

Foundations of Machine Learning AI2000 and AI5000

FoML-10

Bias Variance Decomposition

Dr. Konda Reddy Mopuri

Department of AI, IIT Hyderabad

July-Nov 2025



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



So far in FoML

- What is ML and the learning paradigms
- Probability refresher
- MLE, MAP, and fully Bayesian treatment
- Linear Regression with basis functions - and regularization
- Model selection



Breaking down the prediction error of a model



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



Frequentist interpretation of the model complexity



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



Expected Loss for Regression

- Regression loss $L(t, y(\mathbf{x})) = [t - y(\mathbf{x})]^2$
for a given $(\mathbf{x}, t) \sim P(\mathbf{x}, t)$



Expected Loss for Regression

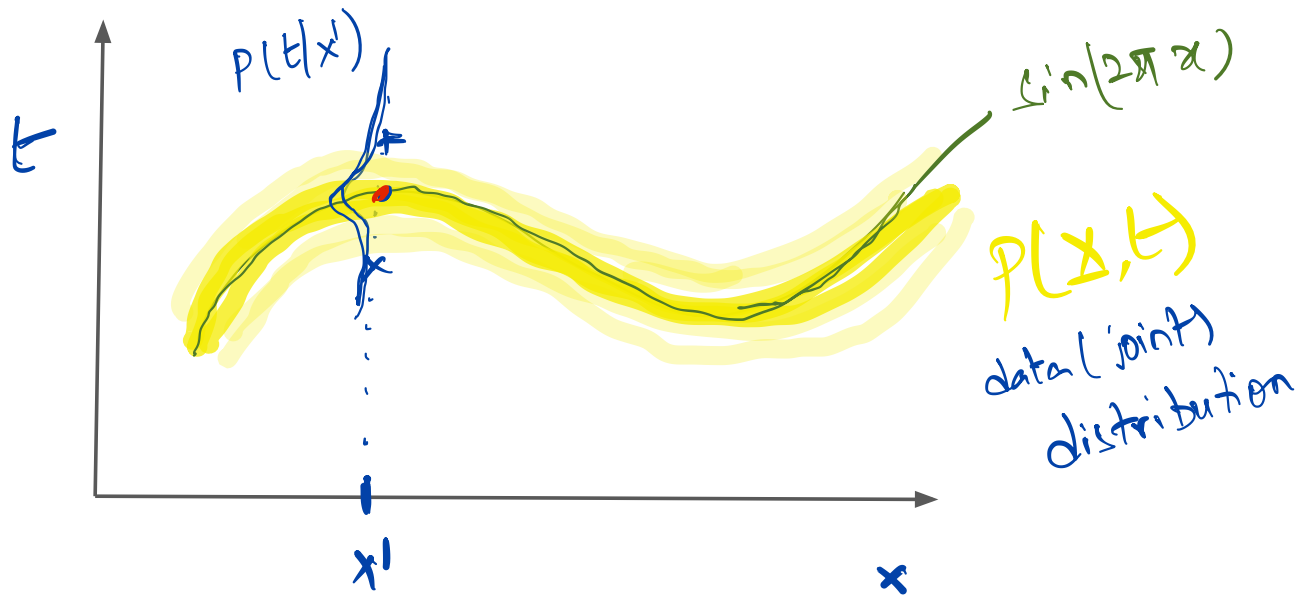
- Regression loss $L(t, y(\mathbf{x})) = [t - y(\mathbf{x})]^2$
for a given $(\mathbf{x}, t) \sim P(\mathbf{x}, t)$

- If we know the data distribution, we can find the

$$\mathbb{E}[L(t, y(\mathbf{x}))] = \iint [t - y(\mathbf{x})]^2 P(\mathbf{x}, t) d\mathbf{x} dt$$



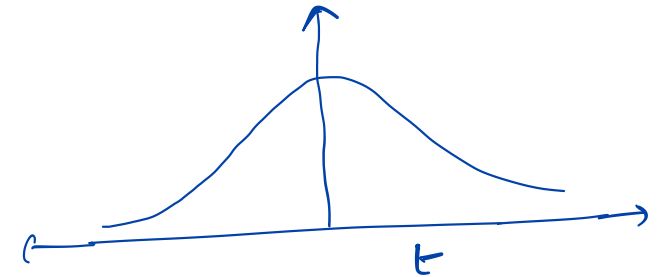
Data and prediction distributions



$$t = \sin 2\pi x + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \bar{\sigma}^2)$$

$$p(t|x)$$



Minimizing the Expected loss at given x

$$L = \int [t - y(x)]^2 p(t|x) dt$$
$$\frac{\partial L}{\partial y(x)} = 2 \int [t - y(x)] p(t|x) dt = 0$$

$$\int t p(t|x) dt = \int y(x) p(t|x) dt$$

$$E[t|x] = y(x)$$

Regression function



Expected Loss for Regression

$$\mathbb{E}[L] = \int \int (y(\mathbf{x}) - t)^2 p(\mathbf{x}, t) dt d\mathbf{x}$$

$$= \mathbb{E}_{\mathbf{x}, t} [(y(\mathbf{x}) - t)^2] = \mathbb{E}_{\mathbf{x}, t} [(y(\mathbf{x}) - \mathbb{E}[\varphi(\mathbf{x})] + \mathbb{E}[\varphi(\mathbf{x})] - t)^2]$$

$$= \underbrace{\mathbb{E}_{\mathbf{x}, t} [(y(\mathbf{x}) - \mathbb{E}[\varphi(\mathbf{x})])^2]}_{\text{Due to the model}} + \underbrace{\mathbb{E}_{\mathbf{x}, t} [(\mathbb{E}[\varphi(\mathbf{x})] - t)^2]}_{\text{Due to the intrinsic noise}}$$

Cross terms
vanish

$$\mathbb{E}_{\mathbf{x}, t} [(y(\mathbf{x}) - \mathbb{E}[\varphi(\mathbf{x})]) \cdot (\mathbb{E}[\varphi(\mathbf{x})] - t)]$$

$$= \mathbb{E}_{\mathbf{x}} [\underbrace{\mathbb{E}_t [(y(\mathbf{x}) - \mathbb{E}[\varphi(\mathbf{x})]) \cdot (\mathbb{E}[\varphi(\mathbf{x})] - t) | \mathbf{x}]}_0]$$



Minimizing the expected loss

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - \mathbb{E}[t/\mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \int \text{var}[t/\mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

- Optimal solution is unknown $y(\mathbf{x}) = \mathbb{E}[t/\mathbf{x}]$

if we model $\mathbb{E}[t/\mathbf{x}]$ using parameters $\underline{\omega}$, then from a Bayesian perspective, we can express the model's uncertainty via a posterior on $\underline{\omega}$

But we make point estimate for $\underline{\omega}$ on a dataset D



Minimizing the expected loss

$$\mathbb{E}[L] = \int \underbrace{\{y(\mathbf{x}) - \mathbb{E}[t/\mathbf{x}]\}^2}_{\text{var}[t/\mathbf{x}]} p(\mathbf{x}) d\mathbf{x} + \underbrace{\int \text{var}[t/\mathbf{x}] p(\mathbf{x}) d\mathbf{x}}$$

- Optimal solution is unknown $y(\mathbf{x}) = \mathbb{E}[t/\mathbf{x}]$
- We only have finite dataset (but not the distribution)

$$D = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_n, t_n)\}$$



Minimizing the expected loss

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - \mathbb{E}[t/\mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \int \text{var}[t/\mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

- Frequentist approach → multiple datasets, multiple models

$$D_1 = \{ \dots \} \quad D_2 = \{ \dots \} \quad \dots \quad D_L = \{ \dots \}$$

$\underbrace{y_1}$ $\underbrace{y_2}$ $\underbrace{y_L}$

$$\mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}[t/\mathbf{x}])^2]$$

Estimate the performance by averaging
the expected loss over different
datasets



Minimizing the expected loss

$$\mathbb{E}[\mathbb{E}_D[L]] = \int \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}[t/\mathbf{x}])^2]p(\mathbf{x})d\mathbf{x} + \int \text{var}[t/\mathbf{x}]p(\mathbf{x})d\mathbf{x}$$

- Bias-Variance decomposition



Minimizing the expected loss

$$\mathbb{E}[\mathbb{E}_D[L]] = \int \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}[t/\mathbf{x}])^2] p(\mathbf{x}) d\mathbf{x} + \int \text{var}[t/\mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

- Bias-Variance decomposition

$$\mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}[t/\mathbf{x}])^2] = \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}_D[y_D(\mathbf{x})] + \mathbb{E}_D[y_D(\mathbf{x})] - \mathbb{E}[t/\mathbf{x}])^2]$$

$$(\text{Bias})^2 = \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}_D[y_D(\mathbf{x})])^2] + \mathbb{E}_D[(\mathbb{E}_D[y_D(\mathbf{x})] - \mathbb{E}[t/\mathbf{x}])^2]$$

Variance
(Bias)²

Variance =

$$\text{Noise} = \int \text{var}[t/\mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$



Example



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL

Data-driven Intelligence
& Learning Lab

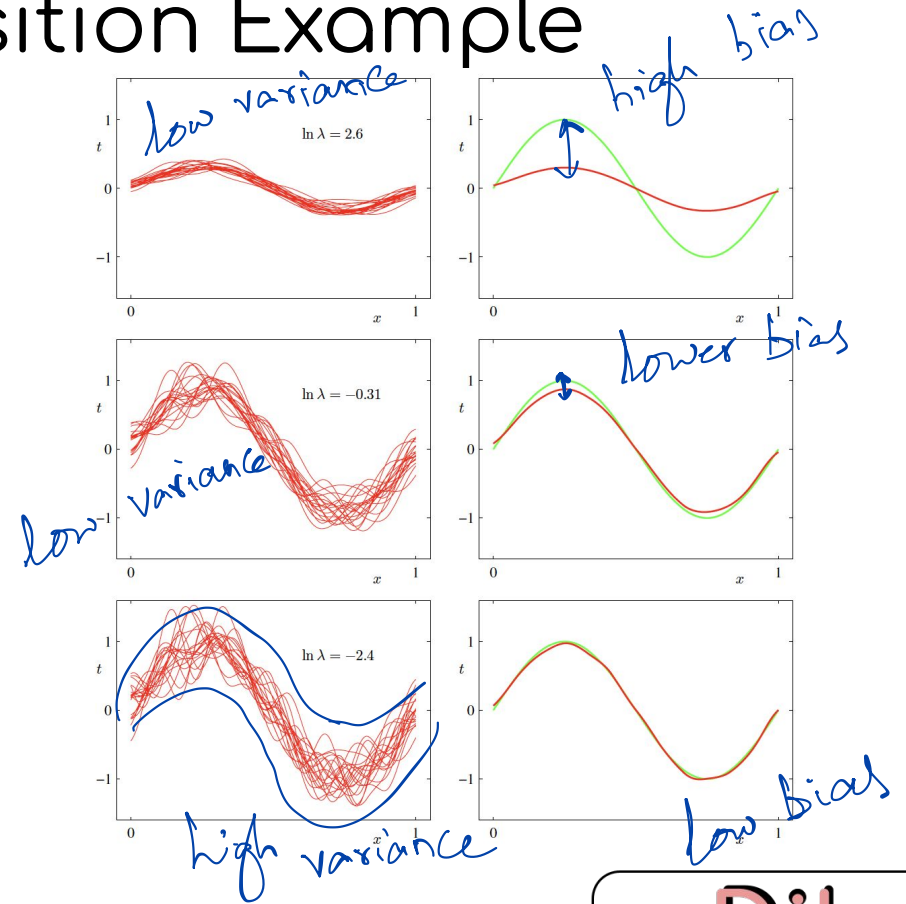
Bias-Variance Decomposition Example

- 100 datasets of size 25
- $x \sim U[0, 1]$
- $t = \sin(2\pi x) + \epsilon$

$$y^{(i)} = \omega^{(i)T} \phi(x)$$

$d = 1 \text{ to } 100$

$$\mathbb{E}_D[y_D(x)] = \bar{y}(x)$$



Bias-Variance Decomposition Example

(quantifying)

Estimating the bias and variance [since we know the ground truth]

$$(\text{bias})^2 = \int \{ \mathbb{E}_D[y_D(x) - \mathbb{E}[t/x] \}^2 p(x) dx$$

$$= \frac{1}{N} \sum_{i=1}^N \left[\bar{y}(x_i) - \mathbb{E}(t/x_i) \right]^2$$

$\frac{1}{L} \sum_{l=1}^L y^{(l)}$

Numerical
integration
over

$\{x_1, x_2, \dots, x_N\}$ ← Monte Carlo approximation

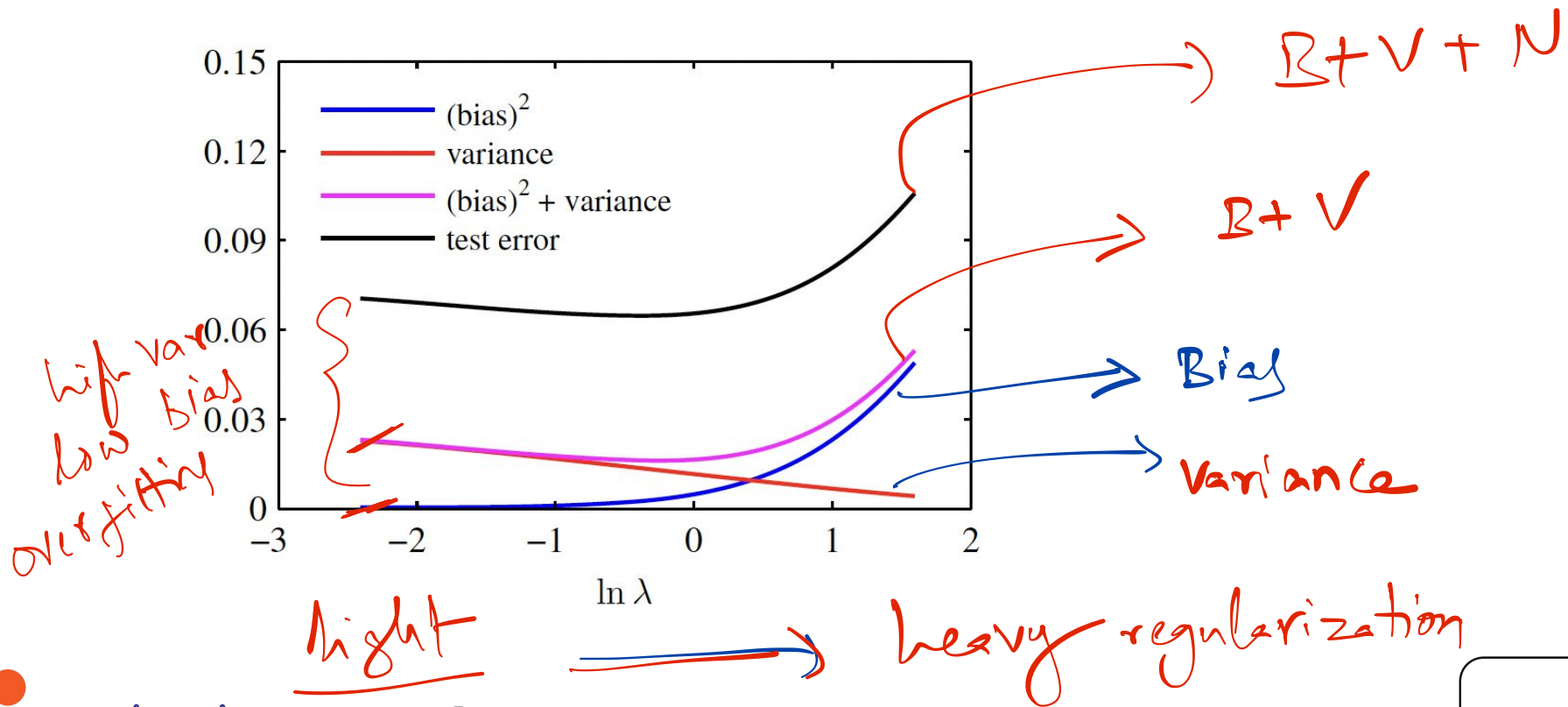
$$\text{variance} = \mathbb{E}_D[\{y_D(x) - \mathbb{E}_D[y_D(x)]\}^2 p(x) dx]$$

$$= \frac{1}{N} \sum_{i=1}^N \frac{1}{L} \sum_{l=1}^L [y^{(l)}(x_i) - \bar{y}(x_i)]^2$$

$$\begin{aligned} \bar{y}(x) &= \mathbb{E}_D[y_D(x)] \\ &= \frac{1}{L} \sum_{l=1}^L y^{(l)}(x) \end{aligned}$$



Bias-Variance Decomposition Example



Bias-Variance Decomposition

- In practice - we don't split our dataset to determine the model complexity
 - Large datasets are better
- Bayesian regression!



Rough work



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



Next Bayesian Regression



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

