

STA 5820

Chapter 2

Statistical Learning

Kazuhiko Shinki

Wayne State University

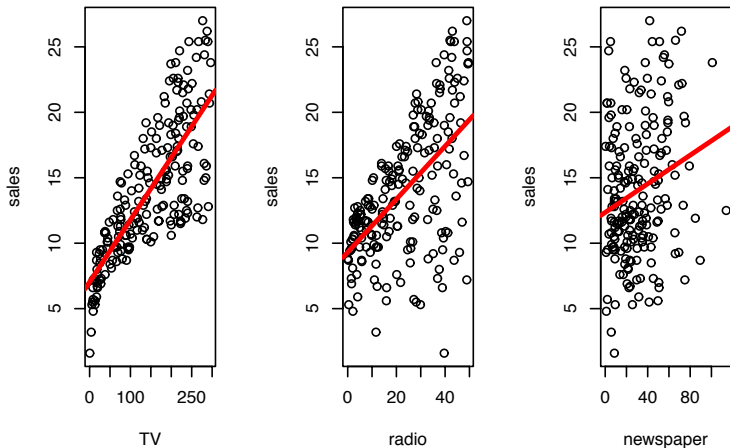
2.1 What is statistical learning?

Example:

- Objective: We want to increase sales (Y) of a particular product.
- Data: sales, advertising budgets for TV (X_1), radio (X_2) and newspaper (X_3) for 200 markets.
- Action: Which advertising budget should we increase: TV, radio or newspaper?

2.1.1 Why estimating f ?

Figure 2-1: Estimate α_i and β_i in a linear model $Y = \alpha_i + \beta_i X_i + \epsilon$ ($i = 1, 2, 3$). i with the largest β_i is the most effective advertising channel.



2.1.1 Why estimating f ?

Other models:

A multiple linear model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon \quad (1)$$

The estimated value of β_i as $\hat{\beta}_i$ ($i = 0, 1, 2, 3$) is denoted by $\hat{\beta}_i$.

More generally,

$$Y = f(X_1, X_2, X_3) + \epsilon \quad (2)$$

An estimated function of f is denoted by \hat{f} .

Conventionally $\mathbf{X} = (X_1, X_2, X_3)$ is assumed deterministic, ϵ is a random variable (so is \mathbf{Y}).

2.1.1 Why estimating f ?

Objectives:

- **Forecasting:** When a new observation $(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)$ is available but Y -value is unknown, Y is estimated by

$$\hat{Y} = \hat{f}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) \quad (3)$$

- **Inference:** In the model (1), we may want to:
 - ▶ estimate the size of β_i ($i = 1, 2, 3$). For example, we want to see how much sales increases as we spend \$1,000 additional dollars on advertisements on newspaper.
 - ▶ test if $\beta_i = 0$ ($i = 1, 2, 3$). For example, we want to see if advertisements on newspaper is statistically meaningful.

2.1.1 Why estimating f ?

A measure of error:

Let $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ be observations and $\hat{y}_i = \hat{f}(\mathbf{x}_i)$ be the predicted value for y_i ($1 \leq i \leq n$).

The **mean square error (MSE)** is defined by

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2. \quad (4)$$

The MSE is the sample version of the expected size of squared error $E[(Y - \hat{Y})^2]$.

2.1.1 Why estimating f ?

Reducible and irreducible errors:

Assume that (i) ϵ 's are independent for different observations.

Consider the situation that f is estimated by observations

$$(X_n, Y_n) = (X_{n,1}, \dots, X_{n,3}, Y_n) \quad (n = 1, \dots, N).$$

Want to evaluate the size of the error (i.e., variance of the error) when Y is forecasted by $\hat{Y} = \hat{f}(X)$ for a new observation (X, Y) .

$$\begin{aligned} E(Y - \hat{Y})^2 &= E(f(X) + \epsilon - \hat{f}(X))^2 \\ &= E(f(X) - \hat{f}(X))^2 + 2E\epsilon(f(X) - \hat{f}(X)) + E\epsilon^2 \\ &= E(f(X) - \hat{f}(X))^2 + \text{Var}\epsilon^2 \end{aligned}$$

Note that the second term in the second line vanishes because ϵ is independent of $\hat{f}(X)$ since \hat{X} is estimated by old observations.

2.1.1 Why estimating f ?

In the most right hand side of the equation,

- $E(f(X) - \hat{f}(X))^2$ is called a **reducible error** since this term gets smaller if \hat{f} is estimated accurately.
- $E\epsilon^2$ is called an **irreducible error** since this does not change regardless of our estimate of f .

In real world problems, the irreducible error exists because Y is partly determined by factors not included in X . For example, sales (Y) may also be influenced by temperature, business cycle, rival companies' products, etc.

2.1.2 How do we estimate f ?

- **Parameteric methods:** The number of parameters is finite and fixed .
 - ▶ Example: $\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$.
- **Non-parameteric Methods:** The model can not be represented by finite and fixed number of parameters.
 - ▶ Example: \hat{Y} is defined by the mean Y -value of 5 nearest observations among $(X_1, Y_1), \dots, (X_N, Y_N)$ (in terms of location of X). (**Nearest neighborhood method**).

2.1.3 The trade-off b/w Accuracy and Interpretability

More complicated model is often hard to justify by logic.

Example:

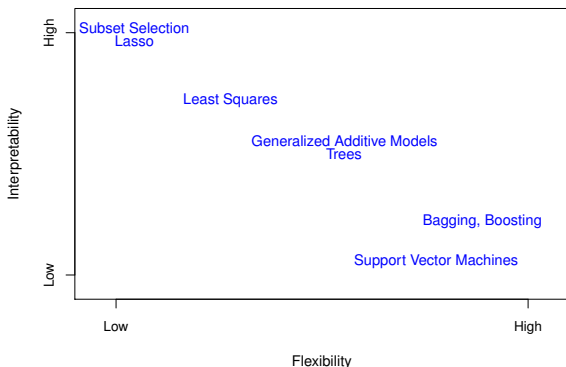
Let X be the diameter of a shell and Y be the weight of the shell.

$Y = \alpha X^3$ is justifiable.

$Y = \alpha X^{2.1}$ is hard to justify even if it is the statistically best model.

2.1.3 The trade-off b/w Accuracy and Interpretability

Figure 2-7: As the flexibility of the method increases, its interpretability decreases.



2.1.5 Regression vs Classification

As in chapter 1,

- **Regression** for quantitative response (typically real-valued).
- **Classification** for qualitative response (e.g., binary such as success/failure, and multi-class such as species { A, B, C }).

2.1.5 Regression vs Classification

In some cases, classification problem is regarded as a regression problem by considering probability.

Example: logistic regression

Want to explain gender (male or female) by some physical variables $\mathbf{X} = (X_1, \dots, X_p)$ such as weight, height, body fat percentage, body water percentage and age. When the **probability of being female** given \mathbf{X} is considered, it can be modeled by a regression model.

$$P(\text{female}|\mathbf{X}) = f(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p) \quad (5)$$

where f is called a link function which maps all real numbers to an interval $(0, 1)$.

2.1.5 Regression vs Classification

Predictors may be quantitative or qualitative (categorical), but the difference is not very important.

Example: dummy variables

Want to explain house sale price by square footage (quantitative), age (quantitative), and architecture style (qualitative; ranch, bungalow or colonial).

You can set up **dummy variables** for styles (bungalow and colonial). Then, usual regression models can be used.

| Price | Sqft | Age | Style | Bungalow | Colonial |
|--------|------|-----|----------|----------|----------|
| 325000 | 1950 | 9 | Ranch | 0 | 0 |
| 349000 | 3120 | 31 | Colonial | 0 | 1 |
| 271000 | 2070 | 35 | Colonial | 0 | 1 |
| 197000 | 1470 | 50 | Ranch | 0 | 0 |
| 149000 | 1210 | 58 | Bungalow | 1 | 0 |

2.2 Assessing Model Accuracy

For regression problems,

- Measuring the Quality of fit (2.2.1).
- Bias-Variance trade-off (2.2.2).

For classification problems, both of the measuring and trade-off issues are discussed in

- Classification setting (2.2.3).

2.2.1 Measuring the Quality of Fit

Training MSE

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(\mathbf{x}_i))^2. \quad (6)$$

is a good measure of fit. However, our ultimate objective is to minimize the expected error size for **new** observations, say (\mathbf{x}_0, y_0) , and not for **training** data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$.

(6) is called a **training MSE**.

2.2.1 Measuring the Quality of Fit

A problem on training MSE

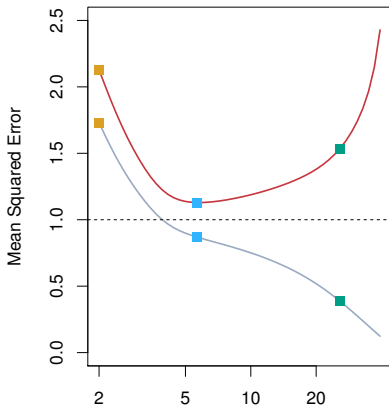
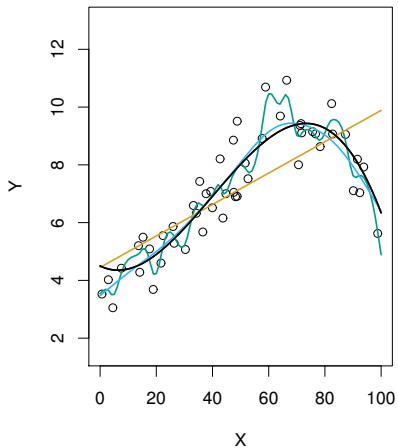
The training MSE gets (often arbitrarily) smaller as the model is more complex, but it does not mean that the expected error size for new observations gets smaller.

For example, consider a data set: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{100}, y_{100})$.

If we use a polynomial of degree 99 (so that there are 100 parameters), we can perfectly fit the data set by a curve. However, it makes poor sense.

2.2.1 Measuring the Quality of Fit

Figure 2-9: (left) Training data (circles) were generated by black curve plus noise. As a function becomes more complex (green), it fits better for training data. However, it does not mean the green curve is better than the simpler blue curve.



2.2.1 Measuring the Quality of Fit

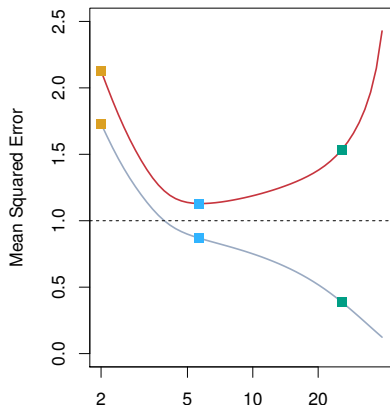
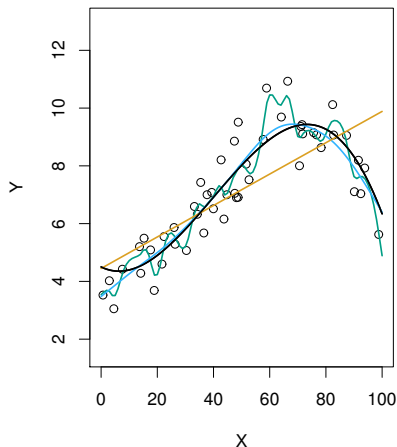
Test MSE

Consider we have (or divide a data set into) a training data set D_1 and a **test** data set D_2 . When we estimate the function f with D_1 and calculate MSE by D_2 , the MSE is called a **test** MSE.

The test MSE is a more accurate measure of the expected error.

2.2.1 Measuring the Quality of Fit

Figure 2-9: (right) The gray curve shows **training MSE** (vertical) given the flexibility (i.e., the number of parameters, horizontal). The red curve shows **test MSE**. The green curve is better in training MSE, but the blue curve is optimal in test MSE.



2.2.2 Bias-Variance trade-off

Example: Bias and Variance

Predict the height (Y) of a male by foot size (X).

Training data (foot size in cm, height in cm):

(27.2, 176), (27.5, 180), (28.0, 176), (28.7, 184), (29.0, 183).

A new observation: the foot size is 28.6cm with height unknown.

Compare two choices for prediction:

- (a) Use the data point **(28.7, 184)** only. $\hat{Y} = 184$.
 - ▶ \hat{Y} has a smaller **bias** since only the closest observation is used.
- (b) Use the mean height of 5 data points. $\hat{Y} = 179.8$.
 - ▶ \hat{Y} has a smaller **variance** since more observations are used.

2.2.2 Bias-Variance trade-off

A mathematical result

Let $\mathbf{y} = \mathbf{f}(\mathbf{x}) + \epsilon$ be the true model, and $\hat{\mathbf{f}}$ is the estimated model from data. Denote $\mathbf{f}_0 := \mathbf{f}(\mathbf{x}_0)$, $\hat{\mathbf{f}}_0 := \hat{\mathbf{f}}(\mathbf{x}_0)$.

The expected error size is decomposed into three components:

$$E(y_0 - \hat{f}_0)^2 = E(f_0 + \epsilon - \hat{f}_0)^2 \quad (7)$$

$$= E(f_0 - E\hat{f}_0 + E\hat{f}_0 - \hat{f}_0 + \epsilon)^2 \quad (8)$$

$$= (E\hat{f}_0 - f_0)^2 + E(\hat{f}_0 - E\hat{f}_0)^2 + E\epsilon^2 \quad (9)$$

$$= (\text{Bias of } \hat{f}_0)^2 + \text{Var}\hat{f}_0 + E\epsilon^2 \quad (10)$$

where the last term is an irreducible error.

Note that these three components are conceptual, and not observed in real data analysis as \mathbf{f}_0 and ϵ are unknown. The decomposition is possible for simulated examples only.

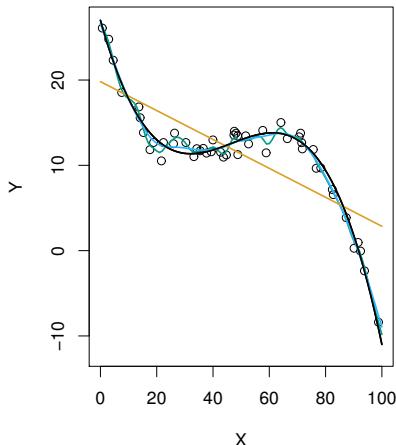
2.2.2 Bias-Variance trade-off

Implications of the mathematical result:

- To minimize expected squared error, the model has to balance between bias and variance.
- There is a trade off between bias and variance.

2.2.2 Bias-Variance trade-off

Fig 2-11: The fitted line (yellow) does not change much when the data set changes (small variance), but for example y is underestimated at $x = 70$ due to too much simplification (large bias).



2.2.3 Classification setting

A measure of error for classification problems

The **error rate** is a natural measure of training error for classification problems.

Suppose that observations are $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$ are observations. $\mathbf{y}_i \in \{1, 2, 3, \dots, G\}$ if there are G classes.

We estimate \mathbf{y}_i by \mathbf{x}_i . Let $\hat{\mathbf{y}}_i$ be the estimated class. Then,

$$\text{error rate} := \frac{1}{n} \sum_{i=1}^n I(\mathbf{y}_i \neq \hat{\mathbf{y}}_i) \quad (11)$$

where $I(\mathbf{y}_i \neq \hat{\mathbf{y}}_i) = 1$ if $\mathbf{y}_i \neq \hat{\mathbf{y}}_i$, and $I(\mathbf{y}_i \neq \hat{\mathbf{y}}_i) = 0$ if $\mathbf{y}_i = \hat{\mathbf{y}}_i$.

2.2.3 Classification setting

Test error

The expected error size $E[I(y_0 \neq \hat{y}_0)]$ is estimated by the **test error rate**, which is also defined by (11) but for the test data set.

A good classifier is one for which the ~~test~~
test error is smallest.

2.2.3 Classification setting

Bayes Classifier

Suppose that the conditional probability of Y given \mathbf{x}_0 , i.e., $P(Y = j|X = \mathbf{x}_0)$, is calculated for $j = 1, \dots, G$ and all \mathbf{x}_0 .

Then, the best classifier is to assign each observation to the most likely class. This is called the **Bayes classifier**.

The error rate for this classifier is called the **Bayesian error rate**. Given \mathbf{X} ,

$$\text{Bayesian error rate} = 1 - E\left(\max_j P(Y = j|X)\right). \quad (12)$$

2.2.3 Classification setting

Example 1:

Suppose that X is height of a person, and Y is either 'child', 'adult female' or 'adult male', and that we have following conditional probabilities.

| Height (cm) | Child | Adult Female | Adult Male |
|--------------------|-------|--------------|------------|
| $X < 145$ | 0.7 | 0.2 | 0.1 |
| $145 \leq X < 165$ | 0.2 | 0.6 | 0.2 |
| $165 \leq X$ | 0.1 | 0.3 | 0.6 |

Then, the red class is assigned for each observation. For example, if a person has height between 145cm and 165cm, 'adult female' is assigned since it is the most likely class.

2.2.3 Classification setting

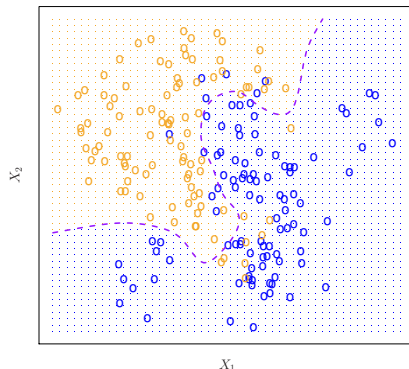
Example 2: Binary case

Consider a binary classification problem. Y is either 0 or 1. Then, assign $Y = 1$ if $P(Y = 1|X = x_0) > 0.5$ and $Y = 0$ if $P(Y = 1|X = x_0) \leq 0.5$.

The conditional probability is given by a statistical model (cf. logistic regression).

2.2.3 Classification setting

Fig 2-13: An example of Bayesian classifier. $\mathbf{X} = (X_1, X_2)$ is 2-dimensional. There are 2 classes: orange and blue. Circles are observations. Dots show classification for each possible value of \mathbf{X} . The black dashed curve is the Bayesian decision boundary, i.e., the set of \mathbf{x} satisfying $P(Y = 1 | \mathbf{X} = \mathbf{x}) = 0.5$.



2.2.3 Classification setting

K-nearest neighbors

Let K be a natural number $(1, 2, 3, \dots)$. The **K-nearest neighbors (KNN)** is a non-parametric approach to estimate conditional probabilities by

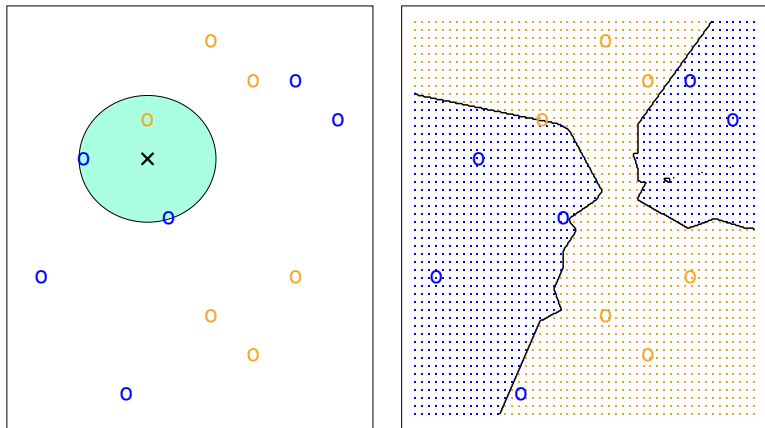
$$P(Y = j|X = x_0) = \frac{1}{K} \sum_{i \in N_0} I(y_i = j). \quad (13)$$

where N_0 is the set of K points closest to x_0 in the training set.

In other words, KNN estimates the probability by empirical results of K most similar cases.

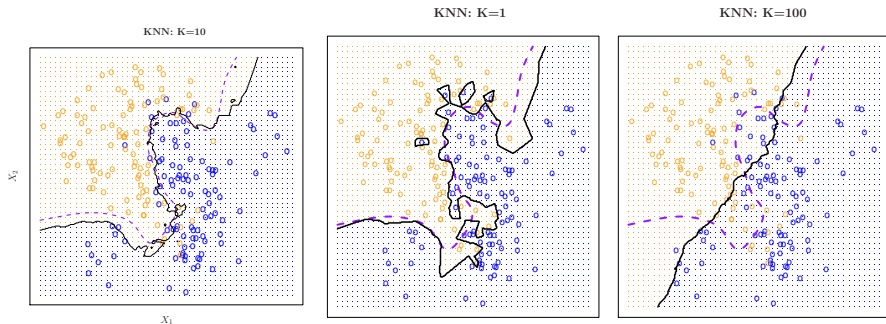
2.2.3 Classification setting

Fig 2-14: An example of KNN ($K = 3$). Left: At the cross, estimates are $P(Y = \text{blue}) = 2/3$ and $P(Y = \text{orange}) = 1/3$. Right: Bayesian classifier by this KNN.



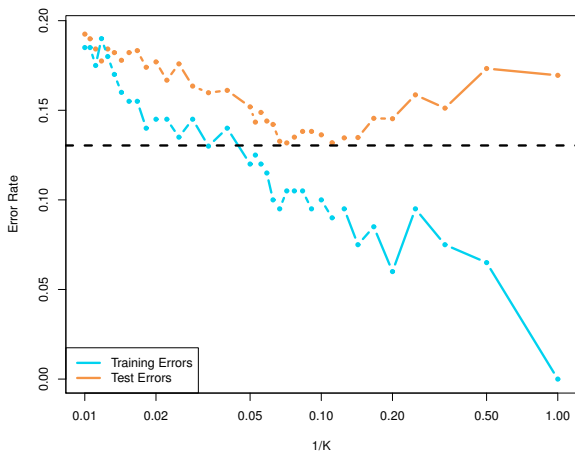
2.2.3 Classification setting

Fig 2-15&16: The result of KNN depends on K (**10, 1, 100**). A small K is more detailed; a large K is more robust. Solid lines: KNN. Purple dashed lines: Bayesian boundary by a benchmark method (no details in book).



2.2.3 Classification setting

Fig 2-17: The error rate by K in KNN. There is an optimal K to minimize the (expected) test error.



2.2.3 Classification setting

Pros and cons of a small K and a large K is similar to bias-variance trade-off in regression problems.

- A small K : small bias, large variance.
- A large K : large bias, small variance.

Memo