

STA 35C: Statistical Data Science III

Lecture 4: Statistical Learning

Dogyoon Song

Spring 2026, UC Davis

Announcements

Homework 1 is due tomorrow (Tue, Apr 7, 11:59 PM)

- Please submit on time and follow the submission instructions
- Discussion section tomorrow
- TA office hours: Tue, 12:30 PM – 2:30 PM (MSB 1143)
- Feel free to post questions on Piazza (the earlier, the better)

Agenda

Last time:

- Bayes' rule
- Random variables
 - Distribution functions
 - Expectation & variance
 - Distribution of multiple RVs
 - Covariance & correlation

Today: Statistical learning

- Motivating examples
- Supervised vs. unsupervised learning
- Prediction vs. inference
- Parametric vs. nonparametric methods

Recap: Bayes' theorem and random variables

Bayes' theorem states that $\text{Posterior} = \text{Likelihood} \times \text{Prior} / \text{Evidence}$, i.e.,

$$P(A | B) = \frac{P(B | A) P(A)}{P(B)} \quad \text{where} \quad P(B) = P(B | A)P(A) + P(B | A^c)P(A^c)$$

- For example, A: nature/model & B: data

Random variable: a map $X : \Omega \rightarrow \mathbb{R}$ assigns a real number to each outcome $\omega \in \Omega$

- The **PMF** of a discrete random variable p_X : $p_X(x) = P(X = x)$
- The **PDF** of a continuous random variable f_X : $P(a \leq X \leq b) = \int_a^b f_X(x) dx$
- The **CDF** of either variable F_X : $F_X(x) = P(X \leq x)$
- **Expectation** and **variance** summarize center and spread
 - Useful identities follow from the linearity of expectation
- Joint, marginal, and conditional distributions
 - X and Y are independent if $p_{X,Y}(x, y) = p_X(x) p_Y(y)$ for all x, y

Covariance and correlation

The **covariance** between X and Y measures how they vary together:

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y].$$

- $\text{Cov}(X, Y) > 0$: large values of X tend to occur with large values of Y
- If X and Y are independent, then $\text{Cov}(X, Y) = 0$; the converse is false in general

Correlation is the normalized version of covariance:

$$\rho(X, Y) := \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}} \in [-1, 1] \quad (\text{when } \text{Var}(X), \text{Var}(Y) > 0).$$

Example (One fair die)

Recall the example from last time: $X = \mathbf{1}\{\text{roll is even}\}$, and $Y = \mathbf{1}\{\text{roll} > 3\}$

| | $Y = 0$ | $Y = 1$ |
|---------|---------|---------|
| $X = 0$ | 1/3 | 1/6 |
| $X = 1$ | 1/6 | 1/3 |

- $\text{Cov}(X, Y) = \frac{1}{12}$
- $\rho(X, Y) = \frac{1}{3}$

Sum of random variables: Expectation and variance are easy

Sums arise naturally when we combine several sources of randomness:

- Total number of heads in n tosses: $S_n = X_1 + \dots + X_n$
- Total score, total revenue, total waiting time
- Total measurement error from several noisy components

The expectation and variance of $X + Y$ are often easy to compute:

- $\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$

$$\begin{aligned}\because \mathbb{E}[X + Y] &= \sum_{x,y} (x + y) p_{X,Y}(x,y) = \sum_x x p_X(x) + \sum_y y p_Y(y) \\ &= \mathbb{E}[X] + \mathbb{E}[Y]\end{aligned}$$

- $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$
 - *Exercise:* Verify this using properties of expectation
 - If X and Y are independent, then $\text{Cov}(X, Y) = 0$, so $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$

Sum of random variables: Distributions can be harder

Even when $\mathbb{E}[X + Y]$ and $\text{Var}(X + Y)$ are easy, the full distribution of $X + Y$ may not be

- We usually need the *joint distribution* of (X, Y) , not just the separate marginals

Example (Same marginals, different joint)

Suppose X_1 and X_2 each take values 0 or 1 with probability $1/2$

- If X_1 and X_2 are independent, then

$$P(X_1 + X_2 = k) = \begin{cases} 1/4 & k = 0, \\ 1/2 & k = 1, \\ 1/4 & k = 2 \end{cases}$$

- If $X_1 = X_2$ always, then

$$P(X_1 + X_2 = k) = \begin{cases} 1/2 & k = 0, \\ 1/2 & k = 2 \end{cases}$$

\implies Marginals do *not* determine the distribution of the sum; the joint structure matters

Pop-up quiz

Consider two tosses of a coin with sample space

$$\Omega = \{HH, HT, TH, TT\}.$$

Define

$$X(\omega) = \text{number of Heads in } \omega.$$

Which statement is correct?

- a). X is not a random variable because both HT and TH map to 1.
- b). X is an event in Ω .
- c). Since X takes values in $\{0, 1, 2\}$, the sample space is $\{0, 1, 2\}$.
- d). $\{X = 1\} = \{HT, TH\}$, so $P(X = 1) = 1/2$.

Answer: D.

A random variable is a function from outcomes to real numbers. Different outcomes are allowed to map to the same value, and $\{X = 1\}$ is an event (a subset of Ω).

Motivating example: Predicting sales

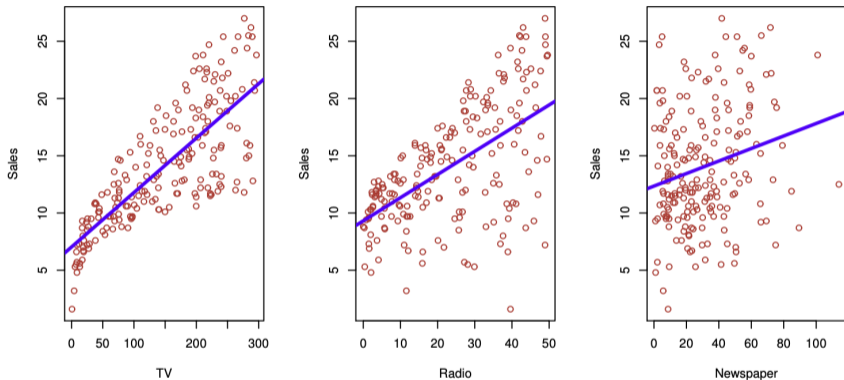


Figure: The Advertising data set shows Sales of a product in 200 different markets against advertising budgets for three media: TV, Radio, and Newspaper [JWHT21, Figure 2.1].

Question: Can we predict Sales from the budgets TV, Radio, and Newspaper?

Motivating example: Income and predictors

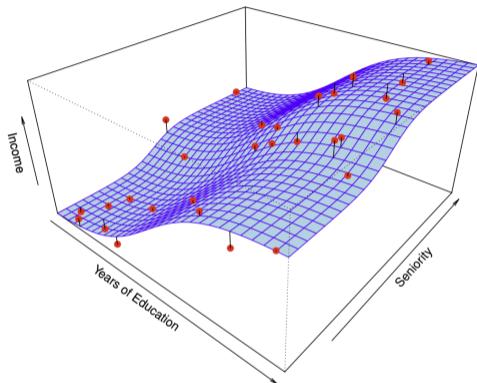


Figure: The simulated **Income** data set displays **Income** of 30 individuals as a function of **Years of education** and **Seniority** [JWHT21, Figure 2.3].

Question: Which of the predictors, e.g., **Years of education** and **Seniority**, are most strongly associated with **Income** and how?

Motivating example: Clustering unlabeled data

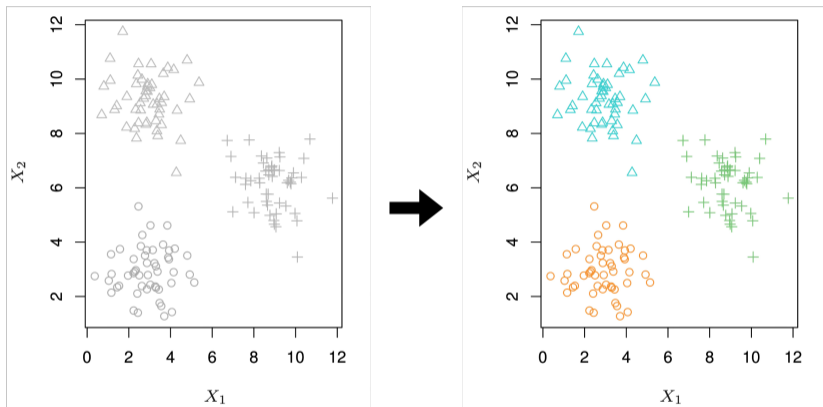


Figure: A synthetic data set with three underlying groups; colors are shown only for illustration. [JWHT21, Figure 2.8, modified].

Question: If there are no labels given, can we still discover meaningful groups?

Statistical learning: Supervised vs. unsupervised learning

Most statistical learning problems fall into two categories: supervised or unsupervised

In **supervised learning**:

- Each predictor observation x_i is accompanied by a response y_i
- “Supervised” because the responses guide (supervise) the analysis
- Many classical statistical learning methods operate in the supervised learning domain
 - *Example*: linear regression, logistic regression, support vector machine, ...

In **unsupervised learning**:

- We have observations x_i but no response y_i
- “Unsupervised” because there is no response to guide the analysis
- Often used to explore relationships among observations or variables
 - *Example*: Cluster analysis, dimension reduction, ...

Some problems lie between these two categories and the distinction can be less clear-cut

Illustration: Supervised vs. unsupervised learning

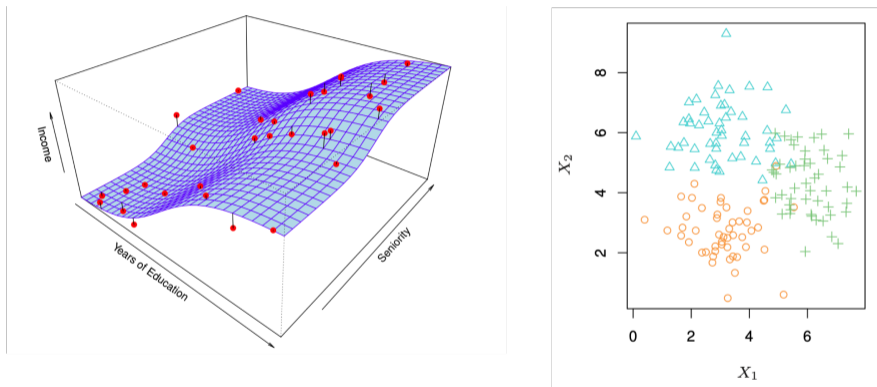


Figure: Supervised vs. unsupervised learning

For now, we focus on supervised learning; will get to unsupervised learning in Weeks 9–10

Statistical learning: Terminology and notation

Response (dependent variable) Y :

- The output variable we want to predict (e.g., **Sales**)

Predictors (independent variables, features) X :

- Input variables used to predict Y (e.g., **TV**, **Radio**, **Newspaper**)
- Often multiple predictors are collectively denoted by $X = (X_1, X_2, \dots, X_p)$

A common working assumption is that there is some relationship between Y and X :

$$Y = f(X) + \epsilon,$$

where

- f is some fixed but **unknown** function.
- ϵ is a **random** noise, which has *mean zero*, and is *independent of X* .

Goal: Learn an estimator \hat{f} of the unknown function f

Why learn f ? Three concrete goals

Predicting Y :

- We often have input variables X but not the corresponding output Y
- With an estimate \hat{f} , we can *predict* Y at new points $X = x$ via $\hat{Y} = \hat{f}(x)$
- *Example*: X = patient's blood sample, Y = risk of a disease or adverse reactions

Identifying relevant predictors:

- We can determine *which* predictors among X_1, \dots, X_p are important in explaining Y
- *Example*: Some predictors may be strongly associated with income, while others add little predictive value

Understanding how Y changes with X :

- If f is not too complex, we can interpret *how* each predictor affects Y
- *Example*: Understanding how predicted sales changes as TV advertising increases, holding the other predictors fixed

These goals mostly fall under two broader objectives: *prediction* and *inference*

Two broad goals: Prediction vs. inference

Prediction

- *Objective:* Make accurate prediction of Y given X
- \hat{f} can be treated as a “black box,” prioritizing predictive accuracy over exact form
- *Examples:*
 - Which individuals, based on demographics, are likely to respond positively to a mailer?
 - Based on blood sample, is a patient at high risk of a severe adverse drug reaction?

Inference

- *Objective:* Understand the association between Y and X
- We cannot treat \hat{f} as a black box; we need to know its exact form
- *Examples:*
 - Which media are linked to higher sales?
 - Which medium generates the largest boost in sales?
 - How much of an increase in sales is attributable to a given increase in TV advertising?

Prediction error: Reducible vs. irreducible

Under the working model $Y = f(X) + \epsilon$, prediction error of $\hat{Y} = \hat{f}(X)$ has two parts:

- **Reducible error:** error from approximating f by \hat{f}
- **Irreducible error:** randomness in ϵ , which no method can eliminate
 - ϵ may include *unmeasured* variables important for predicting Y
 - ϵ may also reflect *inherent* fluctuations (e.g., day-to-day or manufacturing variation)

At a fixed input $X = x$,

$$\mathbb{E} \left[(Y - \hat{f}(x))^2 \mid X = x \right] = \underbrace{(f(x) - \hat{f}(x))^2}_{\text{reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{irreducible}}, \quad (\because \mathbb{E}[\epsilon \mid X = x] = 0)$$

Goal: Estimate f that (1) minimizes the reducible error and (2) is interpretable

How do we estimate f ?

We will explore multiple approaches to estimate f from data

We use *training data* $\{(x_1, y_1), \dots, (x_n, y_n)\}$ to fit \hat{f}

Our goal: ensure \hat{f} generalizes well to future data (X, Y)

Most statistical learning methods are either **parametric** or **non-parametric**

- **Parametric (model-based) approach:**

- Step 1: Assume a functional form (model) of f (e.g., linear $Y = \alpha + \beta X$)
- Step 2: Use the training data to fit model parameters

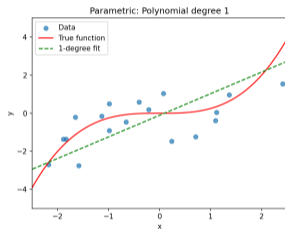
- **Non-parametric approach:**

- Make no explicit assumption about the functional form of f
- Instead, seek an \hat{f} that fits data closely while remaining sufficiently smooth

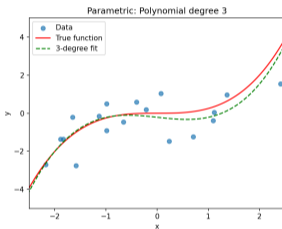
Illustration of parametric methods

Suppose that $Y = f(X) + \epsilon$ where $f(x) = \frac{1}{4}x^3$

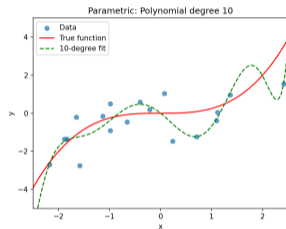
Parametric methods: e.g., polynomial regression $\hat{f}(x) = \sum_{j=0}^d \beta_j x^j$



(a) Linear regression (deg 1)



(b) Polynomial regression (deg 3)



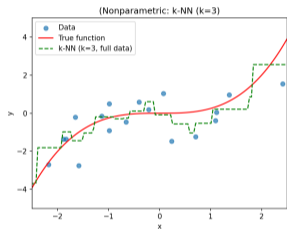
(c) Polynomial regression (deg 10)

- **Strength:** Estimating parameters β_0, \dots, β_d is easier than estimating an arbitrary function f
- **Limitation:** The assumed model may not match the true functional form of f
- **Tradeoff:** Choosing a more flexible model can reduce bias but risks *overfitting*

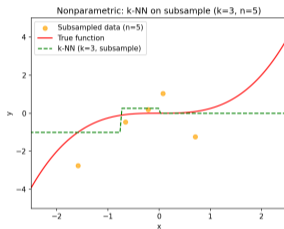
Illustration of non-parametric methods

Suppose that $Y = f(X) + \epsilon$ where $f(x) = \frac{1}{4}x^3$

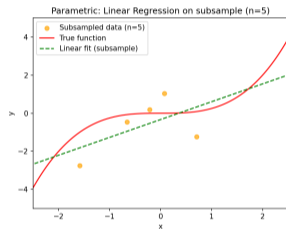
Non-parametric methods: e.g., k -nearest neighbors (k -NN)



(a) k -nearest neighbor ($k = 3$)



(b) k -NN ($k = 3$, subsampled)



(c) Linear regression (subsampled)

- **Strength:** Highly flexible; avoids the danger of using a wrong functional form
- **Limitation:** Requires more data to accurately estimate f & interpretation is more difficult
- **Tradeoff:** Flexibility can increase the overfitting risk & computation can explode at query

Tradeoff: Prediction accuracy vs. model interpretability

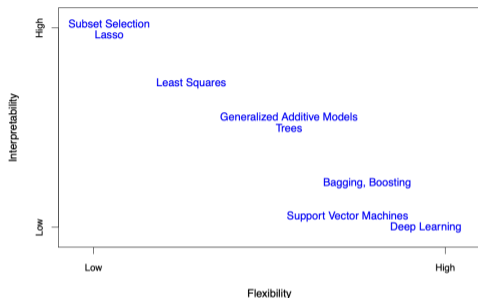


Figure: A representation of the tradeoff between flexibility and interpretability [JWHT21, Figure 2.7].

More flexible methods can capture a much wider range of relationships between X and Y , but we may still prefer more restrictive approaches because of:

- **Interpretability:** Restrictive (parametric) models are typically easier to interpret
- **Sample complexity:** Flexible models often require more observations
- **Risk of overfitting:** Very flexible methods can fit noise ϵ rather than true f

Wrap-up

- Statistical learning studies how predictors X can be used to predict or understand a response Y
- A common working model is $Y = f(X) + \epsilon$
 - f captures systematic signal
 - ϵ captures noise
- Supervised learning uses labeled responses Y ; unsupervised learning looks for structure without them
- Two main goals are prediction and inference; the favored models depend on goals
- Parametric methods are usually simpler and more interpretable; non-parametric methods are more flexible but often need more data and can overfit

Next lecture: linear regression

References



Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani.

An Introduction to Statistical Learning: with Applications in R, volume 112 of *Springer Texts in Statistics*.

Springer, New York, NY, 2nd edition, 2021.