

Parametric regression (cont'd) and nonparametric methods

STATS 303 Statistical Machine Learning

Spring 2022

Lecture 6

multivariate regression

$$y = f(x) + \varepsilon, \quad x \in \mathbb{R}^D, \quad f: \mathbb{R}^D \rightarrow \mathbb{R}$$

Estimate the "true model" $f(x)$ with $\underbrace{g(x|w)}$

$$\mathcal{X} = \{\mathcal{X}_n\}_{n=1}^N = \{(x_n, y_n)\}_{n=1}^N$$

assume to be linear:
 $g(x|w) = w_0 + w_1 x_1 + \dots + w_D x_D$

Following the assumption that $\varepsilon \sim \mathcal{N}(0, \sigma^2)$, similarly to the 1D case, we will have to minimize an error function

$$E(w_0, w_1, \dots, w_D | \mathcal{X})$$

$$= \frac{1}{2} \sum_{n=1}^N \left(y_n - \underbrace{(w_0 + w_1 x_{n1} + w_2 x_{n2} + \dots + w_D x_{nD})}_{\text{n-th entry of } Xw} \right)^2$$

Write $X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1D} \\ 1 & x_{21} & x_{22} & \dots & x_{2D} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & x_{N2} & \dots & x_{ND} \end{bmatrix} \in \mathbb{R}^{N \times (D+1)}$

data matrix

$$w = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{pmatrix} \in \mathbb{R}^{D+1}$$

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \in \mathbb{R}^N$$

$$E(w | \mathcal{X}) = \frac{1}{2} \|y - Xw\|^2$$

$$\text{Setting } \frac{\partial \epsilon}{\partial w} = \frac{1}{2} \cdot 2 \cdot X^T (Xw - y) = \underbrace{X^T (Xw - y)}_{X^T X w = X^T y} = 0$$

$$\text{yields } w = (X^T X)^{-1} X^T y =: \hat{w}_{ls}$$

Ridge regression

$$\hat{w}_{\text{map}} = \underset{w}{\operatorname{argmax}} \log p(w|x)$$

$$= \underset{w}{\operatorname{argmax}} \log p(x|w) + \log p(w)$$

$$\text{Assume } p(w) \sim \mathcal{N}(0, \tilde{\sigma}^2 I)$$

$$\text{Then } p(w) = \prod_{j=0}^D N(w_j | 0, \tilde{\sigma}^2)$$

$$\propto \prod_{j=0}^D \exp\left(-\frac{w_j^2}{2\tilde{\sigma}^2}\right)$$

$$\text{That implies } \log p(w) = -\sum_{j=0}^D \frac{w_j^2}{2\tilde{\sigma}^2} + \text{const.}$$

$$\hat{w}_{\text{map}} = \underset{w}{\operatorname{argmin}} \underbrace{\frac{1}{2} \|y - Xw\|^2 + \frac{\lambda}{2} \|w\|^2}_{E_{\text{ridge}}}$$

"ridge regression"

E_{ridge}

Setting $\frac{\partial E_{\text{ridge}}}{\partial w} = X^T(Xw - y) + \lambda w = 0$

yields $X^T X w + \lambda I w = X^T y$

That is $(X^T X + \lambda I) w = X^T y$

which gives $w = \underline{(X^T X + \lambda I)^{-1} X^T y} =: \hat{w}_{\text{ridge}}$

Next, let's compare \hat{w}_{ridge} with \hat{w}_{ls}

Apply SVD to X , say $X = U \Delta V^T$,

where $\Delta = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_D)$ and $U^T U = V^T V = I$

Then $X^T = V \Delta U^T$ and

$$X^T X = V \Delta \underbrace{U^T U}_I \Delta V^T = V \Delta^2 V^T$$

That implies

$$X \hat{w}_{\text{ridge}} = U \Delta V^T (V \Delta^2 V^T + \lambda I)^{-1} V \Delta U^T y$$

$$= U \Delta V^T (V (\Delta^2 + \lambda I) V^T)^{-1} V \Delta U^T y$$

$$= U \Delta \underbrace{V^T}_I V (\Delta^2 + \lambda I)^{-1} \underbrace{V^T}_I V \Delta U^T y$$

$$= U \underbrace{\Delta (\Delta^2 + \lambda I)^{-1} \Delta}_{\text{diagonal}} U^T y = U \Delta^2 (\Delta^2 + \lambda I)^{-1} U^T y$$

$$= \left[\sum_{j=0}^D u_j \frac{\sigma_j^2}{\sigma_j^2 + \lambda} u_j^T y \right] X \hat{w}_{\text{ridge}}$$

On the other hand,

$$X \hat{w}_{\text{ls}} = X (X^T X)^{-1} X^T y$$

$$= U \Delta \underbrace{V^T V}_I \Delta^{-2} \underbrace{V^T V}_I \Delta U^T y$$

$$= U \underbrace{\Delta \Delta^{-2} \Delta}_I U^T y$$

$$= U U^T y = \left[\sum_{j=0}^D u_j u_j^T y \right] X \hat{w}_{\text{ls}}$$

" Ridge regression gives a shrunken version of estimation compared with regular linear regression. "

LASSO regression

$$\hat{w}_{\text{LASSO}} = \underset{w}{\operatorname{argmin}} \quad \frac{1}{2} \|y - Xw\|^2 + \lambda \|w\|_1$$

$$\text{where } \|w\|_1 = \sum_{j=0}^D |w_j|$$

$$\left[\begin{array}{l} \text{MAP estimator with a different prior} \\ p(w) \propto \exp\left(-\alpha \sum_{j=0}^D |w_j|\right) \end{array} \right. \text{Laplacian prior} \left. \right]$$

"No closed-form solution."

bias-variance decomposition

decomposition of expected loss

true
output
↓
our
model
↓

- In general, suppose we have a loss function $L(t, y(\mathbf{x}))$. The overall average expected loss is

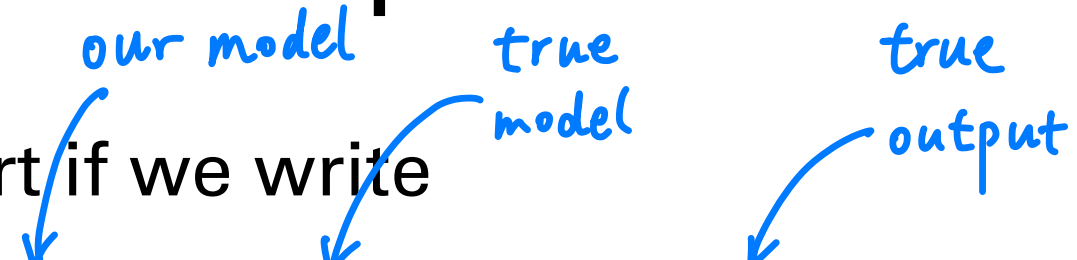
$$\mathbb{E}[L] = \iint L(t, y(\mathbf{x})) p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$

- In regression, the loss is given by

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

decomposition of expected loss

- It does not hurt if we write



The diagram shows three handwritten blue annotations with arrows pointing to terms in the equation: 'our model' points to $y(\mathbf{x})$, 'true model' points to $\mathbb{E}[t|\mathbf{x}]$, and 'true output' points to t .

$$\begin{aligned}\{y(\mathbf{x}) - t\}^2 &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 \\ &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 + 2\{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}\{\mathbb{E}[t|\mathbf{x}] - t\} + \{\mathbb{E}[t|\mathbf{x}] - t\}^2\end{aligned}$$

where the expectation is taken over t .

decomposition of expected loss

$$\begin{aligned} \{y(\mathbf{x}) - t\}^2 &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 \\ &= \underbrace{\{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2}_{\textcircled{1}} + 2\underbrace{\{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}\{\mathbb{E}[t|\mathbf{x}] - t\}}_{\textcircled{2}} + \underbrace{\{\mathbb{E}[t|\mathbf{x}] - t\}^2}_{\textcircled{3}} \end{aligned}$$

- Let's calculate

$$\begin{aligned} \mathbb{E}[L] &= \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt \\ &= \iint \textcircled{1} \, p(\mathbf{x}, t) \, d\mathbf{x} \, dt + \iint \textcircled{2} \, p(\mathbf{x}, t) \, d\mathbf{x} \, dt + \\ &\quad \iint \textcircled{3} \, p(\mathbf{x}, t) \, d\mathbf{x} \, dt \end{aligned}$$

decomposition of expected loss

$$\begin{aligned}\{y(\mathbf{x}) - t\}^2 &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 \\ &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 + 2\{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}\{\mathbb{E}[t|\mathbf{x}] - t\} + \{\mathbb{E}[t|\mathbf{x}] - t\}^2\end{aligned}$$

- Let's calculate

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

First,

$$\begin{aligned}& \iint (y(x) - \mathbb{E}[t|x])^2 p(x, t) dx dt \\ &= \int \underbrace{\left(\int p(x, t) dt \right)}_{p(x)} (y(x) - \mathbb{E}[t|x])^2 dx \\ &= \int (y(x) - \mathbb{E}[t(x)])^2 p(x) dx\end{aligned}$$

decomposition of expected loss

$$\begin{aligned}\{y(\mathbf{x}) - t\}^2 &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 \\ &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 + 2\{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}\{\mathbb{E}[t|\mathbf{x}] - t\} + \{\mathbb{E}[t|\mathbf{x}] - t\}^2\end{aligned}$$

- Let's calculate

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

Second, $2 \iint (y(x) - \mathbb{E}[t|x])(\mathbb{E}[t|x] - t) p(x, t) dx dt$

$$= 2 \iint (y(x) - \mathbb{E}[t|x])(\mathbb{E}[t|x] - t) p(t|x) p(x) dx dt$$

$$= 2 \int (y(x) - \mathbb{E}[t|x]) \left(\int (\mathbb{E}[t|x] - t) p(t|x) dt \right) p(x) dx$$

$$= 0 \quad \text{because} \quad \mathbb{E}[t|x] - \int t p(t|x) dt = \mathbb{E}[t|x] - \mathbb{E}[t|x]$$

decomposition of expected loss

$$\begin{aligned}\{y(\mathbf{x}) - t\}^2 &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 \\ &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 + 2\{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}\{\mathbb{E}[t|\mathbf{x}] - t\} + \{\mathbb{E}[t|\mathbf{x}] - t\}^2\end{aligned}$$

- Let's calculate

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

Third,

$$\iint (\underbrace{\mathbb{E}[t|\mathbf{x}]}_{\text{true model}} - \underbrace{t}_{\text{true output}})^2 p(\mathbf{x}, t) d\mathbf{x} dt = \text{"noise"}$$

③

decomposition of expected loss

- We have

$$\mathbb{E}[L] = \underbrace{\int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x}}_{\substack{\text{?} \\ \text{our} \\ \text{model} \quad \text{true} \\ \text{model}}} + \underbrace{\int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt}_{\substack{\text{noise} \quad \text{"unavoidable"}}}$$

where

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt$$

- NOTE, our $y(\mathbf{x})$ depends on the dataset (sample) \mathcal{D} !!! Let's denote it as $y(\mathbf{x}) = y(\mathbf{x}; \mathcal{D})$ to explicitly emphasize this fact.

Remark: \mathcal{D} was denoted as χ before

decomposition of expected loss

- We expand $\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2$ as

$$\begin{aligned} & \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ & \quad + 2\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}. \end{aligned}$$

decomposition of expected loss

- We expand $\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2$ as

$$\begin{aligned} & \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ & \quad + 2\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}. \end{aligned}$$

- Taking the expectation over the ensemble of datasets, we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] \\ &= \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]}_{\text{variance}} \end{aligned}$$

decomposition of expected loss

- Therefore, we have

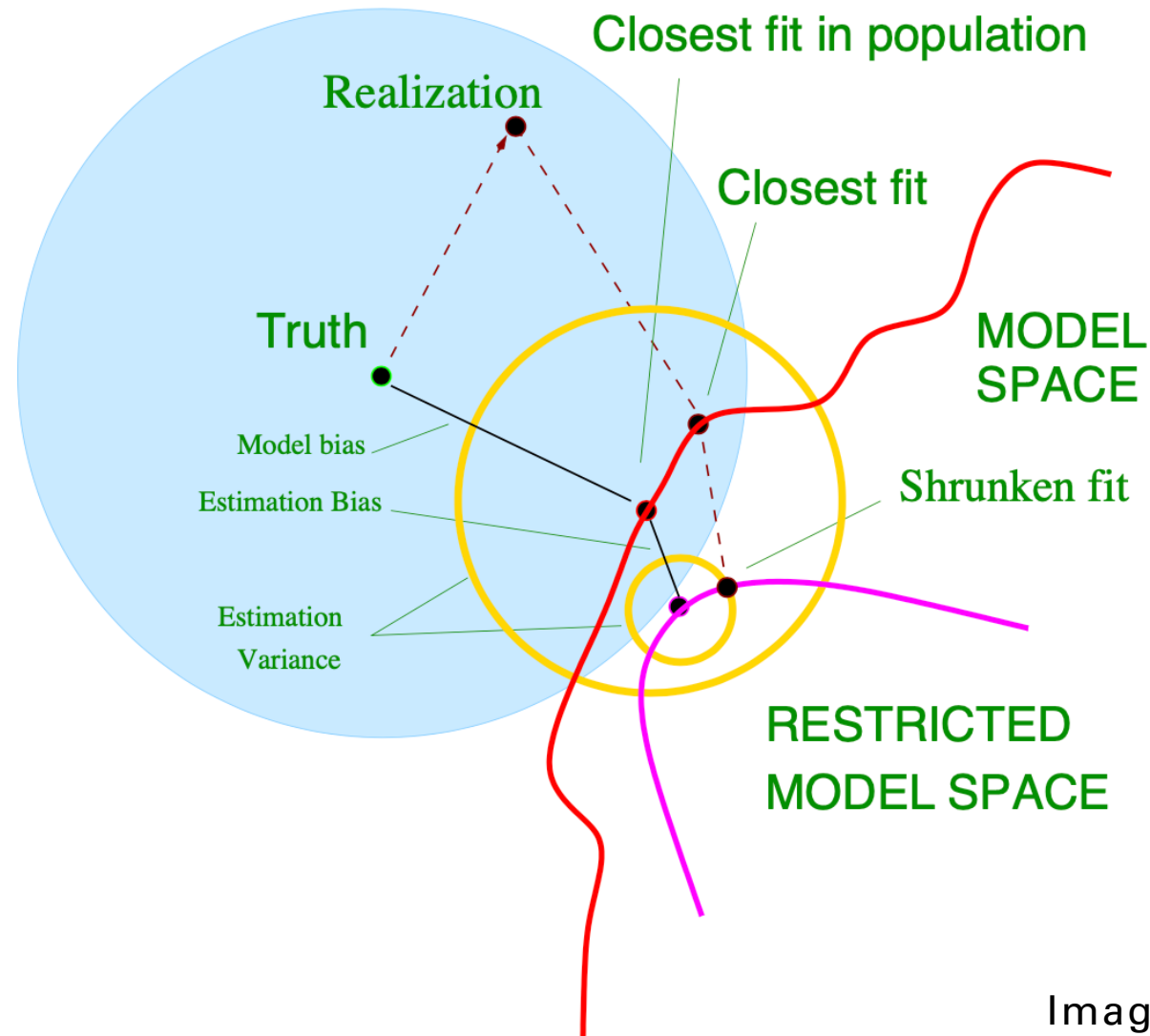
$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$

$$(\text{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) \, d\mathbf{x}$$

$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) \, d\mathbf{x}$$

$$\text{noise} = \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$

decomposition of expected loss





Questions?

Reference

- *Multivariate regression:*
 - [Al] Ch.5.8
 - [HaTF] Ch.3.2
 - [Bi] Ch.3.2 (bias-var)
- *Ridge and LASSO:*
 - [HaTF] Ch.3.4
- *Overview nonparametric methods:*
 - [Al] Ch.8.1-8.3
 - [Bi] Ch.2.5.1
- *Nonparametric classification and regression:*
 - [Al] Ch.8.4-8.6, 8.8
 - [HaTF] Ch.6.6, 6.1.1-6.1.2
 - [Bi] Ch.2.5.2, 6.3.1