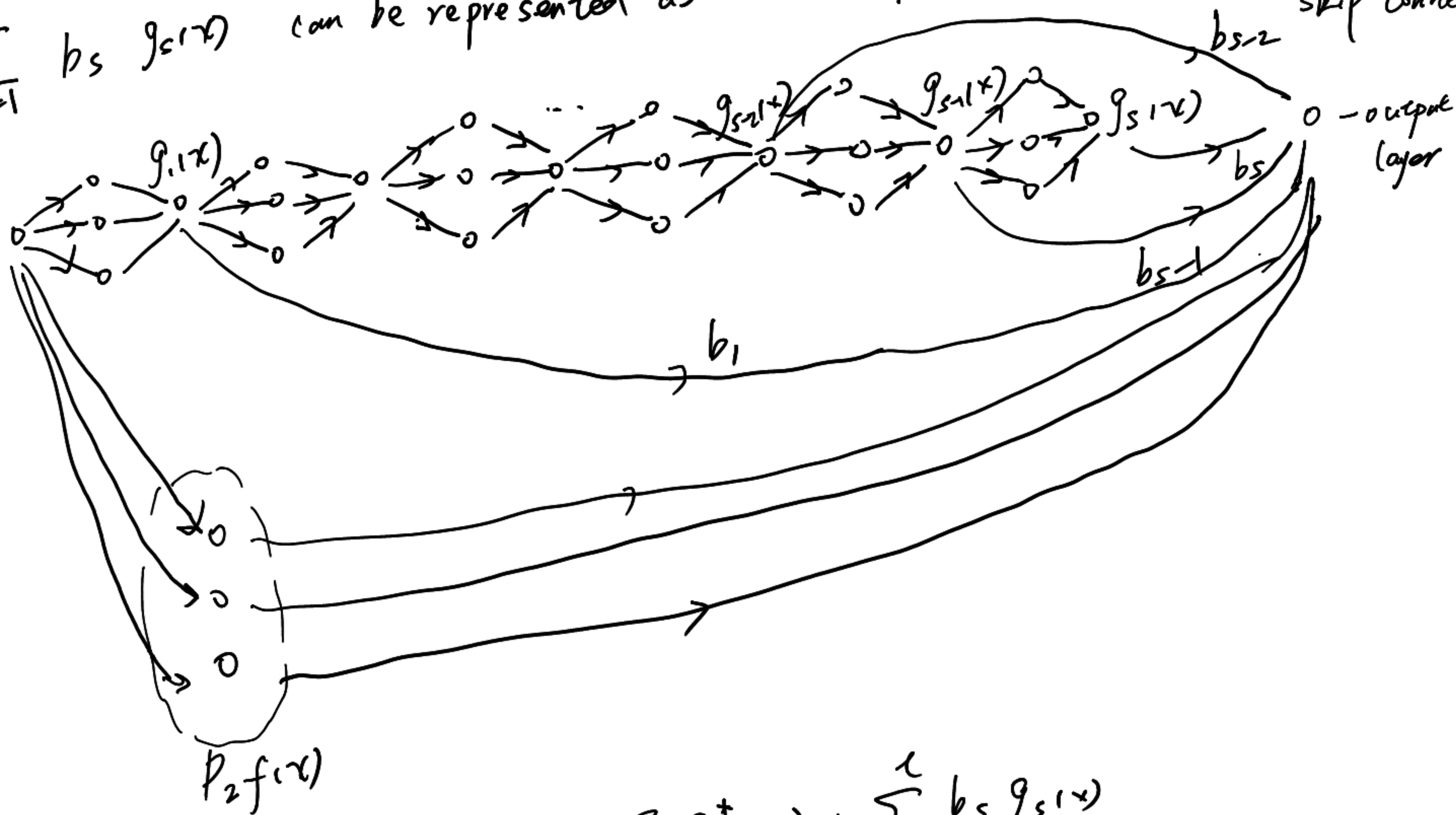
$$= \sigma(2x-2) + \sigma(2x) - 2\sigma(2x-1)$$



$t(2x-1)$

So $g_\ell(x) = g \circ \cdots \circ g(x)$ is $O(\ell)$ layer-NN with width 3.

$\underbrace{\phantom{g \circ \cdots \circ g}}_{\ell}$

From (b), $P_2 f^*(x) - P_{2^\ell} f^*(x) = \sum_{s=1}^{\ell} \frac{g_s(x)}{2^{2s}}$

$P_{2^\ell} f^*(x) = P_2 f^*(x) - \sum_{s=1}^{\ell} \frac{g_s(x)}{2^{2s}}$    Let $b_s = -\frac{1}{2^{2s}}$

$\sum_{s=1}^{\ell} b_s g_s(x)$ can be represented as $O(\ell)$ depth, $O(1)$ width NN with skip connection.



As in figure above, $P_{2^\ell} f^*(x) = P_2 f^*(x) + \sum_{s=1}^{\ell} b_s g_s(x)$

is a NN with width 6 and depth $2\ell+1$

with skip connection. #

## 3 Dropout for linear regression

(a) $f(x;\beta) = \beta^T x$    $\beta, x \in \mathbb{R}^d$    $\hat{R}(\beta) = \frac{1}{n}\sum_{i=1}^{n}(f(x_i;\beta) - y_i)^2$

$\tilde{\beta} = P\beta$    $w_j := \frac{1}{n}\sum_{i=1}^{n}x_{ij}^2$.

     Show that $\hat{R}_{drop}(\beta) = \hat{R}(\tilde{\beta}) + \frac{1-P}{P}\sum_{j=1}^{d}w_j\,\tilde{\beta}_j^2$

**Proof.**

$$RHS = \frac{1}{n}\sum_{i=1}^{n}(P\beta^T x_i - y_i)^2 + \frac{P(1-P)}{n}\sum_{j=1}^{d}\sum_{i=1}^{n}x_{ij}^2\beta_j^2$$

$$LHS = \mathbb{E}_{z\sim\pi}\left[\hat{R}(\beta\odot z)\right]$$

$$= \mathbb{E}_{z\sim\pi}\,\frac{1}{n}\sum_{i=1}^{n}(f(x_i;\beta\odot z) - y_i)^2$$

$$= \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}_{z\sim\pi}(f(x_i;\beta\odot z) - y_i)^2$$

$$= \frac{1}{n}\sum_{i=1}^{n}y_i^2 - \frac{2}{n}\sum_{i=1}^{n}y_i\,\mathbb{E}_{z\sim\pi}f(x_i;\beta\odot z) + \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}_{z\sim\pi}f^2(x_i;\beta\odot z)$$

$$= \frac{1}{n}\sum_{i=1}^{n}y_i^2 - \frac{2}{n}\sum_{i=1}^{n}y_i\,P\beta^T x_i + \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}_{z\sim\pi}\left(\sum_{j=1}^{d}z_j x_{ij}\beta_j\right)^2$$

$$= \frac{1}{n}\sum_{i=1}^{n}y_i^2 - \frac{2}{n}\sum_{i=1}^{n}y_i\,P\beta^T x_i + \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}_{z\sim\pi}\left(\sum_{j=1}^{d}z_j x_{ij}^2\beta_j^2 + 2\sum_{j<k}z_j z_k x_{ij}x_{ik}\beta_j\beta_k\right)$$

$$= \frac{1}{n}\sum_{i=1}^{n}y_i^2 - \frac{2}{n}\sum_{i=1}^{n}y_i\,P\beta^T x_i + \frac{1}{n}\sum_{i=1}^{n}P\sum_{j=1}^{d}x_{ij}^2\beta_j^2$$

$$+ \frac{2}{n}\sum_{i=1}^{n}\sum_{j<k}P^2 x_{ij}x_{ik}\beta_j\beta_k$$

$$RHS = LHS \iff \frac{1}{n}\sum_{i=1}^{n}P^2\left(\sum_{j=1}^{n}\beta_j x_{ij}\right)^2 + \frac{P-P^2}{n}\sum_{i=1}^{n}\sum_{j=1}^{d}x_{ij}^2\beta_j^2 =$$

$$\frac{1}{n}\sum_{i=1}^{n}P\sum_{j=1}^{d}x_{ij}^2\beta_j^2 + \frac{2}{n}\sum_{i=1}^{n}\sum_{j<k}P^2 x_{ij}x_{ik}\beta_j\beta_k$$

It suffices to prove $\forall$ $i\le i\le n$,

$$P\left(\sum_{j=1}^{d}\beta_j x_{ij}\right)^2 + (1-P)\sum_{j=1}^{d}x_{ij}^2\beta_j^2 = \sum_{j=1}^{d}x_{ij}^2\beta_j^2 + 2\sum_{j<k}P x_{ij}x_{ik}\beta_j\beta_k$$

     This holds obviously. #

(b) This shows dropout training of $f(x; \beta)$

is regularied training of $f(x; p\beta)$ ( ridge regression)

with controlling parameter $\lambda = \frac{1-p}{p}$

i.e. $\hat{R}_{dropout}(\beta) = \hat{R}(p\beta) + \lambda \cdot T(\beta)$

which explains why it improves generalization.

Parameter $p$ controls regularization parameter $\lambda$

If $p = 1$, no dropout

The smaller $p$ is, the stronger regularization is. #