

Linear Models for Classification

Three Approaches to Classification Problems

1. Deterministic model
2. Probabilistic generative model
3. Probabilistic deterministic model

Linear Models for Classification

Goal: Take an input vector \mathbf{x} and assign it to one of K classes C_1, \dots, C_K .

The input space is divided into regions whose boundaries are called **decision boundaries** or **decision surfaces**.

Linear classification models yield linear decision boundaries, which are $(D - 1)$ -dimensional hyperplanes within the D -dimensional input space.

4.1 Discriminant Functions

A **discriminant** is a function that takes an input vector \mathbf{x} and assigns it to one of K classes C_1, \dots, C_K . A **linear discriminant** is a function whose decision surfaces are hyperplanes.

The simplest representation of a linear discriminant function is

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

where \mathbf{w} is a weight vector and w_0 is a bias.

Binary Classification ($K = 2$)

If $y(\mathbf{x}) \geq 0$, assign C_1 . Otherwise, assign C_2 .

The decision surface is defined by $y(\mathbf{x}) = 0$, which is a $(D - 1)$ -dimensional hyperplane in the D -dimensional input space.

If \mathbf{v} is any vector that lies within the decision surface, then $\mathbf{w}^T \mathbf{v} = 0$, so \mathbf{w} is orthogonal to every vector in the decision surface. \mathbf{w} determines the orientation of the decision surface.

The bias term w_0 determines the location of the decision boundary.

Classification Involving Multiple Classes ($K > 2$)

Use a single K -class discriminant comprising K linear functions of the form

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}, \text{ for } k = 1, \dots, K.$$

Assign \mathbf{x} to class C_k if $y_k(\mathbf{x}) > y_j(\mathbf{x})$ for all $j \neq k$. The decision boundary between class C_k and class C_j is given by $y_k(\mathbf{x}) = y_j(\mathbf{x})$, which is a $(D - 1)$ -dimensional hyperplane defined by

$$(\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0.$$

Classification Involving Multiple Classes ($K > 2$)

The decision regions of such a discriminant are **convex**. To see this, if $\mathbf{x}_A, \mathbf{x}_B$ are two points lying in the same decision region \mathcal{R}_k , the line $\hat{\mathbf{x}}$ connecting \mathbf{x}_A and \mathbf{x}_B can be expressed as

$$\hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B, \text{ where } 0 \leq \lambda \leq 1.$$

By linearity of the discriminant function,

$$y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda) y_k(\mathbf{x}_B).$$

Since $y_k(\mathbf{x}_A) > y_j(\mathbf{x}_A)$ and $y_k(\mathbf{x}_B) > y_j(\mathbf{x}_B)$ for all $j \neq k$, $y_k(\hat{\mathbf{x}}) > y_j(\hat{\mathbf{x}})$, so $\hat{\mathbf{x}}$ lies inside \mathcal{R}_k .

Learning the Parameters of Linear Discriminant Functions

Three approaches:

1. Least squares
2. Fisher's linear discriminant
3. Perceptron algorithm

Least Squares

For ease of notation, let

$$\tilde{\mathbf{w}}_k = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix}, \quad \tilde{\mathbf{x}} = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}, \quad \tilde{\mathbf{W}} = \begin{bmatrix} \tilde{\mathbf{w}}_1 & \tilde{\mathbf{w}}_2 & \cdots & \tilde{\mathbf{w}}_K \end{bmatrix}.$$

Then we can combine the K equations

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}, \text{ for } k = 1, \dots, K$$

into

$$y(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}}.$$

Least Squares

Note: D = input space dimension, N = number of samples in the training data

We determine the parameter matrix \tilde{W} by minimizing a sum-of-squares error function.

Let

$$T = \begin{bmatrix} \tilde{t}_1^T \\ \tilde{t}_2^T \\ \vdots \\ \tilde{t}_N^T \end{bmatrix}$$

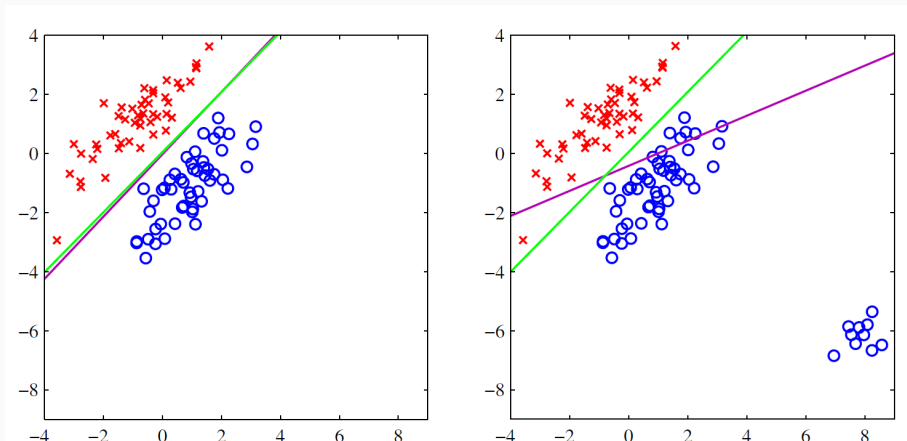
where the training data set is $\{\mathbf{x}_n, \mathbf{t}_n\}$ for $n = 1, \dots, N$. Then the sum-of-squares error function is

$$E(\tilde{W}) = \frac{1}{2} \text{Tr}\{(\tilde{X}\tilde{W} - T)^T(\tilde{X}\tilde{W} - T)\}.$$

The solution for \tilde{W} is

$$\tilde{W} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T T = \tilde{X}^\dagger T.$$

Drawbacks of Least Squares



Least squares is highly sensitive to outliers.

Fisher's Linear Discriminant

We can view linear classification as dimensionality reduction.

First look at binary classification ($K = 2$).

Consider taking a D -dimensional input vector \mathbf{x} and projecting it down to one dimension using

$$y = \mathbf{w}^T \mathbf{x}.$$

Projection onto one dimension leads to a loss of information, and classes that are well separated in the original D -dimensional space may become overlapping in one dimension.

By adjusting the components of the weight vector \mathbf{w} , we can select a projection that maximizes the class separation.

Fisher's Linear Discriminant

Suppose there are N_1 points of class C_1 and N_2 points of class C_2 . The mean vectors of the two classes are

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in C_2} \mathbf{x}_n.$$

The simplest measure of the separation of the classes, when projected onto \mathbf{w} , is the separation of the projected class means. So we can choose \mathbf{w} to maximize

$$m_2 - m_1 = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)$$

where

$$m_k = \mathbf{w}^T \mathbf{m}_k$$

is the mean of the projected data from class C_k .

Fisher's Linear Discriminant

Constrain \mathbf{w} to have unit length, so that $\sum_i w_i^2 = 1$. Perform constrained maximization using Lagrange multiplier. We get $\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$. This can still have problems, such as having considerable class overlap in the projected space.

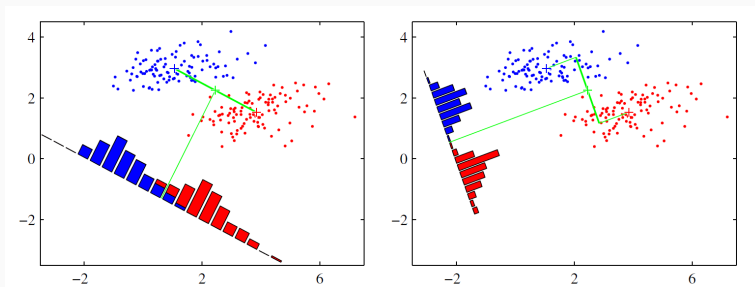


Figure 1: Left: projection with class overlap; Right: projection based on the Fischer linear discriminant

Fisher's Linear Discriminant

Maximize a function that will give a large separation between the projected class means while also giving a small variance within each class.

The within-class variance of the transformed data from class C_k is

$$s_k^2 = \sum_{n \in C_k} (y_n - m_k)^2$$

where $y_n = \mathbf{w}^T \mathbf{x}$.

Define the total within-class variance for the whole data set to be $s_1^2 + s_2^2$.

The Fisher criterion is defined as the ratio of the between-class variance to the within-class variance:

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}.$$

Fisher's Linear Discriminant

We can rewrite $J(\mathbf{w})$ as

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \quad (4.26)$$

where \mathbf{S}_B is the *between-class* covariance matrix and is given by

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \quad (4.27)$$

and \mathbf{S}_W is the total *within-class* covariance matrix, given by

$$\mathbf{S}_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T. \quad (4.28)$$

Differentiating (4.26) with respect to \mathbf{w} , we find that $J(\mathbf{w})$ is maximized when

$$(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}. \quad (4.29)$$

From (4.27), we see that $\mathbf{S}_B \mathbf{w}$ is always in the direction of $(\mathbf{m}_2 - \mathbf{m}_1)$. Furthermore, we do not care about the magnitude of \mathbf{w} , only its direction, and so we can drop the scalar factors $(\mathbf{w}^T \mathbf{S}_B \mathbf{w})$ and $(\mathbf{w}^T \mathbf{S}_W \mathbf{w})$. Multiplying both sides of (4.29) by \mathbf{S}_W^{-1} we then obtain

$$\mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1). \quad (4.30)$$

Fisher's Linear Discriminant

Fisher's linear discriminant gives a specific choice of direction for projection of the data down to one dimension.

Projected data can subsequently be used to construct a discriminant, by choosing a threshold y_0 so that we classify a new point as belonging to C_1 if $y(x) \geq y_0$ and classify it as belonging to C_2 otherwise.

Perceptron Algorithm

For binary classification where the input vector \mathbf{x} is transformed using a nonlinear transformation to give a feature vector $\phi(\mathbf{x})$. We construct a **generalized linear model** is of the form

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$$

where the nonlinear activation function $f(\cdot)$ is given by a step function

$$f(a) = \begin{cases} 1 & \text{if } a \geq 0 \\ -1 & \text{if } a < 0. \end{cases} \quad (1)$$

Perceptron Algorithm

Find a weight vector \mathbf{w} such that for \mathbf{x}_n in class C_1 will have $\mathbf{w}^T \phi(\mathbf{x}_n) > 0$, whereas \mathbf{x}_n in class C_2 have $\mathbf{w}^T \phi(\mathbf{x}_n) < 0$.

The perceptron criterion associates zero error with any pattern that is correctly classified, whereas for a misclassified \mathbf{x}_n it tries to minimize the quantity $-\mathbf{w}^T \phi(\mathbf{x}_n) t_n$.

The perceptron criterion is given by

$$E(\mathbf{w}) = - \sum_{n \in M} \mathbf{w}^T \phi_n t_n$$

where M is the set of misclassified data.

Bayesian Probability Recap

Bayes' Theorem:

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$

Bayes' theorem can be used to convert a **prior probability** into a **posterior probability** by incorporating the evidence from the observed data.

Bayesian Probability Recap

Bayes' Theorem:

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$

Bayes' theorem can be used to convert a **prior probability** into a **posterior probability** by incorporating the evidence from the observed data.

Bayesian view of probability: probabilities provide a quantification of uncertainty.

We can use probability theory to describe the uncertainty in model parameters such as \mathbf{w} , or in the choice of model itself.

Bayesian Probability Recap

We capture our assumptions about the model parameter \mathbf{w} , before observing the data, in the form of a **prior probability distribution** $p(\mathbf{w})$.

The effect of the observed data $D = \{t_1, \dots, t_N\}$ is expressed through the conditional probability $p(D|\mathbf{w})$.

Bayes's theorem gives us

$$p(\mathbf{w}|D) = \frac{p(D|\mathbf{w})p(\mathbf{w})}{p(D)} \quad (\text{posterior} \propto \text{likelihood} \times \text{prior})$$

which allows us to evaluate the uncertainty in \mathbf{w} after we have observed D in the form of the **posterior probability** $p(\mathbf{w}|D)$.

$p(D|\mathbf{w})$ on the right hand side can be viewed as a function of \mathbf{w} , which is called the **likelihood function**. It expresses how probable the observed data set is for different settings of the parameter vector \mathbf{w} .