

Discussion 1

Review

Least squares & nearest neighbours

RSS & MSE & EPE

Variable Types and Terminology

Input: a variable X . If X is a vector, its j -th element is X_j

an observation x_i
(scalar or vector)

Typically, we use i to denote the index of **observations**, while use j to denote the index of **variables**.

Model

X **Scalar**

N observations

x_1
\vdots
x_i
\vdots
x_N

$\mathbf{x} \in \mathbb{R}^N$

$X_1 \quad \dots \quad X_j \quad \dots \quad X_p$ **Vector**

N observations

x_1^T
\vdots
x_i^T
\vdots
x_N^T

$\mathbf{X} \in \mathbb{R}^{N \times p}$

p variables

Simple Approach 1: Least Squares

- Training procedure:
Method of *least-squares*
- $N = \text{\#observations}$
- Minimize the *residual sum of squares*

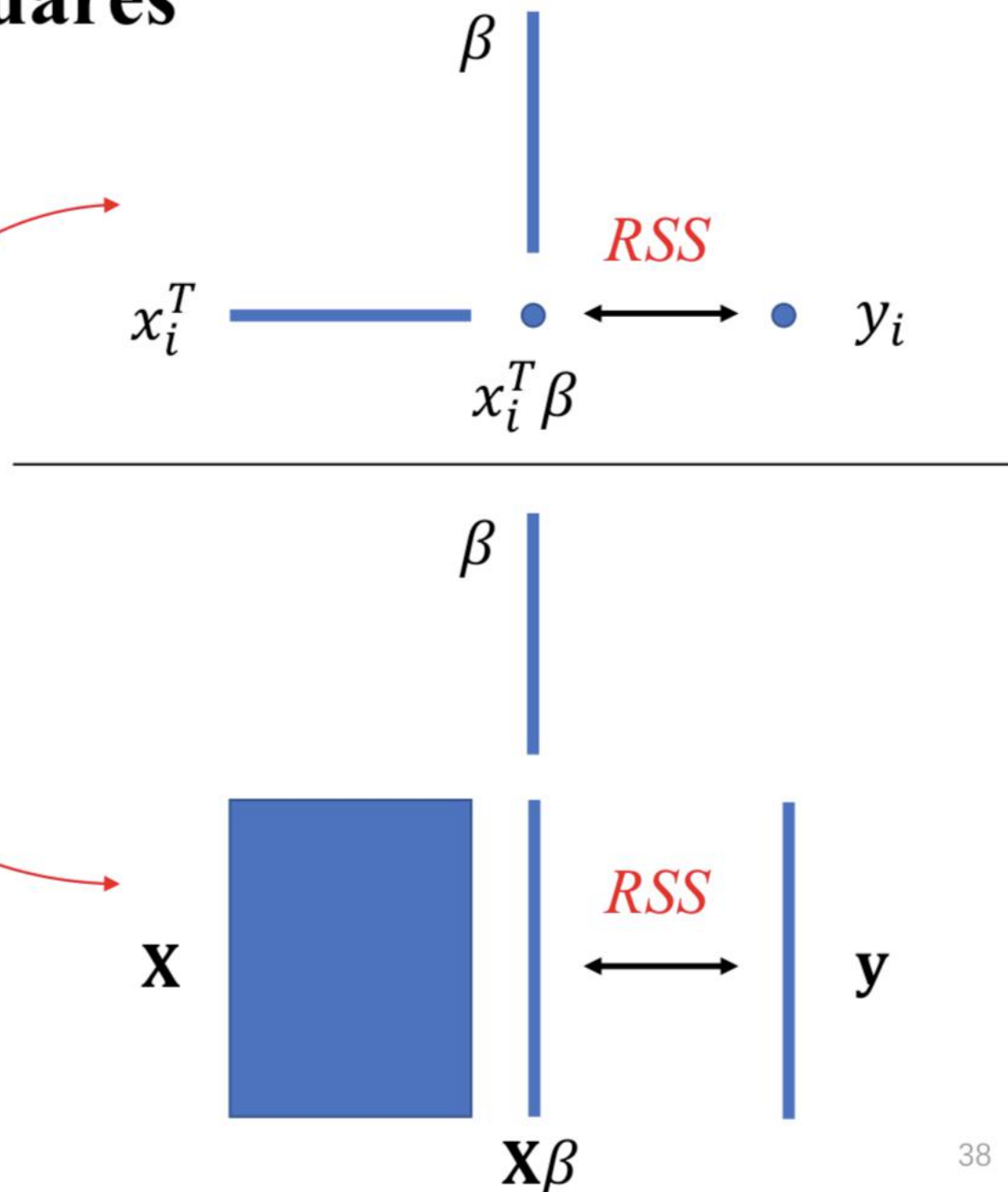
$$\text{RSS}(\beta) = \sum_{i=1}^N (y_i - x_i^T \beta)^2$$

Or equivalently,

$$\begin{aligned}\text{RSS}(\beta) &= (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \\ &= \|\mathbf{y} - \mathbf{X}\beta\|_2^2\end{aligned}$$

- This quadratic function always has a global minimum, but it may not be unique.

Q: What is the difference among x_i , x_i^T , \mathbf{x} , \mathbf{X} and \mathbf{X} ?



- (scalar to scalar) $df = f'(x)dx$
- (scalar to vector) $df = \sum_i \frac{\partial f}{\partial x_i} dx_i = \frac{\partial f^T}{\partial \mathbf{x}} d\mathbf{x}$
- $\frac{\partial A\mathbf{x}}{\partial \mathbf{x}} = A^T, \frac{\partial \mathbf{x}^T A}{\partial \mathbf{x}} = A$

$$RSS(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$$

$$= \mathbf{y}^T \mathbf{y} - 2\beta^T \mathbf{X}^T \mathbf{y} + \beta^T \mathbf{X}^T \mathbf{X} \beta$$

$$\frac{\partial RSS(\beta)}{\partial \beta} = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \beta = 0$$

$$\Rightarrow \mathbf{X}^T \mathbf{X} \beta = \mathbf{X}^T \mathbf{y}$$

$$\Rightarrow \hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{RSS}(\beta) = (\mathbf{X}\beta - \mathbf{y})^T \mathbf{W}(\mathbf{X}\beta - \mathbf{y})$$

$$\begin{aligned} \nabla_{\beta} \text{RSS}(\beta) &= \frac{\partial \text{RSS}(\beta)}{\partial \beta} \\ &= \frac{\partial}{\partial \beta} (\mathbf{X}\beta - \mathbf{y})^T \mathbf{W}(\mathbf{X}\beta - \mathbf{y}) \\ &= 2\mathbf{X}^T \mathbf{W}(\mathbf{X}\beta - \mathbf{y}) \\ &= 0 \\ \Rightarrow \hat{\beta} &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{y}. \end{aligned}$$

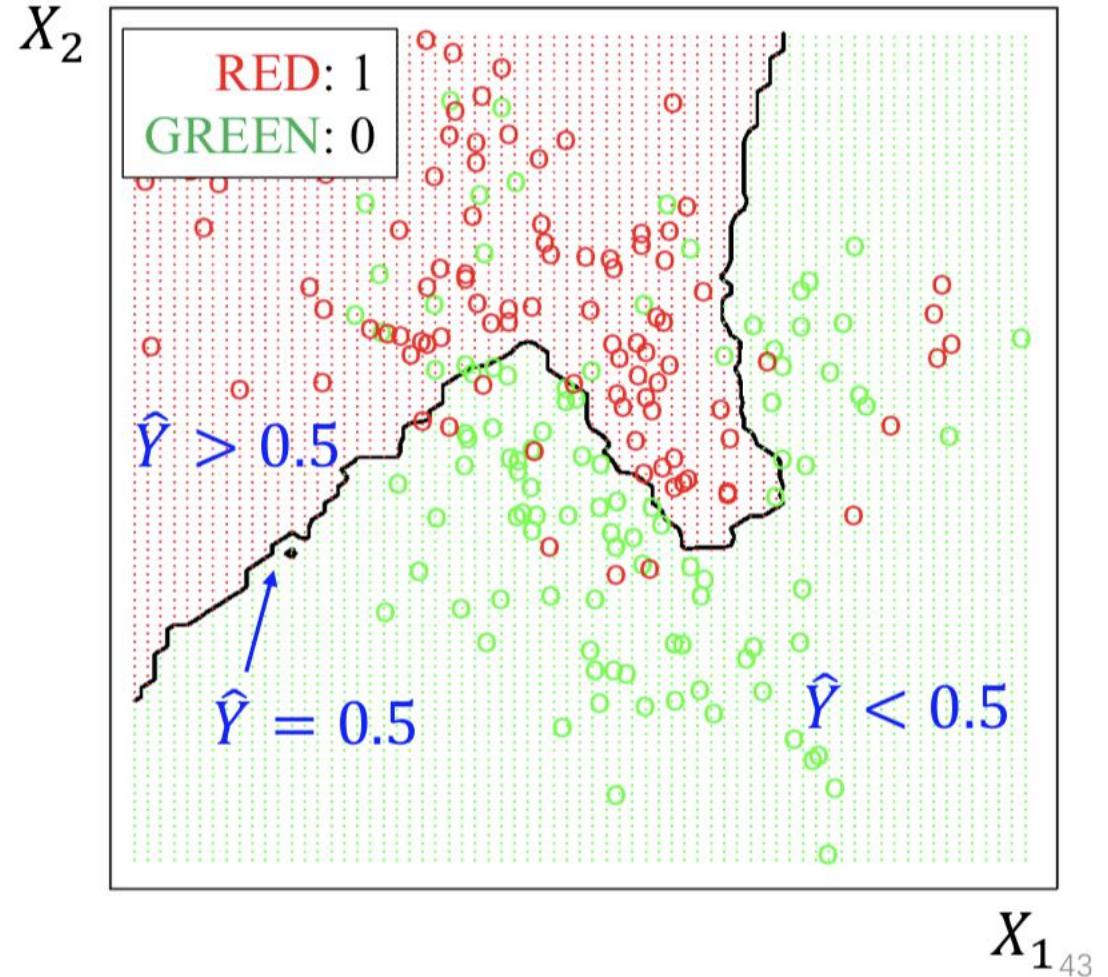
Simple Approach 2: Nearest Neighbors

- Use observations in the training set closest to the given input.

$$\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i.$$

- $N_k(x)$ is the set of the k **closest** points to x is the training sample
- **Average** the outcome of the k closest training sample points
- **Fewer misclassifications**

15-nearest neighbors averaging



RSS: training error

$$\begin{aligned}\text{RSS}(\beta) &= \sum_{i=1}^N (y_i - x_i^T \beta)^2 \\ &= (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \\ &= \|\mathbf{y} - \mathbf{X}\beta\|_2^2\end{aligned}$$

EPE: generalization error, out-of sample error

$$\begin{aligned}\text{EPE}(f) &= \mathbb{E}(Y - f(X))^2 \\ &= \int (y - f(x))^2 \Pr(dx, dy).\end{aligned}$$

$$\text{SE: } L(Y, f(X)) = (Y - f(X))^2.$$

MSE: mean squared error, in-sample error

$$\begin{aligned}\text{MSE}(x_0) &= \mathbb{E}_{\mathcal{T}} [f(x_0) - \hat{y}_0]^2 \\ &= \mathbb{E}_{\mathcal{T}} [\hat{y}_0 - \mathbb{E}_{\mathcal{T}}(\hat{y}_0)]^2 \\ &\quad + [\mathbb{E}_{\mathcal{T}}(\hat{y}_0) - f(x_0)]^2 \\ &= \text{Var}_{\mathcal{T}}(\hat{y}_0) + \text{Bias}^2(\hat{y}_0)\end{aligned}$$

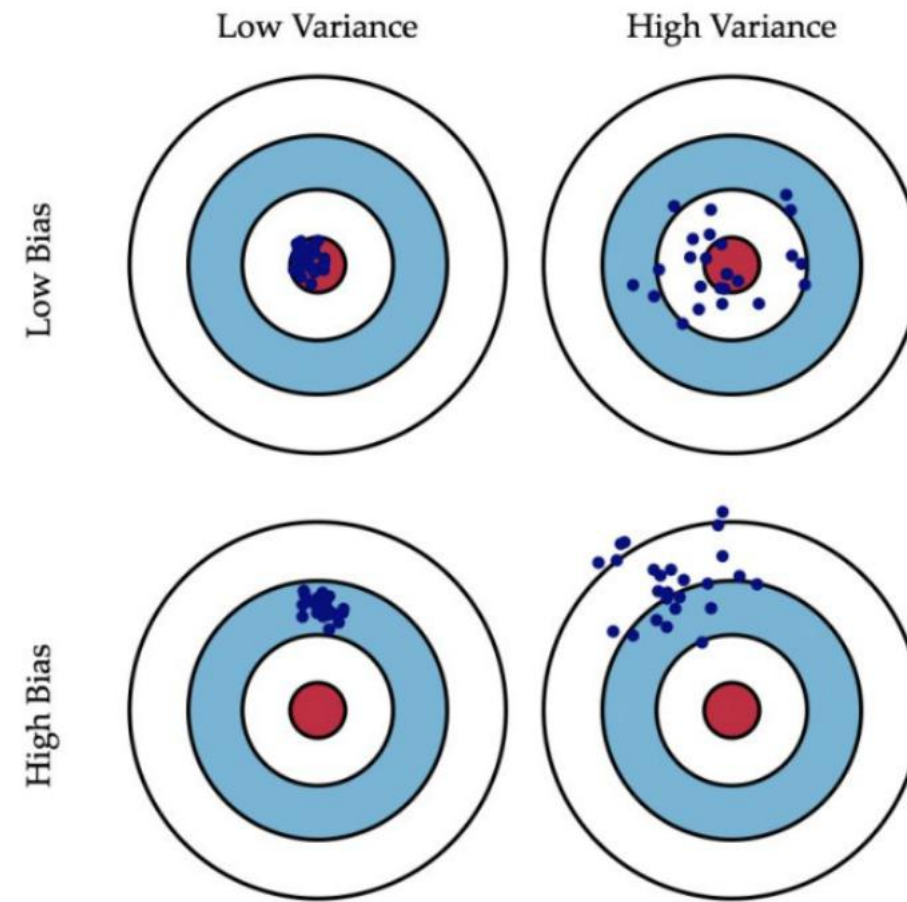


Figure 1: The difference between Bias and Var.

Regression function: $f(x) = E(Y|X = x)$.

$$\begin{aligned}\frac{\partial}{\partial f} \mathbb{E}_{Y|X}[(Y - f)^2 | X = x] &= \frac{\partial}{\partial f} \int [y - f]^2 \Pr(y|x) dy \\ &= \int \frac{\partial}{\partial f} [y - f]^2 \Pr(y|x) dy \\ &\Rightarrow 2 \int y \Pr(y|x) dy = 2f \int \Pr(y|x) dy \\ &\Rightarrow 2\mathbb{E}[Y|X = x] = 2f \\ &\Rightarrow \hat{f}(x) = \mathbb{E}[Y|X = x].\end{aligned}$$

$$\hat{G}(x) = \operatorname{argmax}_{k \in \mathcal{G}} \Pr(G = k | X = x)$$

$$\text{EPE} = \mathbb{E}_X \sum_{k=1}^K L[\mathcal{G}_k, \hat{G}(X)] \Pr(\mathcal{G}_k | X)$$

- Bayesian methods
- Formula for joint probabilities

$$\begin{aligned}\Pr(X, Y) &= \Pr(Y|X) \Pr(X) \\ &= \Pr(X|Y) \Pr(Y)\end{aligned}$$

- Bayes's theorem

$$\Pr(Y|X) = \frac{\Pr(X|Y) \Pr(Y)}{\Pr(X)}$$

Diagram illustrating Bayes's theorem with labels and arrows:

- Likelihood** points to $\Pr(X|Y)$
- Prior probability for Y** points to $\Pr(Y)$
- Posterior probability for Y** points to $\Pr(Y|X)$
- Evidence** points to $\Pr(X)$