

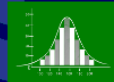
Introduction to Data Science

Winter Semester 2018/19

Oliver Ernst

TU Chemnitz, Fakultät für Mathematik, Professur Numerische Mathematik

Lecture Slides



Contents I

① What is Data Science?

② Learning Theory

2.1 What is Statistical Learning?

2.2 Assessing Model Accuracy

③ Linear Regression

3.1 Simple Linear Regression

3.2 Multiple Linear Regression

3.3 Other Considerations in the Regression Model

3.4 Revisiting the Marketing Data Questions

3.5 Linear Regression vs. K -Nearest Neighbors

④ Classification

4.1 Overview of Classification

4.2 Why Not Linear Regression?

4.3 Logistic Regression

4.4 Linear Discriminant Analysis

4.5 A Comparison of Classification Methods

⑤ Resampling Methods

Contents II

5.1 Cross Validation

5.2 The Bootstrap

6 Linear Model Selection and Regularization

6.1 Subset Selection

6.2 Shrinkage Methods

6.3 Dimension Reduction Methods

6.4 Considerations in High Dimensions

6.5 Miscellanea

7 Nonlinear Regression Models

7.1 Polynomial Regression

7.2 Step Functions

7.3 Regression Splines

7.4 Smoothing Splines

7.5 Generalized Additive Models

8 Tree-Based Methods

8.1 Decision Tree Fundamentals

8.2 Bagging, Random Forests and Boosting

Contents III

9 Support Vector Machines

- 9.1 Maximal Margin Classifier
- 9.2 Support Vector Classifiers
- 9.3 Support Vector Machines
- 9.4 SVMs with More than Two Classes
- 9.5 Relationship to Logistic Regression

10 Unsupervised Learning

- 10.1 Principal Components Analysis
- 10.2 Clustering Methods

⑨ Support Vector Machines

- 9.1 Maximal Margin Classifier
- 9.2 Support Vector Classifiers
- 9.3 Support Vector Machines
- 9.4 SVMs with More than Two Classes
- 9.5 Relationship to Logistic Regression

9 Support Vector Machines

9.1 Maximal Margin Classifier

9.2 Support Vector Classifiers

9.3 Support Vector Machines

9.4 SVMs with More than Two Classes

9.5 Relationship to Logistic Regression

Support Vector Machines

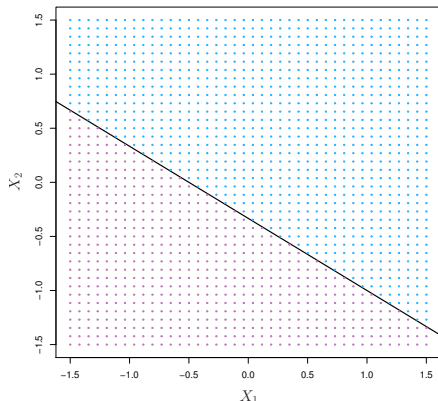
Hyperplanes for classification

- Assume a (true) bivariate classification model has a linear decision boundary.
- The points **on** the decision boundary are characterized by an equation of the form

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 = 0.$$

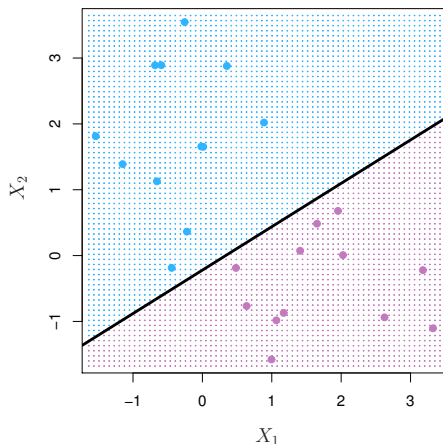
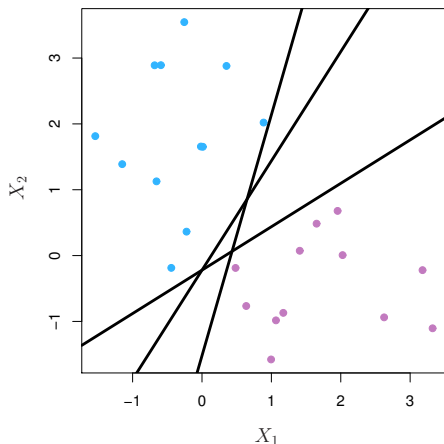
- The points on either side are characterized by

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 \begin{cases} > 0 & \text{one side,} \\ < 0 & \text{other side.} \end{cases}$$



Support Vector Machines

Hyperplanes for classification



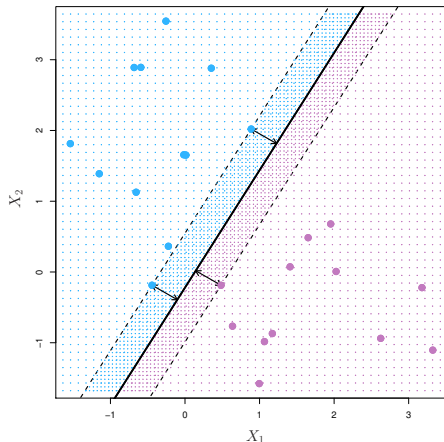
Left: Given a 2D data set of binary observations that can be split by a line, the line is in general not unique.

Right: Once a line has been fixed, this defines an associated classification model.

Support Vector Machines

Maximal margin classifier

- Compare with: KNN (Chapter 2), logistic regression, LDA, QDA (Chapter 4), decision trees (Chapter 8).
- Given binary observations, which linear decision boundary to choose?
- Separating line farthest from all training observations: **maximal margin hyperplane** or **optimal separating hyperplane**.
- Minimal distance from this line to closest observation: **margin**.
- Margin: half-width of largest slab separating the two observation set classes.
- Associated classification model: **maximal margin classifier**.



Support Vector Machines

Maximal margin classifier

- Large margin instills confidence in the classifier.
- In previous example: 3 points at marginal distance to separating hyperplane (there will always be at least 2, why?) These observations are called the **support vectors** of the maximal margin hyperplane, as moving these would move the latter.
- “Although the maximal margin classifier is often successful, it can also lead to overfitting when p is large.”

Support Vector Machines

Constructing the maximal margin classifier

- **Given:** training set $\{(\mathbf{x}_j, y_j)\}_{j=1}^n$, $\mathbf{x}_j \in \mathbb{R}^p$, $y_j \in \{-1, 1\}$.
- **Goal:** determine coefficients $\beta_0, \beta_1, \dots, \beta_p$ such that

$$\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p} \begin{cases} > 0 & \text{if } y_i = 1, \\ < 0 & \text{if } y_i = -1, \end{cases} \quad i = 1, \dots, n,$$

or, equivalently, $y_i(\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}) > 0$, $i = 1, \dots, n$, such that the margin between the resulting hyperplane and the data set is maximal.

- Formally:

maximize M as a function of $\beta_0, \beta_1, \dots, \beta_p$, M such that

$$\sum_{j=1}^p \beta_j^2 = 1, \tag{9.1a}$$

$$y_i(\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}) \geq M, \quad i = 1, \dots, n. \tag{9.1b}$$

Support Vector Machines

Constructing the maximal margin classifier

- (9.1) is a constrained optimization problem.
- Constraint (9.1b) ensures all observations on correct side of separating hyperplane if $M > 0$.
- Constraint (9.1a) not strictly necessary, merely scales M to coincide with margin.
- For details on solving the optimization problem see the reading list or a class on constrained optimization.
- Optimal separating hyperplane produces function $\hat{f}(\mathbf{x}) = \beta_0 + \boldsymbol{\beta}^\top \mathbf{x}$, $\boldsymbol{\beta} := (\beta_1, \dots, \beta_p)^\top$ from which the classifier is derived as

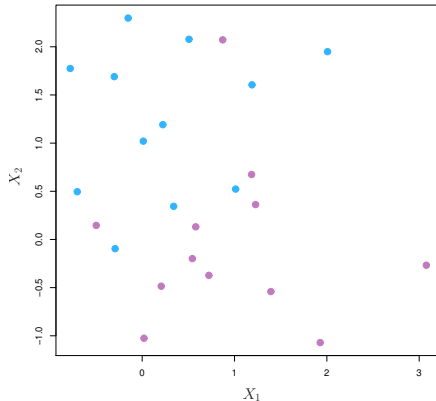
$$\hat{F}(\mathbf{x}) := \text{sign } \hat{f}(\mathbf{x}).$$

- All training observations outside margin, not necessarily so for test observations. Intuition: large margin on training data will result in good separation of test data.

Support Vector Machines

Non-separable case

- A maximal margin classifier (or even a separating hyperplane) need not exist for a given data set.
- In this case optimization problem (9.1) has no solution with $M > 0$.
- Will extend concept of separating hyperplane using so-called **soft margin** to almost separate the classes.
- This leads to **support vector classifier**.



⑨ Support Vector Machines

9.1 Maximal Margin Classifier

9.2 Support Vector Classifiers

9.3 Support Vector Machines

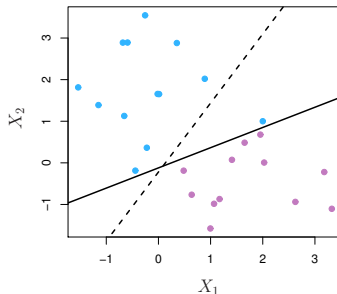
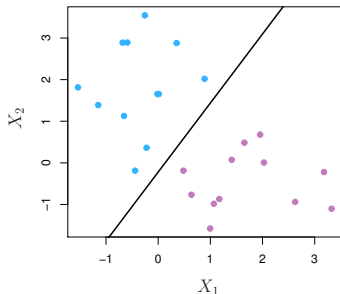
9.4 SVMs with More than Two Classes

9.5 Relationship to Logistic Regression

Support Vector Classifiers

Shortcomings of exact hyperplane classifiers

- Even if separating hyperplane exists, associated classifier possibly undesirable.
- Example: single additional blue point dramatically changes maximal margin hyperplane, leading to tiny margin.

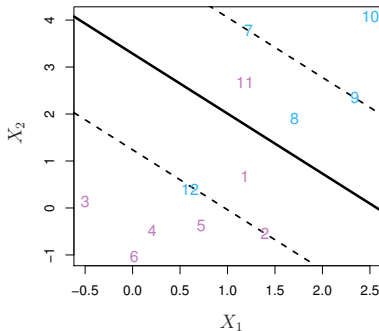
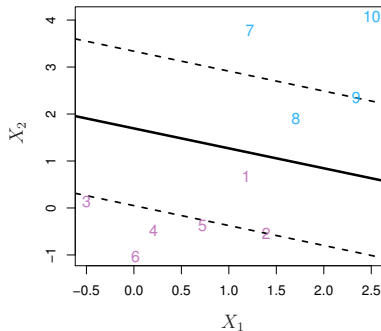


- Reasonable to accept hyperplane classifier which does **not** perfectly separate training data in the interest of greater robustness to individual observations and better classification of **most** training observations.

Support Vector Classifiers

Support vector classifiers

- **Support vector classifier (soft margin classifier)**: allow some observations to be on wrong side of margin or hyperplane. (Soft margin: can be violated by some training observations).
- Example: hyperplane (solid), margins (dashed) and observations.



Observations 11 (violet) and 12 (blue) are added in right panel.

Support Vector Machines

Support vector classifiers

- Support vector classifier solves following optimization problem:

maximize M as a function of $\beta_0, \beta_1, \dots, \beta_p, M, \epsilon_1, \dots, \epsilon_n$
subject to

$$\sum_{j=1}^p \beta_j^2 = 1, \quad (9.2a)$$

$$y_i(\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}) \geq M(1 - \epsilon_i), \quad (9.2b)$$

$$\epsilon_i \geq 0, \quad \sum_{i=1}^n \epsilon_i \leq C, \quad i = 1, \dots, n, \quad (9.2c)$$

where $C \geq 0$ is a tuning parameter.

$\epsilon_1, \dots, \epsilon_n$ are known as **slack variables**, allowing some observations to be on the wrong side of the margin or hyperplane.

- Classify test observation \mathbf{x} based on sign of $\hat{f}(\mathbf{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$.

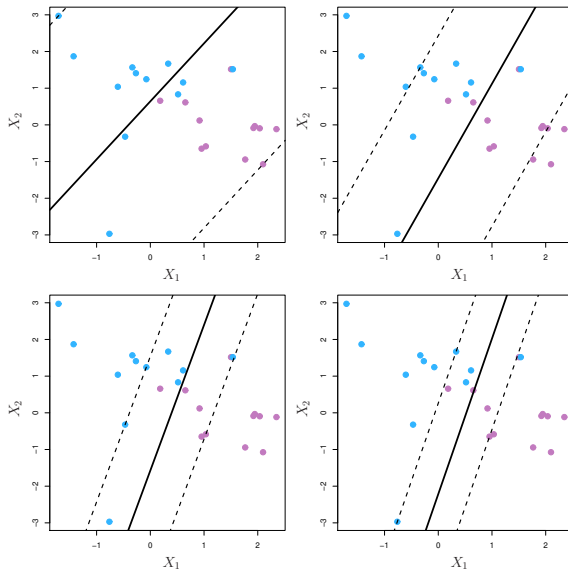
Support Vector Classifiers

Remarks on optimization problem

- Slack variable ϵ_i indicates location of \mathbf{x}_i relative to margin:
 - $\epsilon_i = 0$: \mathbf{x}_i on correct side of margin,
 - $\epsilon_i > 0$: \mathbf{x}_i on wrong side of margin,
 - $\epsilon_i > 1$: \mathbf{x}_i on wrong side of hyperplane.
- C : bounds total permitted violation, budget for amount margin can be violated. $C = 0$ implies $\epsilon_i = 0$ for all i , i.e., recover maximal margin hyperplane optimization.
 C controls bias-variance trade-off.
- Example on next slide: Support vector classifier fit using 4 different values of C (value decreases from top left to bottom right).

Support Vector Classifiers

Example



Support Vector Classifiers

Support vectors

- Only observations on or violating margin will affect hyperplane/classifier.
- Observations on correct side orf margin have no effect on support vector classifier: changing their position without violating the margin would not change the classifier.
- The remaining observations (on or in violation of margin) are called **support vectors**.
- C large: margin wide, many observations violate margin, many support vectors, many observations involved determining hyperplane, classifier has low variance, potentially high bias (top left, pevious figure)
- C small: fewer support vectors, lower bias, higher variance (bottom right, pevious figure).
- Compare LDA classifier: depends on mean of **all** observations within each class as well as within-class covariance matrix computed using all observations.
- By contrast: logistic regression, unlike LDA, has low sensitivity to observations far from decision boundary. (Will see later: logistic regression and support vector classifier closely related.)

⑨ Support Vector Machines

9.1 Maximal Margin Classifier

9.2 Support Vector Classifiers

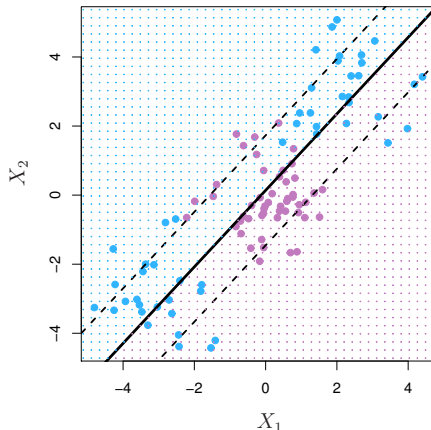
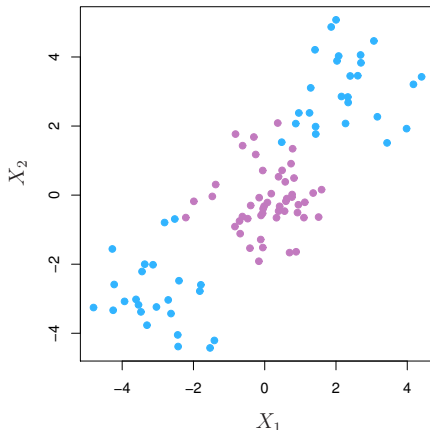
9.3 Support Vector Machines

9.4 SVMs with More than Two Classes

9.5 Relationship to Logistic Regression

Support Vector Machines

Nonlinear decision boundaries



Some data sets require decision boundaries which are curved, hence a support vector classifier will perform poorly.

Support Vector Machines

Nonlinear decision boundaries

- In Chapter 7: extended model flexibility by adding nonlinear terms.
- Analogously: in addition to features X_1, \dots, X_p , add X_1^2, \dots, X_p^2 . Then (9.2) becomes

maximize M as a function of $\beta_0, \beta_{1,1}, \beta_{1,2}, \dots, \beta_{p,1}, \beta_{p,2}, M, \epsilon_1, \dots, \epsilon_n$
subject to

$$\sum_{j=1}^p \sum_{k=1}^2 \beta_{j,k}^2 = 1, \quad (9.3a)$$

$$y_i \left(\beta_0 + \sum_{j=1}^p \beta_{j,1} x_{i,j} + \sum_{j=1}^p \beta_{j,2} x_{i,j}^2 \right) \geq M(1 - \epsilon_i), \quad (9.3b)$$

$$\epsilon_i \geq 0, \quad \sum_{i=1}^n \epsilon_i \leq C, \quad i = 1, \dots, n, \quad (9.3c)$$

- Could also add interaction terms $X_j X_k$ in this way.

Support Vector Machines

New idea

- Enlarging the feature space in this way quickly makes the computations unmanageable.
- The key computational ingredient involves using **kernels**.
- Details of solving optimization problem (9.2) for support vector classifier involve **inner product** of observations

$$\langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle = \sum_{j=1}^p x_{i,j} x_{i',j}$$

rather than observations themselves.

- Can show that linear support vector classifier has representation

$$\hat{f}(\mathbf{x}) = \beta_0 + \sum_{i=1}^n \alpha_i \langle \mathbf{x}, \mathbf{x}_i \rangle \quad (9.4)$$

in terms of n parameters $\{\alpha_i\}_{i=1}^n$ (one per training observation).

Support Vector Machines

New idea

- Estimating parameters $\beta_0, \alpha_1, \dots, \alpha_n$ involves all $n(n-1)/2$ pairwise inner products of training observations.
- It turns out that in (9.4) α_i is only nonzero if \mathbf{x}_i is a support point.
- Denoting the index set of the support points by \mathcal{S} , (9.4) becomes

$$\hat{f}(\mathbf{x}) = \beta_0 + \sum_{i \in \mathcal{S}} \alpha_i \langle \mathbf{x}, \mathbf{x}_i \rangle. \quad (9.5)$$

- This typically involves far fewer points than (9.4).
- Expanding the inner products in (9.5) establishes the relationship between the parameters α_i and the original coefficients β_j .
- **Summary:** computation and evaluation of linear classifier rests on evaluating inner products of point in feature space.
- Replace inner products $\langle \mathbf{x}, \mathbf{x}_i \rangle$ with different function $K(\mathbf{x}_i, \mathbf{x}_{i'})$ of \mathbf{x}_i and $\mathbf{x}_{i'}$ referred to as a **kernel function** or simply **kernel**.

Support Vector Machines

Kernels

- Kernels quantify the degree of similarity or strength of relationship between two points \mathbf{x}_i and $\mathbf{x}_{i'}$ of feature space.
- Could simply take inner product

$$K(\mathbf{x}_i, \mathbf{x}_{i'}) = \sum_{j=1}^p x_{i,j} x_{i',j},$$

(**linear kernel**), recovering support vector classifier. Here similarity of observations quantified using standard correlation.

- Alternatively,

$$K(\mathbf{x}_i, \mathbf{x}_{i'}) = \left(1 + \sum_{j=1}^p x_{i,j} x_{i',j} \right)^d$$

known as **polynomial kernel** of degree $d \in \mathbb{N}$. Leads to nonlinear decision boundary for support classifier, involves higher order terms in original features.

Support Vector Machines

Kernels

- Combination of support vector classifier with nonlinear kernel referred to as **support vector machine** (SVM).
- In this case the model has the form

$$\hat{f}(\mathbf{x}) = \beta_0 + \sum_{i \in \mathcal{S}} \alpha_i K(\mathbf{x}, \mathbf{x}_i). \quad (9.6)$$

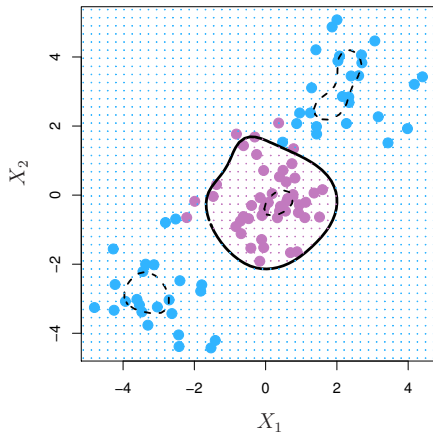
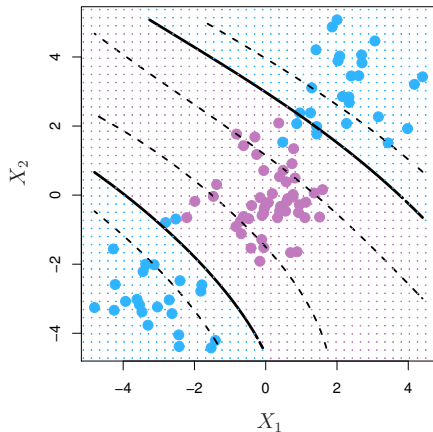
- Next picture: same data set as previous picture, SVM classifier using polynomial kernel (left panel). Right panel shows SVM using **radial kernel**

$$K(\mathbf{x}_i, \mathbf{x}_{i'}) = \exp \left(-\gamma \sum_{j=1}^p (x_{i,j} - x_{i',j})^2 \right) = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|_2^2), \quad (9.7)$$

$\gamma > 0$ a parameter.

Support Vector Machines

Kernels



Support Vector Machines

Kernels

- If test observation $\mathbf{x} \in \mathbb{R}^p$ far from training observation \mathbf{x}_i in the sense that $\|\mathbf{x} - \mathbf{x}_i\|_2$ large, then $K(\mathbf{x}, \mathbf{x}_i) = \exp(-\gamma\|\mathbf{x} - \mathbf{x}_i\|_2^2)$ will be very small.
- Therefore in (9.6) observation \mathbf{x}_i will have almost no influence on $\hat{f}(\mathbf{x})$.
- Recall: class label prediction depends on sign of $\hat{f}(\mathbf{x})$.
- Therefore, observations far away from \mathbf{x} have little influence in class prediction for \mathbf{x} . (Radial kernel has very **local** behavior.)
- Advantage of using kernel over simply enlarging feature space with nonlinear expressions in predictor variables: explicit enlarging of feature space avoided, leading to great computational savings (particularly if feature space infinite-dimensional).

⑨ Support Vector Machines

9.1 Maximal Margin Classifier

9.2 Support Vector Classifiers

9.3 Support Vector Machines

9.4 SVMs with More than Two Classes

9.5 Relationship to Logistic Regression

Support Vector Machines

SVMs with more than two classes

- Up to now: maximal margin/support vector classifiers/machines for binary classification.
- Concept rests on separating hyperplane idea, does not readily generalize to $K > 2$ classes.
- Of many proposals for this problem, two most popular are **one-versus-one** and **one-versus-all** approaches.

Support Vector Machines

One-versus-one classification

- For $K > 2$ classes, there are $\binom{K}{2} = K(K - 1)/2$ possible class pairs.
- For each pair, construct separate SVM to compare just these two classes, e.g., for pair (k_1, k_2) , code k_1 as $+1$ and k_2 as -1 .
- To each test observation, apply all $K(K - 1)/2$ SVMs, record how often observation assigned to each class.
- Finally, classify observation as belonging to class to which it was assigned most often.

Support Vector Machines

One-versus-all classification

- For $K > 2$ classes, fit K SVMs, each comparing one class with remaining $K - 1$ classes.
- Denote by $\beta_{0,k}, \beta_{1,k}, \dots, \beta_{p,k}$ the coefficients resulting from fitting SVM to determine membership in k -th class, coded as $+1$, against membership in one of the $K - 1$ remaining classes, coded as -1 .
- For test observation $\mathbf{x} = (x_1, \dots, x_p)^\top$, assign to class k for which

$$\beta_0 + \beta_{1,k}x_1 + \dots + \beta_{p,k}x_p$$

is largest, indicating high level of confidence that \mathbf{x} belongs to class k .

⑨ Support Vector Machines

9.1 Maximal Margin Classifier

9.2 Support Vector Classifiers

9.3 Support Vector Machines

9.4 SVMs with More than Two Classes

9.5 Relationship to Logistic Regression

Support Vector Machines

Relationship to logistic regression

- SVMs introduced in mid-1990s.
- Immediately successful due to performance, marketing, novelty of approach.
- Approach notably different from previously established classification methods such as logistic regression, LDA.
- Introduced kernel idea for expanding feature space and allowing nonlinear decision boundaries.
- In the meantime, deep connections between SVMs and classical statistical learning methods have been discovered.

Support Vector Machines

Relationship to logistic regression

- Can show: optimization problem (9.3) for fitting support vector classifier $f(X) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$ can be reformulated as

$$\arg \min_{\beta_0, \beta_1, \dots, \beta_p} \left\{ \sum_{i=1}^n \max\{0, 1 - y_i f(x_i)\} + \lambda \sum_{j=1}^p \beta_j^2 \right\}, \quad (9.8)$$

with $\lambda \geq 0$ a tuning parameter.

- λ large: β_0, \dots, β_p small, more margin violations tolerated, low variance, high bias.
- λ small: fewer violations tolerated high variance, low bias.
Corresponds to low value of C in (9.3c).
- Term $\lambda \sum_{j=1}^p \beta_j^2$ is exactly the ridge regression penalty term, provides regularization, controls bias-variance trade-off for support vector classifier.

Support Vector Machines

Relationship to logistic regression

- Formulation (9.8) has form “Loss + Penalty” (data misfit + regularization)

$$\arg \min_{\beta_0, \beta_1, \dots, \beta_p} (L(\mathbf{x}, \mathbf{y}, \boldsymbol{\beta}) + \lambda P(\boldsymbol{\beta})) \quad \begin{array}{l} L : \text{loss function} \\ P : \text{penalty function} \\ (\mathbf{x}, \mathbf{y}) : \text{data.} \end{array} \quad (9.9)$$

Loss function measures degree to which model fits data.

Penalty function weighted by regularization parameter λ .

- Ridge regression and the lasso both employ loss function

$$L(\mathbf{x}, \mathbf{y}, \boldsymbol{\beta}) = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{i,j} \right)^2,$$

with $P(\boldsymbol{\beta}) = \sum_{j=1}^p \beta_j^2$ for ridge regression and $\sum_{j=1}^p |\beta_j|$ for the lasso.

Support Vector Machines

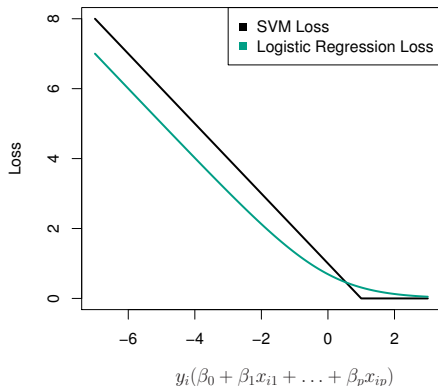
Relationship to logistic regression

- In (9.9), loss function has form

$$L(\mathbf{x}, \mathbf{y}, \boldsymbol{\beta}) = \max\{0, 1 - y_i(\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p})\}$$

known as **hinge loss**.

- Hinge loss function closely related to logistic regression loss function.



Support Vector Machines

Relationship to logistic regression

- In support classifier, only support vectors determine classifier, not observations on correct side of margin.
- Reason: loss function zero whenever $y_i(\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}) \geq 1$. Here margin corresponds to value 1 and $\sum \beta_j^2$ determines margin width.
- Logistic regression loss function does not vanish anywhere, but very small away from decision boundary.
- SVM's behave better for well-separated classes, logistic regression preferred when more overlap present.
- Importance of parameter C not initially realized.
- Nonlinear kernels possible also in logistic regression, but not as common in practice.
- Extension of SVM to regression: **support vector regression**. Uses different loss function than LS regression, where only residuals above positive threshold contribute (extension of margin concept from classification to regression).