

Lecture 7: Linear Classification

Isabel Valera

Machine Learning Group
Department of Mathematics and Computer Science
Saarland University, Saarbrücken, Germany

10.05.2021

Outline

- 1 Bibliography
- 2 Introduction
- 3 LDA
- 4 Logistic regression
- 5 Summary
- 6 Appendix

Main references

- Bishop - Chapter 4
- ESL - Chapter 4

Outline

- 1 Bibliography
- 2 Introduction**
- 3 LDA
- 4 Logistic regression
- 5 Summary
- 6 Appendix

Classification setting

Binary classification setting: We consider problems where the output (target) variable takes values in $\mathcal{Y} = \{-1, +1\}$. Moreover, we assume we have access to training data $(X_i, Y_i)_{i=1}^n$, which is an i.i.d. sample from the probability measure P on $\mathcal{X} \times \mathcal{Y}$. Note that here we treat each observation (X_i, Y_i) in the training dataset as a random variable.

Bayes classifier, which decides according to $y^*(x) = \text{sign}(\mathbb{E}[Y|X = x])$, is optimal.

Goal: Learn a mapping function $\hat{y}(X)$ that minimizes the probability of error, i.e., $P(\text{error}) = R(\hat{y}) = \mathbb{E}[\mathbb{1}_{\hat{y}(X)Y \leq 0}]$. To this end, we often assume that $\hat{y}(X) = \text{sign}(f(X))$ and focus on learning (using training data) the discriminant function $f(X)$ that minimizes a surrogate loss that is convex and upperbounds the 0-1-loss (refer to Lecture 3).

Linear classification considers a family of functions $f \in \mathcal{F}$ that are linear, i.e., it takes the form $\left\{ \langle w, x \rangle + b, \text{ with } x, w \in \mathbb{R}^d, b \in \mathbb{R} \right\}$.

Linear discriminant function

Linear discriminant function: $f(x) = \{ \langle w, x \rangle + b \}$

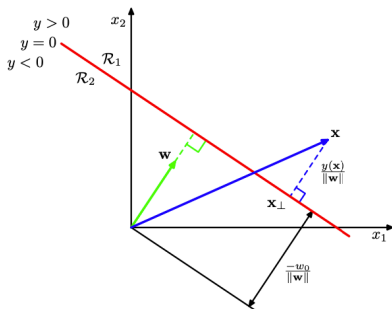


Figure: Example of linear discriminant function in $\mathcal{X} = \mathbb{R}^2$

- Any pair of points x_a and x_b lying on the decision surface satisfy that $f(x_a) = f(x_b) = 0$ and that $\langle w, x_a - x_b \rangle = 0$, as w is orthonormal to any vector lying in the decision surface.
- $f(x)$ gives a signed measure of the perpendicular distance $r = \frac{f(x)}{\|w\|}$ of the point x to the decision surface.
- Thus, we may use such a (hyper-)plane to divide \mathbb{R}^d into two regions \mathcal{R}_1 and \mathcal{R}_2 .

Linear Classification

Let $\mathcal{X} = \mathbb{R}^d$ be the input space, then a linear binary classifier $\hat{y} : \mathbb{R}^d \rightarrow \{-1, 1\}$ has the form

$$\hat{y}(x) = \text{sign}(f(x)) = \text{sign}(\langle w, x \rangle + b) = \begin{cases} 1 & \text{if } \langle w, x \rangle + b > 0, \\ -1 & \text{if } \langle w, x \rangle + b \leq 0. \end{cases}$$

Separation of the input space \mathbb{R}^d into two half spaces.

Observation: From now on we will include the bias term directly in the weight vector w .

Decision boundary

Decision boundary

Definition

Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a function and $\hat{y}(x) = \text{sign}(f(x))$ be the resulting classifier with output in $\mathcal{Y} = \{-1, 1\}$, then we call the set

$$\{x \in \mathcal{X} \mid f(x) = 0\},$$

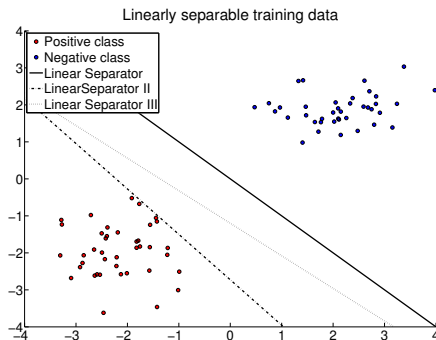
the **decision boundary** of the classifier \hat{y} .

A training set $D = (x_i, y_i)_{i=1}^n$ is **linearly separable** if there exists a weight vector w and an offset b such that,

$$y_i f(x_i) = y_i (\langle w, x_i \rangle) > 0, \quad \forall i = 1, \dots, n,$$

\Rightarrow There exists a **hyperplane** $\{x \in \mathbb{R}^d \mid \langle w, x \rangle = 0\}$ which separates the sets $X_+ = \{x_i \in D \mid y_i = 1\}$ and $X_- = \{x_i \in D \mid y_i = -1\}$.

Example



A training sample of a two-class problem in \mathbb{R}^2 . The two classes are linearly separable and three different decision hyperplanes are shown which separate the two classes. (Image by Prof. Hein)

Basis functions

As in linear regression we may apply basis functions to map the feature vectors $\mathcal{X} = \mathbb{R}^d$ into a possibly larger **feature space** \mathbb{R}^D , i.e.,

$$\mathbf{x} \in \mathbb{R}^d \longrightarrow (\phi_1(\mathbf{x}), \dots, \phi_D(\mathbf{x})),$$

such that we define the linear binary classifier as

$$\hat{y}(\mathbf{x}) = \text{sign}(f(\mathbf{x})) = \text{sign}(\langle \mathbf{w}, \Phi(\mathbf{x}) \rangle)$$

The discriminant function is linear in the parameters \mathbf{w} but not necessarily linear in the input space!

Observation: In the following we will make no distinction between using or not using basis functions Φ .

Methods for linear classification

Three linear methods: $\hat{y}(x) = \text{sign}(f(x)) = \text{sign}(\langle w, \Phi(x) \rangle)$.

- **Linear Discriminant Analysis:**

- Loss: Squared loss, $L(y, f(x)) = (y - f(x))^2$
- Regularization: none

- **Logistic Regression:**

- Loss: Logistic loss, $L(y, f(x)) = \log(1 + \exp(-y f(x)))$
- Regularization: usually none, but there exist regularized versions.

- **Support Vector Machines** (Lecture 14).

- Loss: hinge loss, $L(y, f(x)) = \max(0, 1 - y f(x))$
- Regularization: L2-regularization, i.e., $\Omega(w) = \|w\|_2^2$

All three methods construct a **linear** classifier but all three have different **objectives**.

Outline

- 1 Bibliography
- 2 Introduction
- 3 LDA**
- 4 Logistic regression
- 5 Summary
- 6 Appendix

Linear Discriminant Analysis (LDA)

Linear Discriminant Analysis (LDA):

- is often called **Fisher Discriminant Analysis** named after its inventor Ronald A. Fisher, the “father” of parametric statistics.
- projects the data $x \in \mathbb{R}^d$ into a lower dimensional space via the inner product with the weight vector, i.e., $\langle w, x \rangle$.
 - In binary classification, such a projection of the feature space \mathbb{R}^D onto the line $L = \{\alpha w \mid \alpha \in \mathbb{R}\}$.
 - Classification of the data by thresholding.

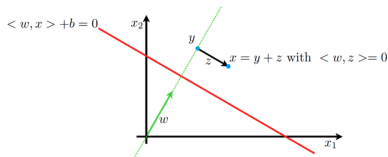


Figure: Projection onto the line
 $L = \{\alpha w \mid \alpha \in \mathbb{R}\}$

- $x = y + z$ such that
 $y \in L = \{\alpha w\}$ and $\langle w, z \rangle = 0$
- For each $y \in L(\alpha w)$ there exists
an α_0 such that $y = \alpha_0 \frac{w}{\|w\|}$ and
 $\|y\| = \alpha_0$.

LDA illustration

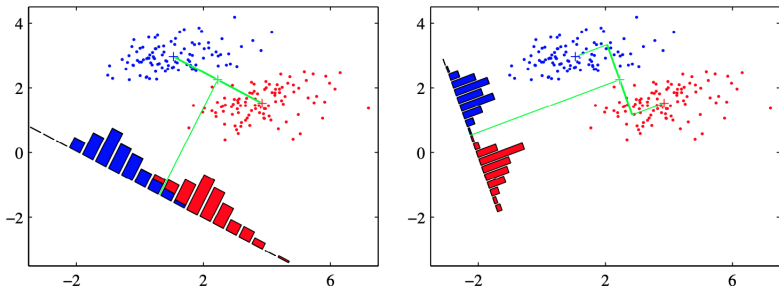


Figure 4.6 The left plot shows samples from two classes (depicted in red and blue) along with the histograms resulting from projection onto the line joining the class means. Note that there is considerable class overlap in the projected space. The right plot shows the corresponding projection based on the Fisher linear discriminant, showing the greatly improved class separation.

Figure: Image from Bishop.

Question: What is the best projection in the sense that it optimally separates the data?

Fisher criterion:

The **Fisher criterion** is defined as

$$J(w) = \frac{\langle w, \mu_+ - \mu_- \rangle^2}{\sigma_{w,+}^2 + \sigma_{w,-}^2},$$

which aims for:

- **Large distance** of the projected class centroids $\langle w, \mu_+ \rangle$ and $\langle w, \mu_- \rangle$. The class **centroids** μ_+ and μ_- of the positive and negative class are defined as:

$$\mu_+ = \frac{1}{n_+} \sum_{\{i \mid Y_i=1\}} X_i, \quad \mu_- = \frac{1}{n_-} \sum_{\{i \mid Y_i=-1\}} X_i,$$

- **Small variances** around the projected class centroids. The **within-class covariances** of the projections of the positive and negative class are given by:

$$\sigma_{w,+}^2 = \sum_{\{i \mid Y_i=1\}} \left(\langle w, X_i \rangle - \langle w, \mu_+ \rangle \right)^2, \quad \sigma_{w,-}^2 = \sum_{\{i \mid Y_i=-1\}} \left(\langle w, X_i \rangle - \langle w, \mu_- \rangle \right)^2.$$

Fisher criterion in matrix formulation

The **between-class covariance** matrix Σ_B is defined as

$$\Sigma_B = (\mu_+ - \mu_-)(\mu_+ - \mu_-)^T,$$

and the total **within-class covariance** matrix Σ_W as

$$\Sigma_W = \sum_{\{i \mid Y_i=1\}} (X_i - \mu_+)(X_i - \mu_+)^T + \sum_{\{i \mid Y_i=-1\}} (X_i - \mu_-)(X_i - \mu_-)^T.$$

Then the **Fisher criterion** $J(w)$ can be written as

$$J(w) = \frac{\langle w, \Sigma_B w \rangle}{\langle w, \Sigma_W w \rangle}.$$

LDA solution

Lemma

The optimal projection $w^* = \arg \max_{w \in \mathbb{R}^d} J(w)$ is given by

$$w^* = \Sigma_W^{-1}(\mu_+ - \mu_-).$$

Proof: We have

$$\nabla_w J(w) = 2 \frac{1}{\langle w, \Sigma_W w \rangle} \Sigma_B w - 2 \frac{\langle w, \Sigma_B w \rangle}{\langle w, \Sigma_W w \rangle^2} \Sigma_W w.$$

We solve $J(w) = 0$ and get

$$\frac{\langle w, \Sigma_W w \rangle}{\langle w, \Sigma_B w \rangle} \Sigma_B w = \Sigma_W w.$$

Now, $\Sigma_B w$ is always proportional to $\mu_+ - \mu_-$ and $\frac{\langle w, \Sigma_W w \rangle}{\langle w, \Sigma_B w \rangle}$ is just a scalar factor. Therefore,

$$w^* \propto \Sigma_W^{-1}(\mu_+ - \mu_-).$$

LDA classification

- **Final classifier:**

$$f(x) = \text{sign}(\langle w^*, x \rangle + b).$$

Determine **the threshold** b by minimizing the training error.

- Optimal Projection can be also derived using **least squares**, as

$$(w', w'_0) = \arg \min_{w \in \mathbb{R}^D, w_0 \in \mathbb{R}} \sum_{i=1}^n (Y_i - \langle w, X_i \rangle - w_0)^2.$$

One can prove that (*exercise!*):

$$w^* \sim w'.$$

Illustration

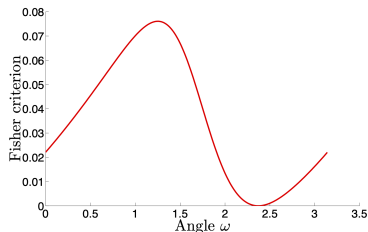
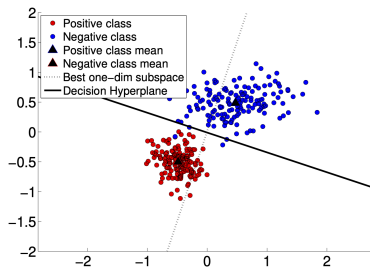


Figure: **Left:** Optimal projection line $\{\alpha w^* + \frac{1}{2}(\mu_+ + \mu_-) \mid \alpha \in \mathbb{R}\}$. **Right:** The Fisher criterion as a function of the angle ω , where ω is a parameterization of all weight vectors $w = (\cos(\omega), \sin(\omega))$ in \mathbb{R}^2 . (Image by Prof. Hein)

Generalization to the multi-class case I

The mean of all feature vectors is denoted by $\mu = \frac{1}{n} \sum_{i=1}^n X_i$, and as before the centroid for each class is denoted as μ_k .

The **between-class covariance** matrix Σ_B is defined as

$$\Sigma_B = \sum_{k=1}^K n_k (\mu_k - \mu)(\mu_k - \mu)^T.$$

and the total **within-class covariance** matrix Σ_W as

$$\Sigma_W = \sum_{k=1}^K \sum_{\{i \mid Y_i=k\}} (X_i - \mu_k)(X_i - \mu_k)^T,$$

The **Fisher criterion** $J(w)$ stays the same, i.e.,

$$J(w) = \frac{\langle w, \Sigma_B w \rangle}{\langle w, \Sigma_W w \rangle}.$$

LDA as Dimensionality Reduction

In general, we project the feature vectors $x \in \mathbb{R}^d$ into a new space $\mathbb{R}^{d'}$, where we often assume $d' = K - 1$.

Thus, multi-class LDA can be seen as a 'supervised' approach for **dimensionality reduction**, where we seek for:

- a lower dimensional $d' \ll d$ representation of the data, which preserves the “interesting” properties of the data
- a lower dimensional representation of the data in which the classifier performs as well as on the original d -dimensional space.

LDA as Dimensionality Reduction

How can we get $d' > 1$ projections from the Fisher criterion?

$$J(w) = \frac{\langle w, \Sigma_B w \rangle}{\langle w, \Sigma_W w \rangle}.$$

The solution is given by the following **generalized eigenvalue problem**:

$$\Sigma_B w = \lambda \Sigma_W w \quad (\text{If } \Sigma_W \text{ is invertible, then } \Sigma_W^{-1} \Sigma_B w = \lambda w).$$

m -dimensional projection is determined by the m eigenvectors corresponding to the m largest eigenvalues.

Note: Σ_B has rank $K - 1$ if class centroids are linearly independent. Thus, the (sorted) eigenvalues for $m > K - 1$ will be zero.

Observation: The above result can be explained by the generalized Rayleigh-Ritz principle (see the Appendix).

Observation: Further details in Chapter 4.3 of ESL book.

Observations

- In general, one needs generally a $K - 1$ -dimensional subspace in order to separate K classes!
- A linear projection of the data will in general lead to a worse Bayes risk (in particular if the data is not linearly separable).
- However, in high dimension the problem might be very difficult to solve (curse of dimensionality) and the hope is that by doing dimensionality reduction we can at least find a relatively good solution in this low-dimensional subspace.
- LDA suffers from overfitting when the number of training observations is on the same order as the number of features.

Outline

- 1 Bibliography
- 2 Introduction
- 3 LDA
- 4 Logistic regression**
- 5 Summary
- 6 Appendix

Logistic Regression

Idea: (Linear) Logistic regression models the conditional likelihood using a sigmoid function, i.e.,

$$p = P(Y = 1|X = x, w) = \frac{1}{1 + e^{-\langle w, \phi(x) \rangle}},$$

and finds the model parameters w using MLE.

Question: Why the sigmoid function? Because we can model the so-called **logistic function** (a.k.a. logit function) $\log(\frac{p}{1-p})$ using a linear function, i.e.,

$$\log\left(\frac{p}{1-p}\right) = \langle w, \phi(x) \rangle,$$

Note: We have included the bias term indirectly in the weight vector.

Proof.

$$\log\left(\frac{p}{1-p}\right) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle \quad (1)$$

$$\left(\frac{p}{1-p}\right) = \exp(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle) \quad (2)$$

$$p = (1-p) \exp(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle) \quad (3)$$

$$p = \frac{1}{(1 + \exp(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle))} \exp(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle) \quad (4)$$

$$p = \frac{\exp(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle)}{1 + \exp(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle)} = \frac{1}{1 + \exp(-\langle \mathbf{w}, \phi(\mathbf{x}) \rangle)}, \quad (5)$$

where we used that

$$1 - p = 1 - \frac{\exp(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle)}{1 + \exp(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle)} = \frac{1}{(1 + \exp(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle))}.$$

Illustration of sigmoid function

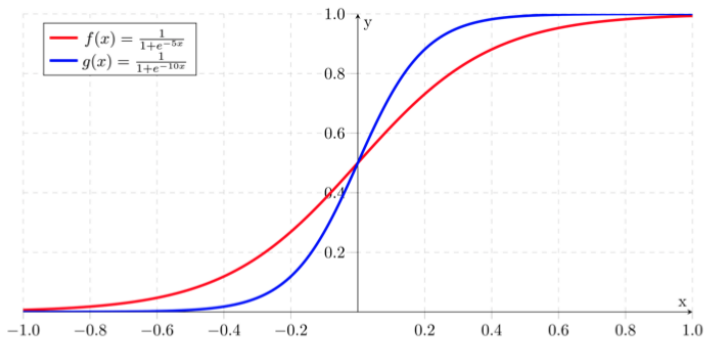


Figure: Illustration of sigmoid function (which is a special case of the logistic function).

Logistic Regression

For a given dataset $D_n = (x_i, y_i)_{i=1}^n$, we can obtain the optimal logistic regression parameters w by maximizing the likelihood:

$$\prod_{i=1}^n P(Y = y_i | X = x_i, w) = \prod_{i=1}^n \frac{1}{1 + e^{-y_i \langle w, \Phi(x_i) \rangle}}.$$

Definition (Logistic regression)

Given a training sample $D_n = (x_i, y_i)_{i=1}^n$ with $x_i \in \mathcal{X}$ and $y_i \in \{-1, 1\}$ and the function space $\mathcal{F} = \{\sum_{i=1}^D w_i \phi_i(x) \mid w \in \mathbb{R}^d\}$ we define **logistic regression** as the mapping

$$D_n \mapsto f_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \log \left(1 + \exp(-y_i \langle w, \Phi(x_i) \rangle) \right). \quad (6)$$

Observation: Logistic regression corresponds to **Empirical risk minimization using the logistic loss, i.e.,**

$L(y, f(x)) = -\log(1 + \exp(-y f(x)))$ where $f(x) = \langle w, \Phi(x) \rangle$.

Solving logistic regression

- Logistic regression has **no analytical solution**, but it can be written as a convex optimization problem w.r.t. the weights w .
- Thus, we can use Newton-type gradient descent method.
- The gradient and the Hessian of the empirical risk:

$$R_{\text{emp}}(w) = \frac{1}{n} \sum_{i=1}^n \log \left(1 + \exp(-y_i \langle w, \Phi(x_i) \rangle) \right),$$

are given by:

$$\nabla_w R_{\text{emp}}(w) = \frac{\partial R_{\text{emp}}}{\partial w_s}(w) = -\frac{1}{n} \sum_{i=1}^n y_i \phi_s(x_i) \frac{\exp(-y_i \langle w, \Phi(x_i) \rangle)}{1 + \exp(-y_i \langle w, \Phi(x_i) \rangle)},$$

$$H(R_{\text{emp}}) = \frac{\partial^2 R_{\text{emp}}}{\partial w_r \partial w_s}(w) = \frac{1}{n} \sum_{i=1}^n \phi_s(x_i) \phi_r(x_i) \frac{\exp(-y_i \langle w, \Phi(x_i) \rangle)}{\left(1 + \exp(-y_i \langle w, \Phi(x_i) \rangle) \right)^2}.$$

Newton-Raphson algorithm I

With stepsize fixed to 1, we can update the weights at each iteration of the algorithm as:

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} - \left(H(R_{\text{emp}}) \right)^{-1} \nabla_{\mathbf{w}} R_{\text{emp}}(\mathbf{w}),$$

With the diagonal matrices W and V with diagonal entries

$$W_{ii} = \frac{\exp(-y_i \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle)}{(1 + \exp(-y_i \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle))^2}, \quad V_{ii} = \frac{\exp(-y_i \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle)}{1 + \exp(-y_i \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle)},$$

we can write the gradient and Hessian $H(R_{\text{emp}})$ of R_{emp} as

$$\nabla_{\mathbf{w}} R_{\text{emp}}(\mathbf{w}) = -\frac{1}{n} \Phi^T V Y, \quad H(R_{\text{emp}})|_{\mathbf{w}} = \frac{1}{n} \Phi^T W \Phi.$$

Newton-Raphson algorithm II

Thus we can write the **Newton-Raphson update** as

$$\begin{aligned} \mathbf{w}_{\text{new}} &= \mathbf{w}_{\text{old}} + \left(\Phi^T \mathbf{W} \Phi \right)^{-1} \Phi^T \mathbf{V} \mathbf{Y} \\ &= \left(\Phi^T \mathbf{W} \Phi \right)^{-1} \Phi^T \mathbf{W} \left(\Phi \mathbf{w}_{\text{old}} + \mathbf{W}^{-1} \mathbf{V} \mathbf{Y} \right) = \left(\Phi^T \mathbf{W} \Phi \right)^{-1} \Phi^T \mathbf{W} \mathbf{Z}, \end{aligned}$$

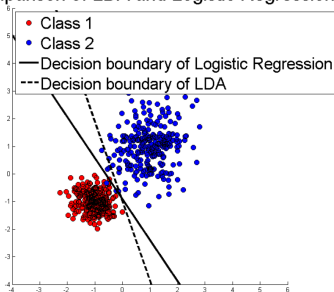
with $\mathbf{Z} = \Phi \mathbf{w}_{\text{old}} + \mathbf{W}^{-1} \mathbf{V} \mathbf{Y}$.

It can be seen as **iterative reweighted least squares problem!**

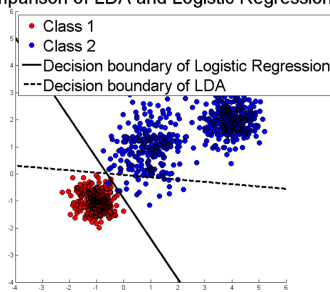
Observation: For further details, refer to Chapter 4.3.3 of Bishop.

Comparison LDA vs Logistic Regression

Comparison of LDA and Logistic Regression



Comparison of LDA and Logistic Regression



Left: Original data, **Right:** Adding the second Gaussian blob should not change the decision boundary. However, LDA changes its decision completely. (Image from Prof. Hein)

Regularized Logistic Regression

- As empirical risk minimization may suffer from overfitting, we can add a regularizer.
- For a linearly separable dataset the solution w^* is unbounded.
- Adding a regularizer on the weights, makes also the numerical solution more stable (Block III) since the involved matrices might be close to singular.

Definition

Given a training sample $D_n = (X_i, Y_i)_{i=1}^n$ with $X_i \in \mathcal{X}$ and $Y_i \in \{-1, 1\}$ and the function space $\mathcal{F} = \{\sum_{i=1}^d w_i \phi_i(x) \mid w \in \mathbb{R}^d\}$ we define **L_2 -regularized logistic regression** as the mapping:

$$D_n \mapsto f_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \log \left(1 + \exp(-y_i \langle w, \Phi(x_i) \rangle) \right) + \lambda \|w\|_2^2,$$

where λ is the regularization parameter.

Multi-class Logistic Regression

We can directly use the **soft-max loss**:

$$\begin{aligned} L(y, f(x)) &= -\log p(y|x, f) \\ &= -\log \left(\frac{\exp(\langle w_y, \Phi(x) \rangle)}{\sum_{k=1}^K \exp(\langle w_k, \Phi(x) \rangle)} \right), \end{aligned}$$

where now we have one weighting vector w_k for each class $k = 1, \dots, K$.

Classification: Once we have trained the classifier, i.e., learned the weight vectors, we assign a new data point x to the class that maximizes $\langle w_k, \Phi(x) \rangle$, which is equivalent to maximize our estimate of the conditional probability $\log p(y|x, f)$.

Outline

- 1 Bibliography
- 2 Introduction
- 3 LDA
- 4 Logistic regression
- 5 Summary**
- 6 Appendix

Summary

- **Linear Discriminant Analysis:**

- Loss: Squared loss, $L(y, f(x)) = (y - f(x))^2$

- **Logistic Regression:**

- Loss: Logistic loss, $L(y, f(x)) = \log(1 + \exp(-y f(x)))$

- In general, logistic regression is much more broadly used than LDA, as it leads to more robust solutions with respect to data perturbations. Least squares loss is influenced heavily by training data which lies far away from the decision boundary, whereas the logistic loss quickly decays far away from the decision boundary and is therefore only marginally influenced by new training data far away from the decision boundary.
- In high dimensions, logistic regression often competes with more complex classification approaches.
- **Support Vector Machines** will be introduced in Lecture 14.

Outline

- 1 Bibliography
- 2 Introduction
- 3 LDA
- 4 Logistic regression
- 5 Summary
- 6 Appendix**

Rayleigh-Ritz principle

Proposition (Rayleigh-Ritz principle)

Let $A \in \mathbb{R}^{d \times d}$ be a symmetric matrix, then

$$\lambda_{\max} = \max_{x \in \mathbb{R}^d} \frac{\langle x, Ax \rangle}{\langle x, x \rangle},$$

is the largest eigenvalue of A and the maximizing argument x_{\max} is the corresponding eigenvector. Equivalently,

$$\lambda_{\max} = \max_{x \in \mathbb{R}^d, \|x\|=1} \langle x, Ax \rangle.$$

Other eigenvalues and eigenvectors can be found as follows. Denote by u_1, \dots, u_r the eigenvectors corresponding to the largest r eigenvalues, then the $r+1$ largest eigenvalue can be found as,

$$\lambda_{r+1} = \max_{x \in \mathbb{R}^d, \langle x, u_s \rangle = 0, s=1, \dots, r} \frac{\langle x, Ax \rangle}{\langle x, x \rangle}.$$