

Summary of Chapter 4: Classification

Introduction

Chapter 4 focuses on **Classification** methods, which aim to predict a *qualitative* response (also often called a *categorical* response). Unlike regression techniques that model quantitative outcomes, classification methods assign observations to one of several discrete classes. Classification techniques appear in a variety of real-world applications, such as determining whether a patient has a particular disease (yes/no), deciding if a credit card transaction is fraudulent (yes/no), or classifying the category of an image (e.g., dog, cat, bird).

This chapter presents multiple approaches to classification. It begins with a rationale for why simple linear regression is usually not suitable for a qualitative response. Then it introduces **Logistic Regression** as a foundational approach for binary classification. The chapter subsequently covers **Linear Discriminant Analysis (LDA)**, **Quadratic Discriminant Analysis (QDA)**, **Naive Bayes**, and **K-Nearest Neighbors (KNN)**. It also outlines how to evaluate and compare these methods, including the use of confusion matrices, sensitivity, specificity, and ROC curves. Towards the end, it discusses *Generalized Linear Models (GLMs)* more broadly, focusing on Poisson Regression as another instance of the GLM family that is suited to count data rather than strictly binary or continuous responses. Finally, the chapter includes labs demonstrating each of these methods in the R programming environment.

In this summary, each major section is reviewed in turn, emphasizing the core concepts, the assumptions behind the models, and the practical performance considerations.

1 Overview of Classification

Classification problems involve a response variable Y taking on qualitative values (classes). The fundamental goal is to assign an observation to the correct category based on a set of predictors $\mathbf{X} = (X_1, X_2, \dots, X_p)$. Examples include:

- Determining a medical diagnosis (e.g., **Stroke**, **Overdose**, or **Seizure**) given a patient's vital signs and other clinical predictors.
- Predicting whether an individual will default on a credit payment (**Yes** vs. **No**) based on financial predictors such as credit balance and income.
- Detecting fraudulent online transactions (**Fraudulent** vs. **Legitimate**).

The data consist of *training samples* $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, where each observation has p predictors \mathbf{x}_i and a class label y_i . We seek to build a model that can correctly classify new, previously unseen observations.

A concept repeatedly emphasized in the chapter is that classification methods aim either:

- to *directly* estimate the posterior probability $P(Y = k \mid \mathbf{X} = \mathbf{x})$ for each class k (e.g., logistic regression, naive Bayes), or
- to *indirectly* model the distribution of predictors \mathbf{X} *within each class* and then apply *Bayes' theorem* (as in LDA and QDA).

1.1 Why Not Use Linear Regression?

One might attempt to fit a linear regression of the form $Y \approx \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p$, even if Y is binary (coded 0 or 1). However, Chapter 4 points out several issues:

1. Predicted values might lie outside $[0, 1]$, making them invalid probabilities.
2. Linear regression treats the response as quantitative, implying equal spacing between classes when more than two classes exist.
3. Statistical properties (like normal residual assumptions) break down when modeling a categorical response using linear regression.

Hence, specialized classification methods are preferred.

2 Logistic Regression

2.1 Logistic Model Formulation

Logistic regression addresses the binary response setting (two classes, often labeled **Yes** vs. **No**, or 1 vs. 0). Instead of modeling Y itself, we model

$$p(\mathbf{x}) = P(Y = 1 \mid \mathbf{X} = \mathbf{x}),$$

the probability of class 1 given predictors \mathbf{x} . Logistic regression posits that

$$\log \left(\frac{p(\mathbf{x})}{1 - p(\mathbf{x})} \right) = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p,$$

where $\log \left(\frac{p}{1-p} \right)$ is known as the *logit* or *log-odds* function. Hence, it ensures $0 \leq p(\mathbf{x}) \leq 1$ for any linear function of \mathbf{x} .

2.2 Coefficient Estimation & Interpretation

The parameters β_j are usually estimated via *maximum likelihood estimation (MLE)*:

$$\ell(\beta_0, \beta_1, \dots, \beta_p) = \prod_{i=1}^n [p(\mathbf{x}_i)]^{y_i} [1 - p(\mathbf{x}_i)]^{(1-y_i)}.$$

Solving this optimization yields $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$, which maximize the likelihood of the observed data. The estimates are often summarized in terms of z -statistics and p -values, analogous to linear regression but using logistic-specific assumptions.

Each β_j can be interpreted in terms of *odds ratio*: a one-unit increase in x_j multiplies the odds $p/(1-p)$ by e^{β_j} , holding other predictors fixed.

2.3 Prediction and Classification

After fitting the model, one obtains predicted probabilities

$$\hat{p}(\mathbf{x}) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p}}.$$

A default or typical rule is to predict **Yes** if $\hat{p}(\mathbf{x}) > 0.5$ and **No** otherwise, though thresholds can be tuned depending on the problem context (e.g., cost of false positives vs. false negatives).

2.4 Example: Default Data

The chapter uses **Default** data to predict whether an individual will default on a credit card. Two main predictors are credit **balance** and **income**. Logistic regression shows that **balance** is strongly associated with default probability. Adjusting the threshold can shift the balance between correctly capturing defaulters (high *sensitivity*) and avoiding false alarms.

3 Linear Discriminant Analysis (LDA)

3.1 Generative Perspective and Bayes' Theorem

LDA provides an alternative approach by first modeling the distribution of the predictors X *within each class*. Under Bayes' theorem, the posterior probability that $Y = k$ given $X = \mathbf{x}$ is

$$P(Y = k | \mathbf{x}) = \frac{\pi_k f_k(\mathbf{x})}{\sum_{\ell=1}^K \pi_\ell f_\ell(\mathbf{x})},$$

where $\pi_k = P(Y = k)$ is the prior probability of class k , and $f_k(\mathbf{x})$ is the class-conditional density. For LDA, we assume

$$\mathbf{x} | (Y = k) \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}),$$

that is, each class has a normal (Gaussian) distribution with class-specific mean vector $\boldsymbol{\mu}_k$ but *common* covariance matrix $\boldsymbol{\Sigma}$.

3.2 LDA Decision Boundaries

Plugging in the Gaussian densities into the Bayes formula shows that the decision rule is based on

$$\delta_k(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k - \frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k + \log \pi_k.$$

LDA classifies to the class k for which $\delta_k(\mathbf{x})$ is largest. Because \mathbf{x} enters linearly ($\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k$), the boundaries are *linear* in x .

3.3 Example: Smarket Data

The chapter applies LDA to **Smarket**, which includes daily financial returns on the S&P 500. It demonstrates how one can estimate π_k , $\boldsymbol{\mu}_k$, and $\boldsymbol{\Sigma}$ from training data, then predict whether the market will go **Up** or **Down** on a future day. LDA often yields decision boundaries that are well-suited to roughly linear separation of classes in the predictor space.

4 Quadratic Discriminant Analysis (QDA)

QDA is similar to LDA but allows each class to have its own covariance matrix Σ_k . That is, we still assume

$$\mathbf{x} \mid (Y = k) \sim \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k),$$

but *no* longer require $\Sigma_1 = \Sigma_2 = \dots = \Sigma_K$. In that case, the decision boundary derived from Bayes' rule is *quadratic* in x , so QDA can capture more complex class boundaries.

However, QDA requires estimating many more parameters (each Σ_k has $p(p+1)/2$ parameters), and so is more prone to high variance if n (the sample size) is not large enough. When the class covariances truly differ, QDA can outperform LDA.

5 Naive Bayes

Another generative classification method is **Naive Bayes**, which also uses Bayes' theorem to compute posterior probabilities $P(Y = k \mid \mathbf{x})$. However, instead of assuming *multivariate* normal distributions for X , it makes a “naive” conditional independence assumption: within each class k , the p predictors X_1, \dots, X_p are treated as independent. Formally:

$$f_k(\mathbf{x}) = \prod_{j=1}^p f_{kj}(x_j).$$

Though naive in many real applications (since predictors often correlate), this assumption drastically reduces the number of parameters. Naive Bayes can work surprisingly well, especially in high-dimensional contexts or when n is small. It does well if the independence assumption is even *approximately* met, or if reducing variance outweighs the bias introduced by the assumption of independence.

6 K-Nearest Neighbors (KNN)

KNN is a fully non-parametric approach:

1. Choose K , the number of neighbors.
2. For a new observation \mathbf{x}_0 , find the K training points closest to \mathbf{x}_0 (usually in Euclidean distance).
3. Predict the class that has the plurality among these K neighbors.

By avoiding assumptions on distributions or linearity, KNN can adapt to highly non-linear decision boundaries. However, it can suffer when p is large (the *curse of dimensionality*) or when n is small relative to p . Choosing K typically involves a bias-variance trade-off: small K yields flexible boundaries (low bias but high variance), while large K yields smoother boundaries (higher bias but lower variance).

7 Comparison of Classification Methods

7.1 Statistical Properties & Decision Boundaries

The chapter dedicates a thorough comparison to LDA, QDA, naive Bayes, logistic regression, and KNN. In short:

- *Logistic regression* and *LDA* share linear log-odds, but logistic regression arises from a direct modeling of $P(Y = k | X)$, while LDA is a generative model that also assumes Gaussian distributions with a shared covariance.
- *QDA* extends LDA by allowing separate covariance matrices, enabling more flexible, *quadratic* decision boundaries.
- *Naive Bayes* relies on an independence assumption for the features within each class but can drastically reduce complexity. It can be viewed as a special case of more general classifiers under certain conditions.
- *KNN* makes minimal distributional assumptions, can capture highly non-linear boundaries, but may struggle with large p or insufficient data.

7.2 Performance Metrics

Central to evaluating classifiers is the confusion matrix:

	Predicted Yes	Predicted No
Actual Yes	TP	FN
Actual No	FP	TN

From this table, measures such as:

- **Accuracy:** $(TP + TN)/(\text{Total})$,
- **Sensitivity** (recall): $TP/(TP + FN)$,
- **Specificity:** $TN/(TN + FP)$,
- **Positive Predictive Value** (precision): $TP/(TP + FP)$,

are used to compare how well each model identifies classes.

7.3 ROC Curves and AUC

Another important tool is the **ROC curve**, which plots *true positive rate* (sensitivity) versus *false positive rate* ($1 - \text{specificity}$) for different thresholds. The **Area Under the ROC Curve (AUC)** measures overall performance across all possible thresholds. An ideal classifier has an AUC near 1, while a random-guess classifier has AUC near 0.5.

7.4 Empirical Comparisons

In numerical experiments, the chapter demonstrates that no single method is universally the best:

- When class boundaries are truly linear, *LDA* or *logistic regression* often excel.
- When boundaries are moderately non-linear, *QDA* or *naive Bayes* might be preferable.
- For very complex boundaries, *KNN* can be superior if sufficient training data are available.

Ultimately, model choice depends on the actual data-generating process, the sample size, the dimensionality of predictors, and considerations of interpretability.

8 Generalized Linear Models (GLMs)

8.1 Beyond Binary Outcomes

The text then broadens to the concept of **generalized linear models**, in which linear regression (for continuous outcomes) and logistic regression (for binary outcomes) are special cases. A GLM is characterized by:

1. A *distribution* in the exponential family (e.g. normal, Bernoulli, Poisson).
2. A *link function* $\eta(\mu)$ relating the mean of the response Y to a linear combination of predictors:

$$\eta(\mathbb{E}[Y]) = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p.$$

8.2 Poisson Regression

An important example is **Poisson regression**, used for *count data* where Y takes integer values in $\{0, 1, 2, \dots\}$. The Poisson distribution for Y with mean λ is:

$$P(Y = k) = \frac{\lambda^k e^{-\lambda}}{k!},$$

and Poisson regression sets

$$\log(\lambda(\mathbf{x})) = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p.$$

This ensures the predicted mean $\lambda(\mathbf{x})$ is always positive, consistent with counts. For instance, the chapter uses **Bikeshare** data to model the number of hourly bike rentals. A linear model might give negative predictions or fail to reflect the mean-variance relationship in count data. By contrast, Poisson regression properly handles that the variance grows with the mean.

9 Lab: Classification Methods

In the final part of the chapter, a series of labs in **R** illustrate how to implement each classification technique:

- **Logistic Regression** on stock market (**Smarket**) data to predict up/down movements based on historical returns.

- **LDA** and **QDA** on the same data, highlighting that these two methods can yield different decision boundaries and performance.
- **Naive Bayes** using the `naiveBayes` function to demonstrate how the conditional independence assumption can simplify computations.
- **KNN** on both `Smarket` and `Caravan` data sets, showing how choosing K and standardizing features can significantly impact results.
- **Poisson Regression** on the `Bikeshare` data, illustrating how to handle count data effectively.

These hands-on examples reinforce the comparative strengths and weaknesses of each approach and underscore the necessity of careful data splitting (training vs. testing), performance measurement, and threshold tuning.

10 Conclusion

Chapter 4 provides a thorough treatment of classification, covering theory, practical considerations, and real-world examples. Central lessons include:

1. **Match the model to the data generation process when possible.** For linear boundaries and moderate data size, LDA/logistic can excel. If covariances differ by class, QDA is a candidate. If independence approximations hold or p is large, naive Bayes can be strong. If n is large or boundaries are highly non-linear, KNN is an option.
2. **Always assess performance on held-out (test) data.** Training error is usually overly optimistic.
3. **Consider interpretability.** Logistic regression yields probability estimates and interpretable coefficients. KNN is highly flexible but does not yield explicit coefficients or easily interpretable decision rules.
4. **Adjust thresholds to suit domain needs.** Costs of false positives and false negatives are often not symmetric.
5. **Generalized Linear Models are a unifying framework.** Logistic and Poisson regression are both examples within this broader class, revealing how a suitable link function and exponential family distribution can systematically address many outcome types.

Classification remains a cornerstone of supervised learning, and these techniques lay a foundation for more advanced methods like trees, random forests, boosting, and support vector machines (covered in later chapters). The chapter's empirical evidence and practical labs highlight that understanding your data, testing multiple methods, tuning hyperparameters (such as K in KNN), and carefully selecting the correct model assumptions are all crucial to success in classification tasks.

Summary of Chapter 5: Resampling Methods

Introduction

Chapter 5 focuses on *resampling methods*, a family of techniques that repeatedly draw samples or subsets from an available dataset and re-fit a model to each of these newly formed samples. The overarching goal of these methods is to gain additional insight about the variability or accuracy of a statistical learning procedure. Two of the most widely used resampling methods, **cross-validation** and **the bootstrap**, are emphasized:

- **Cross-validation** is frequently employed to estimate test error (and thus assess model performance) or to determine an appropriate level of model complexity.
- **The bootstrap** is used to quantify the uncertainty (e.g., the standard error) of an estimate and to provide interval estimates such as confidence intervals.

The chapter underscores the fact that, although computationally more intensive than fitting a model only once, resampling methods have become quite practical with modern computing power.

1 Cross-Validation

Cross-validation (CV) is introduced as a tool for obtaining an estimate of test error or predictive accuracy without requiring an explicitly held-out test set. It addresses a key problem in applied statistical learning: the discrepancy between *training error* (the error obtained by fitting a model to the training set) and *test error* (the error on new, unseen data).

1.1 Motivation and Definitions

- When no large external test set is available, cross-validation creates multiple *pseudo-training* and *pseudo-validation* subsets from the original dataset.
- The idea is to systematically hold out or “*validate on*” portions of the dataset in turn, fit the model on the remaining part, and then compute an estimate of how well the model predicts that held-out portion.

1.2 Validation Set Approach

- This is the simplest form of CV, sometimes called the *hold-out method*.
- The dataset is randomly split into two sets: a **training set** and a **validation set**.
- The model is fitted on the training set and then evaluated on the validation set to estimate the test error.
- Drawbacks:
 1. The estimate of test error can be highly variable, because results depend on how exactly the data are split.
 2. Only a fraction of the dataset is used to fit the model (the other portion being set aside), potentially causing higher bias in the fitted model.

1.3 Leave-One-Out Cross-Validation (LOOCV)

- LOOCV is an extreme version of the validation set approach:
 1. We repeatedly fit the model n times, each time leaving out exactly one observation from the training set of size $n - 1$.
 2. The omitted single observation acts as the validation set, and the prediction error is computed for that lone observation.
 3. We then average these n validation errors to get the LOOCV estimate of the test error.
- Advantages:
 - Much less bias than a single validation set split, because almost all data is used for training each time.
 - No randomness in how folds or splits are made (the single left-out observation shifts deterministically from one iteration to the next).
- Disadvantage:
 - Potentially higher variance of the error estimate, because the fits are trained on extremely similar subsets of size $n - 1$ and hence the individual predictive fits can be highly correlated.
 - Computational cost can be high if each model is expensive to fit, because n separate fits are required (though for linear models a convenient formula reduces the cost).

1.4 k -Fold Cross-Validation

- This generalizes the idea of LOOCV, splitting the data into k roughly equal-sized parts (*folds*) rather than leaving out exactly one observation per iteration.
- For each fold:
 1. One fold is designated as the validation set.
 2. The model is fitted on the other $k - 1$ folds combined.
 3. The resulting error on the held-out fold is recorded.
- The k errors are averaged to produce the overall k -fold CV error.
- A common choice is $k = 5$ or $k = 10$.
- Bias-variance perspective:
 - LOOCV ($k = n$) can have low bias but comparatively high variance in the estimate.
 - k -fold CV with moderate k reduces correlation among the training sets and hence can reduce variance, at the cost of a small increase in bias, relative to LOOCV.

1.5 Cross-Validation for Classification

- The same concepts can be applied when the response Y is qualitative, except that we measure misclassification error instead of mean squared error.
- LOOCV or k -fold CV is equally valid for classification error estimates.

2 The Bootstrap

Whereas cross-validation is mostly concerned with estimating (and tuning) the test error of a statistical learning method, the **bootstrap** is a flexible tool that can be used to assess the variability or uncertainty of an estimator.

2.1 Bootstrap Basics

- The bootstrap, in essence, simulates sampling from the original dataset *with replacement*.
- In practice:
 1. We treat our observed dataset of n points as a stand-in for the population.
 2. We sample n points *with replacement* from our original dataset to create a *bootstrap dataset*.
 3. We compute the quantity of interest (e.g., a regression coefficient, a correlation, etc.) on this bootstrap dataset.
 4. We repeat B times, collecting B estimates, from which we can compute a standard deviation (or standard error) or form confidence intervals.

2.2 Illustrations

Estimating the Standard Error of an Estimate. A running example in the chapter uses a toy problem of estimating α in a two-asset portfolio. The bootstrap is used to evaluate how much the estimate of α might vary if we slightly changed the dataset.

Bootstrap for Regression. Similarly, we can apply the bootstrap to linear regression to assess the variability of the estimated coefficients, especially when no formula-based standard errors are readily available (or when the standard formula assumptions may be violated).

3 Lab: Cross-Validation and the Bootstrap

The chapter concludes with practical labs (in R) demonstrating the usage of:

- `cv.glm()` from the `boot` library to compute k -fold and LOOCV estimates for generalized linear models.
- Validation set approaches manually coded or with specialized functions.
- The `boot()` function for implementing the bootstrap, illustrated through examples such as estimating the standard error of regression coefficients.

These labs highlight how to implement and compare the various resampling approaches discussed in the chapter. They use datasets such as the **Auto** data (for MPG regression) and the **Portfolio** data (for the two-asset α example). The code emphasizes the mechanics of splitting data, computing cross-validation errors, and bootstrapping estimates in practical settings.

Summary of Key Points

- **Resampling** is crucial in modern statistics for accurate estimation of model performance and parameter uncertainty.
 - **Cross-validation:**
 - Provides a data-driven approach to estimate test error and select model complexity (e.g., for linear regression or classification).
 - Variants such as LOOCV and k -fold CV strike different balances between bias and variance in the error estimates.
 - **Bootstrap:**
 - Provides a widely applicable way to quantify uncertainty (e.g., compute standard errors, confidence intervals) for virtually any estimator or predictive model.
 - Achieves this by sampling with replacement from the observed data, thus emulating what would happen if new datasets were repeatedly drawn from the population.
 - By using these techniques, practitioners can better understand how a learning procedure might perform on unseen data and how reliable their estimates of model parameters truly are.
-

Summary of Chapter 6: Linear Model Selection and Regularization

1 Overview

Chapter 6 discusses methods for extending and improving upon the classical linear regression framework. In particular, it focuses on approaches that address two main issues:

1. **Prediction Accuracy:** When the number of predictors p is not small compared to the number of observations n , or when $p > n$, least squares estimates can suffer from high variance and can overfit the data. Methods that reduce variance (at the cost of a small increase in bias) often improve the predictive performance.
2. **Model Interpretability:** Traditional least squares does not naturally produce sparse models (i.e., models where some coefficients are exactly zero). Consequently, identifying the truly relevant predictors among a large set can be difficult. Certain regularization methods automatically set some coefficient estimates to zero, aiding interpretability.

The chapter covers three major classes of approaches:

- *Subset Selection:* Identifying a subset of predictors that best explain the response.
- *Shrinkage (Regularization):* Estimating coefficients by shrinking them toward zero; important examples are *ridge regression* and the *lasso*.
- *Dimension Reduction:* Transforming the original predictors into a smaller set of derived features (e.g., *principal components regression* and *partial least squares*).

This summary reviews these methods and addresses some unique considerations that arise in high-dimensional settings.

2 Subset Selection

In subset selection methods, the goal is to find a subset of predictors (from among the p available) that best models the response. The central challenge is deciding which subsets of predictors to include. The chapter describes:

2.1 Best Subset Selection

- **Idea:** Fit separate least squares regressions for all possible subsets of the p predictors. Then choose the “best” subset of each size $k = 0, 1, \dots, p$.
- **Algorithm:**
 1. For $k = 0$ to p , consider all $\binom{p}{k}$ predictor subsets of size k .
 2. Fit a least squares model for each subset.
 3. Choose the best model (lowest RSS or highest R^2) among the models of size k .
 4. Select a single model from among these $p + 1$ candidates using criteria such as cross-validation, C_p , BIC, adjusted R^2 , or validation-set error.

- **Issues:** The number of models grows exponentially (2^p). This becomes infeasible for even moderate p . The method also risks overfitting if p is large relative to n .

2.2 Stepwise Selection Methods

Because best subset selection quickly becomes computationally infeasible, simpler *stepwise* approaches are often used.

Forward Stepwise Selection

- Starts with the null model (no predictors).
- Adds predictors one at a time, at each step choosing the predictor that improves the model fit the most (e.g., yields the smallest RSS).
- Continues until all predictors are in the model (or until some stopping criterion).
- Far less computationally expensive than best subset selection.

Backward Stepwise Selection

- Starts with the full model (all p predictors).
- Iteratively removes the least useful predictor (e.g., the one whose removal increases the RSS the least).
- Requires $n > p$ so that the full model can be fit initially.

Hybrid Approaches

- Methods that add variables like forward selection but allow for occasional backward removal at each step.

2.3 Choosing the Optimal Model Size

Once all candidate models have been fit, we must select a single best model among the $p + 1$ (or fewer, in stepwise) possibilities. This can be done via:

- **Cross-Validation or Validation-Set Error:** Train each model on one portion of the data and assess its prediction error on the held-out portion.
- C_p , **AIC**, **BIC**, **Adjusted R^2 :** Each is an estimate of test error (or test error-related quantity) that penalizes model complexity differently.

Often, one selects the model that yields the minimum cross-validation error or the smallest C_p , BIC, etc.

3 Shrinkage (Regularization) Methods

Rather than selecting a subset of predictors, *shrinkage* methods involve fitting a model containing all predictors but shrinking the coefficient estimates toward zero. The two most important approaches are **ridge regression** and the **lasso**.

3.1 Ridge Regression

- **Definition:** Estimates coefficients β_j by minimizing

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2.$$

- The second term $\lambda \sum \beta_j^2$ is a *shrinkage penalty* that penalizes large coefficient estimates.
- **Effect:** As λ grows, the coefficient estimates get shrunk toward zero, reducing variance but possibly increasing bias. This can lower test error if the bias-variance trade-off is improved.
- **No Variable Selection:** Ridge regression never sets coefficients exactly to zero (unless λ is infinite). Thus, it does not yield a sparse model.

3.2 The Lasso

- **Definition:** Estimates coefficients by minimizing

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|.$$

- The penalty is now an ℓ_1 norm ($\sum |\beta_j|$), rather than the ℓ_2 norm used by ridge regression.
- **Variable Selection:** The ℓ_1 penalty forces some coefficients to be *exactly* zero when λ is sufficiently large, producing a sparse model. This yields better interpretability (fewer predictors).

3.3 Selecting the Tuning Parameter

Both ridge regression and the lasso require choosing λ . Typically one uses cross-validation:

1. Fit the model for a grid of λ values.
2. For each λ , compute the cross-validation error.
3. Select the λ that yields the smallest cross-validation error.
4. Refit on the full data using this chosen λ .

3.4 Bayesian Interpretation

- Ridge regression corresponds to assuming a Gaussian prior for the coefficients (β_j) , centered at zero.
- Lasso corresponds to a double-exponential (Laplace) prior centered at zero.
- In each case, the solution to the penalized least squares problem can be interpreted as the mode of the posterior distribution under those priors.

4 Dimension Reduction Methods

These techniques reduce the dimension of the problem by projecting the p -dimensional predictors onto a smaller M -dimensional subspace, and then regressing the response on these M derived predictors.

4.1 Principal Components Regression (PCR)

- **Principal Components Analysis (PCA):** Identify new orthogonal directions (principal components) that capture maximal variance in the predictor space.
- **PCR:** Select the first M principal components Z_1, \dots, Z_M of the predictors, and then fit a linear model to predict Y from these components.
- **Dimension Reduction:** Only $M < p$ directions are used, which can mitigate overfitting by dramatically reducing the number of fitted parameters.
- **Drawback:** The principal components are chosen *unsupervised*, i.e., without considering Y . Thus, it is not guaranteed that the components capturing largest variance in predictors are also most relevant for the response.

4.2 Partial Least Squares (PLS)

- **Similar to PCR,** but it is a *supervised* approach. The directions onto which we project the data are chosen to maximize variation in the predictors *and* correlation with the response.
- PLS can sometimes perform better than PCR because it uses Y in constructing the directions.
- Both PCR and PLS require selecting the number of components M , typically via cross-validation.

5 Considerations in High Dimensions

5.1 High-Dimensional Data

In modern applications, p (the number of predictors) can be much larger than n (the number of observations). This is called the *high-dimensional* setting. Traditional methods such as ordinary least squares can fail badly here, because:

- Perfect fits (zero residuals) on the training set are easy to achieve but do not generalize (huge test error).
- Multicollinearity is extreme: any predictor can be expressed approximately (or exactly) as a linear combination of other predictors.

5.2 Pitfalls in High Dimensions

- **Overfitting** becomes a severe threat as p grows.
- Traditional metrics like training set R^2 or residual sum of squares (RSS) can be *misleading*; they may be arbitrarily small due to overfitting.
- **Cross-validation or independent test sets** are essential to obtain honest estimates of predictive performance.
- Model interpretability is more challenging: if many different subsets of predictors yield similar fits, it may be hard to pinpoint a single “true” subset.

5.3 Methods for High Dimensions

- The methods covered in this chapter (subset selection, ridge, lasso, PCR, PLS) are often especially helpful in high-dimensional contexts.
- **Regularization** (ridge or lasso) often drastically reduces variance and can yield better prediction.
- **Sparse solutions** (from the lasso) can improve interpretability by identifying the most relevant predictors.

6 Lab: Linear Models and Regularization Methods

In the lab section of Chapter 6, various **R** techniques for implementing these methods are illustrated:

- **Subset Selection:** Using `regsubsets` (from the `leaps` package) for best subset, forward stepwise, and backward stepwise selection.
- **Shrinkage:** Using `glmnet` to fit ridge regression and the lasso over a grid of λ values and applying cross-validation with `cv.glmnet`.
- **PCR & PLS:** Using `pcr` and `pls` (from the `pls` package) with cross-validation to choose the number of components M .
- **Comparison of Test Errors:** Demonstrating how to split data into training and testing subsets, compute test MSE, and compare different methods.

The lab emphasizes practical details in **R** for model fitting, hyperparameter selection (e.g., λ in ridge/lasso, or M in PCR/PLS), and error assessment via cross-validation.

7 Summary of Key Points

1. **Subset Selection Methods** reduce the model size by including only a subset of predictors. These can be computationally expensive or infeasible for large p .
2. **Ridge Regression** shrinks coefficients by adding an ℓ_2 penalty; it does not set them to zero, but typically lowers variance.
3. **The Lasso** uses an ℓ_1 penalty, performing automatic variable selection by setting some coefficients exactly to zero.
4. **PCR and PLS** reduce dimension by working with linear combinations of predictors. PCR is unsupervised (focuses on predictor variance), whereas PLS is supervised (also looks at the response).
5. In **high-dimensional settings**, careful use of regularization or dimension reduction is essential to avoid severe overfitting and to manage interpretability.
6. **Cross-validation** is crucial for tuning parameter selection and reliable estimation of test error.