- $\hat{h}_{n,\theta_b}$: trained single model on dataset with initialization $\theta_b$.
- $\{(x_i, \hat{s}(x_i))\}$, artificial dataset, where $\hat{s}(x) = \mathbb{E}_{\theta_b}\left[s_{\theta_b}(x)\right]$, $s_{\theta_b}$ is an $\theta_b$-initialized NN.
- $\varphi'_{n,\theta_b}$: auxillary network trained on $\{(x_i, \hat{s}(x_i))\}$, denote $\varphi_{n,\theta_b}(x) = \varphi'(x) - \hat{s}(x)$.
- $h^* = \hat{h}_{n,\theta_b} - \varphi_{n,\theta_b}$.

In regression:

$$AU = E(Var(y))$$

$$EU = var(E(y))$$

The key here is when generating artificial dataset, we have

$$\hat{s}(x) = E(S_{\theta_b}(x))$$

$$var(s) = var(S_{\theta_b}(x))$$

The single model $\hat{h}_{n,\theta_b}(x)$:

$$\hat{h}_{n,\theta_b}(x) = S_{\theta_b}(x) + \boxed{K(x,\underline{x})^T (K(\underline{x},\underline{x}) + \lambda_n n I)^{-1}} (y - S_{\theta_b}(\underline{x}))$$

$$h^*(x) = \hat{s}(x) + K(x,\underline{x})^T (K(\underline{x},\underline{x}) + \lambda_n n I)^{-1} (y - \hat{s}(\underline{x}))$$

$$\varphi'(x) = S_{\theta_b}(x) + K(x,\underline{x})^T (K(\underline{x},\underline{x}) + \lambda_n n I)^{-1} (\hat{s}(\underline{x}) - S_{\theta_b}(\underline{x}))$$

$$\Rightarrow \quad \hat{h}^* = \hat{h}_{n,\theta_b} - \varphi'(x) + \hat{s}(x).$$

$$\hat{h}^*(x) = \hat{s}(x) + k(x,x)^T(k(x,x) + \lambda n I)^{-1}(y - \hat{s}(X))$$

$$\downarrow$$

$$K$$

$$\text{var}(\hat{h}^*) = \text{var}\left(\hat{s}(x) + k(y - \hat{s}(X))\right)$$

$$= \text{var}\left(\hat{s}(x)\right) + \text{var}\left(k(y - \hat{s}(X))\right) + 2\text{cov}(\cdot,\cdot)$$

$$= \frac{1}{m}\text{var}(s(x)) + \frac{1}{m}k \cdot \text{var}(s(X)) \cdot k^T + 2\text{cov}(\cdot,\cdot)$$

$$\hat{s} = \frac{1}{m}\sum_{i=1}^{m} S_{\theta_{b_i}}(x)$$

$$\text{Var}(\hat{s}) = \frac{1}{m^2}\text{var}\left(\sum S_i(x)\right)$$

$$= \frac{1}{m} \cdot \text{var}(s)$$

$$\text{Var}(s(x)) \approx \frac{1}{m-1} \sum_{k=1}^{m} (s_{\theta k}(x) - \hat{s}(x))^2$$

$$\text{Var}(s(\underline{X})) \approx \frac{1}{m-1} \sum_{k=1}^{m} \left[ s_{\theta k}(\underline{X}) - \hat{s}(\underline{X}) \right] \cdot \left[ s_{\theta k}(\underline{X}) - \hat{s}(\underline{X}) \right]^{\top}$$

$$\text{Cov}(s(x), s(\underline{X})) \approx \frac{1}{m-1} \sum_{k=1}^{m} \left[ s_{\theta k}(x) - \hat{s}(x) \right] \left[ s_{\theta k}(\underline{X}) - \hat{s}(\underline{X}) \right]^{\top}$$

---

$$\text{Var}(\hat{h}^*) = \frac{1}{m} \cdot \Big[ \underset{①}{\text{Var}(s(x))} + \underset{②}{\underline{k} \cdot \text{Var}(s(\underline{X})) \underline{k}^{\top}} - \underset{③}{2k \, \text{Cov}(s(x), s(\underline{X}))} \Big]$$

① *test data*

② training data

③ measure relation between test data $x$ and training $\underline{X}$

① dimension reduction

② decomposition.
   $\text{Cov} = L^{\top} L$

③ Approximation

Epistemic.

$$\hat{h}^*(x) = \hat{s}(x) + k(x,x)^T (k(x,x) + \lambda_n nI)^{-1} (y - \hat{s}(X))$$

$$\downarrow$$
$$k$$

$$\left(\hat{h}^*(x) - \hat{s}(x)\right)^{|x|} = k^{|x n|} \cdot \left(y - \hat{s}(X)\right)^{n \times 1}$$

$$\downarrow$$
$$h \qquad\qquad\qquad t$$

$$\Rightarrow \quad h = k \cdot t$$

$$\Rightarrow \quad h \cdot t^T = k \cdot t \cdot t^T$$

$$k = h \cdot t^T (t \cdot t^T)^{-1}$$

$$\text{②} = k \cdot \text{var}(s(\underline{X})) \cdot k$$

$$= k(x, \underline{X})^T (k(\underline{X}, \underline{X}) + \lambda_n n I)^{-1} \cdot \text{var}(s(\underline{X}))$$

$$\cdot \left[ (k(\underline{X}, \underline{X}) + \lambda_n n I)^{-1} \right]^T \cdot k(x, \underline{X})$$

$$\text{var}(\hat{h}^*) = \frac{1}{m} \cdot \Big[$$

$$\overset{\textcircled{\tiny 0}}{\text{var}(s(x))} + \overset{\textcircled{\tiny 2}}{k \cdot \text{var}(s(X)) \, k^\top} - \overset{\textcircled{\tiny 3}}{2k \, \text{cov}(s(x), s(X))} \Big]$$

When training data is sufficient $\to \infty$ :

$\overset{\textcircled{\tiny 1}}{}$ $\text{var}(s(x))$ , not change

$\overset{\textcircled{\tiny 2}}{}$ $k \, \text{var}(s(X)) \, k^\top$ , $k(x, X)^\top (k(X, X) + \lambda_{nn} I)^{-1} \cdot \text{var}(s(X)) \cdot ((k(X, X) + \lambda_{nn} I)^{-1})^\top \cdot k(x, X)$

$\overset{\textcircled{\tiny 3}}{}$ $-2k \, \text{cov}(s(x), s(X))$ , $-2k(x, X) \cdot (k(X, X) + \lambda_{nn} I)^{-1} \text{cov}(s(x), s(X))$

---

$\textcircled{\tiny 2} + \textcircled{\tiny 3}$

$$= k \cdot \text{var}(s(X)) \cdot k^\top$$

$$- 2k \cdot \boxed{\text{cov}(s(x), s(X))}$$

For $K = k(x, \underline{X})(k(\underline{X}, \underline{X}) + \lambda n I)^{-1}$    $\sim k(x, \underline{X}) \cdot O(\frac{1}{n})$

$$\lim_{n \to \infty} \frac{1}{n} k(\underline{X}, \underline{X}) \to K$$

$(k(\underline{X}, \underline{X}) + \lambda n I)^{-1}$
$= \left( n \left( \frac{1}{n} k(\underline{X}, \underline{X}) + \lambda I \right) \right)^{-1}$
$= \frac{1}{n} \left( \frac{1}{n} k(\underline{X}, \underline{X}) + \lambda I \right)^{-1} \sim O(\frac{1}{n})$

if $x$ is a point in $\underline{X}$, then
$$k(x, \underline{X}) \approx k(x_i, \underline{X}) \sim O(1)$$

Then $k(x_i, \underline{X}) \cdot (k(\underline{X}, \underline{X}) + \lambda n I)^{-1}$

$\sim O(\frac{1}{n})$

worst case: $O(\frac{1}{\sqrt{n}})$

So

$$\frac{②+③}{Var(s(x))\, Var(s(\bar{X}))} = O\left(\frac{1}{\sqrt{n}}\right) \cdot O\left(\frac{1}{\sqrt{n}}\right) \cdot 1 \cdot \frac{1}{Var(s(x))}$$

$$- 2\, O\left(\frac{1}{\sqrt{n}}\right)$$

Now, only thing to do is:

$$Var(s(\bar{X})) \cdot O\left(\frac{1}{\sqrt{n}}\right)$$

if $\longrightarrow 0$ ✓