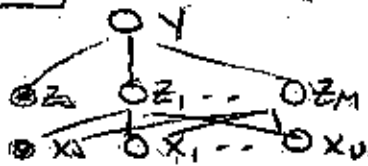


11/6/2013

TODAY

BIAS, VARIANCE,
RESIDUAL

TODAY:

$$E = \frac{1}{2} E_{st} \left[|y(\vec{x}) - t|^2 \right]$$

$$= \frac{1}{2} E \left[|y(\vec{x}) - \langle t | \vec{x} \rangle|^2 \right] + \frac{1}{2} E \left[|\langle t | \vec{x} \rangle - t|^2 \right]$$

DOES THIS REALLY
GO TO ZERO? HOW FAST?CONDITIONAL VARIANCE
OF t GIVEN x

$$+ \frac{1}{2} E \left[(y(\vec{x}) - \langle t | \vec{x} \rangle) (\langle t | \vec{x} \rangle - t) \right]$$

ASSUME THIS IS ≈ 0

$$y(\vec{x}) = \sum_{j=1}^M w_j \phi_j(\vec{x}) \xrightarrow[N \rightarrow \infty]{M \rightarrow \infty} \langle t | \vec{x} \rangle$$

BUT FOR FIXED N, M , $W = (\Phi^T \Phi)^{-1} \Phi^T T$

$$\Phi(n, j) = \phi_j(\vec{x}^{(n)}) \quad T(n) = t^{(n)}$$

SUPPOSE $E[\phi_j \phi_k] = \delta_{jk}$ THEN $\Phi^T \Phi \approx N I$

$$W \approx \frac{1}{N} \Phi^T T \approx \frac{1}{N} \sum_{n=1}^N \vec{\phi}^{(n)} t^{(n)}$$

$$y(\vec{x}) = y(\vec{\phi}) = \vec{\phi}^T \vec{w} = \vec{\phi}^T \left(\frac{1}{N} \sum_{n=1}^N t^{(n)} \vec{\phi}^{(n)} \right)$$

WE CAN TALK OF $D = \{ \vec{\phi}^{(1)}, t^{(1)}, \dots, \vec{\phi}^{(N)}, t^{(N)} \}$

IN OTHER WORDS

$$E_D [y(\vec{x})] = E_D [y(\vec{\phi})] = \vec{\phi}^T \underbrace{E_D \left[\frac{1}{N} \sum_{n=1}^N t^{(n)} \vec{\phi}^{(n)} \right]}_{\text{CORRELATION OF } t, \vec{\phi}}$$

IN GENERAL WE CAN TALK OF

$$\begin{aligned} E_D [y(\vec{x})] &= \int \dots \int y(\vec{x}) p(\vec{x}^{(1)}, t^{(1)}, \dots, \vec{x}^{(N)}, t^{(N)}) d\vec{x}^{(1)} \dots dt^{(N)} \\ &= \int \dots \int y(\vec{x}) \prod_{n=1}^N p(\vec{x}^{(n)}, t^{(n)}) d\vec{x}^{(1)} \dots dt^{(N)} \end{aligned}$$

$$\begin{aligned} E_{D, \vec{x}} [y(\vec{x})] &= \int \dots \int y(\vec{x}) p(\vec{x}, \underbrace{\vec{x}^{(1)}, t^{(1)}, \dots, \vec{x}^{(N)}, t^{(N)}}_{\text{TRAINING TOKENS}}) d\vec{x} \dots dt^{(N)} \\ &= \int \dots \int y(\vec{x}) p(\vec{x}) \prod_{n=1}^N p(\vec{x}^{(n)}, t^{(n)}) d\vec{x} \dots dt^{(N)} \end{aligned}$$

TEST TAKEN

$$E_{D, \vec{x}, t} [|y(\vec{x}) - t|^2] = \int \dots \int |y(\vec{x}) - t|^2 p(\vec{x}, t) \prod_{n=1}^N p(\vec{x}^{(n)}, t^{(n)}) d\vec{x} \dots dt^{(N)}$$

TEST TAKEN TRAINING TOKENS

$$= \underbrace{E_{D, \vec{x}} [|y(\vec{x}) - E_D [y(\vec{x})]|^2]}_{\text{NETWORK VARIANCE}} + \underbrace{E_{\vec{x}} [|E_D [y(\vec{x})] - \langle t | \vec{x} \rangle|^2]}_{\text{NETWORK BIAS}} + \underbrace{E_{\vec{x}, t} [|\langle t | \vec{x} \rangle - t|^2]}_{\text{RESIDUAL (CONDITIONAL) TARGET VARIANCE}}$$

+ CROSS-TERMS THAT GO TO ZERO

VARIANCE

EXAMPLE

$$y(\vec{x}) = \sum_{j=1}^M w_j \phi_j(\vec{x}) = \vec{w}^T \vec{\Phi}$$

$$\vec{w} = \frac{1}{N} \sum_{n=1}^N t^{(n)} \vec{\Phi}^{(n)}$$

SUPPOSE $\vec{\Phi}^{(n)} \sim \mathcal{N}(0, \mathbf{I})$

$$\text{AND } P_T(t^{(n)}) = \begin{cases} \frac{1}{2} & t^{(n)} = 1 \\ \frac{1}{2} & t^{(n)} = -1 \end{cases}$$

THEN

$$t^{(n)} \vec{\Phi}^{(n)} \sim \mathcal{N}(0, \mathbf{I})$$

$$\text{AND } \text{Var}(\vec{w}) = \frac{1}{N^2} \sum_{n=1}^N \text{Var}(t^{(n)} \vec{\Phi}^{(n)}) = \frac{1}{N} \mathbf{I}$$

$$\text{SO } \boxed{\vec{w} \sim \mathcal{N}(0, \frac{1}{N} \mathbf{I})}$$

$$\begin{aligned} & \text{Var}(y(\vec{x}) - \mathbb{E}_D[y(\vec{x})]) \\ &= \text{Var}\left(\sum_{j=1}^M w_j \phi_j(\vec{x}) - \mathbb{E}_D\left(\sum_{j=1}^M w_j \phi_j(\vec{x})\right)\right) \end{aligned}$$

$$\begin{aligned} &= \sum_{j=1}^M \phi_j(\vec{x}) \underbrace{(w_j - \mathbb{E}_D[w_j])}_{\sim \mathcal{N}(0, \frac{1}{N})} \sim \mathcal{N}\left(0, \frac{\sum_{j=1}^M \phi_j^2(\vec{x})}{N}\right) \end{aligned}$$

$$\approx \mathcal{N}\left(0, \frac{M}{N}\right)$$

$$\text{SO } \mathbb{E}_{D, \vec{x}} \left[(y(\vec{x}) - \mathbb{E}_D[y(\vec{x})])^2 \right] \approx \frac{M}{N}$$

MORE HIDDEN
NODES
= WORSE VARIANCE

MORE DATA
= LOWER
VARIANCE

BIAS

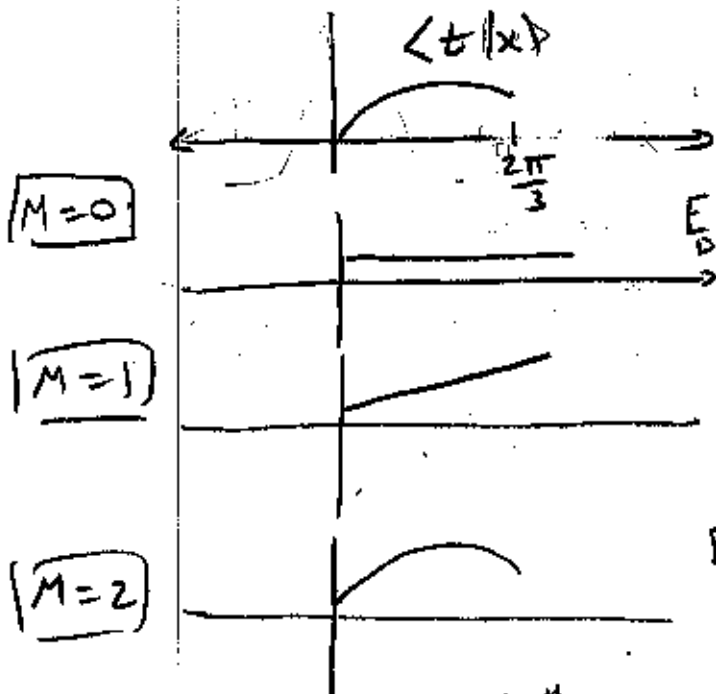
$$E_x \left[\left(E_D [y(x)] - \langle t|x \rangle \right)^2 \right]$$

EXAMPLE

$$x = \text{SCALAR} \sim U\left(-\frac{\pi}{2}, \frac{\pi}{2}\right)$$

$$\phi_m(x) = x^m, \quad y(x) = \sum_{m=0}^M w_m \phi_m(x)$$

$$\langle t|x \rangle = \sin(x)$$



$M=0$

$$E_D [y(x)] = E_D [w_0]$$

$M=1$

$$E_D [y(x)] = E_D [w_0 + w_1 x]$$

$M=2$

$$E_D [y(x)] = E_D [w_0 + w_1 x + w_2 x^2]$$

$$E_D \left[\sum_{m=0}^M w_m x^m \right] = \operatorname{argmin}_x E_x \left[\left(\sin(x) - \sum_{m=0}^M w_m x^m \right)^2 \right]$$

$$\Rightarrow \vec{w} = R^{-1} \vec{r}$$

$$\min_x E_x \left[\left(\sin(x) - \sum_{m=0}^M w_m x^m \right)^2 \right] = \delta - \vec{r}^T R^{-1} \vec{r}$$

$$\text{FOR: } \delta = E[\sin^2(x)]$$

$$R(m, n) = E[x^{m+n}]$$

$$\vec{r}(m) = E[x^m \sin(x)]$$

MONOTONICALLY NON-INCREASING
FUNCTION OF M ,
INDEPENDENT OF N !

so

$$MSE = E_{\theta, \beta, t} [|y(x) - t|]^2$$

$$= E_{\beta} [\text{BIAS}^2(x)] + \text{VARIANCE} + \text{RESIDUAL}$$

