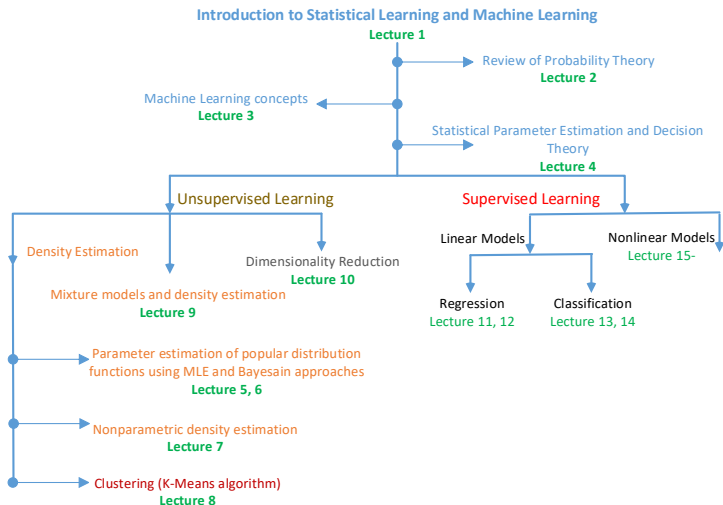


Statistical Learning and Machine Learning

Lecture 14 - Linear Models for Classification 2

October 10, 2021

Course overview and where do we stand



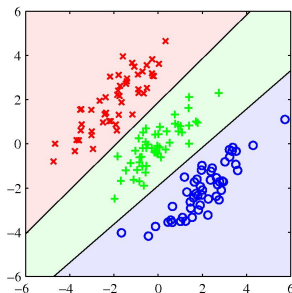
Objectives of the lecture

- Linear models for classification
 - Probabilistic generative models: linear discriminant analysis
 - Probabilistic discriminative models: Logistic regression

Linear classification

The **goal** of classification is to assign an input vector \mathbf{x} to one of the K classes \mathcal{C}_k , $k = 1, \dots, K$:

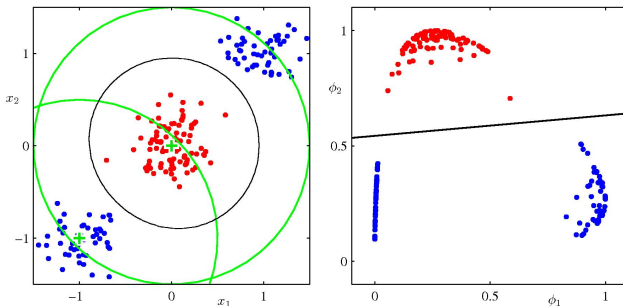
- The input space is divided into K **decision regions**, separated by **decision boundaries** or decision surfaces
- **Linear models**: are those for which decision boundaries are **linear functions** of \mathbf{x} i.e., $\mathbf{w}^T \mathbf{x} + w_0$



Decision boundaries for a linearly separable data set

Generalized linear classification

- **Generalized linear classification** corresponds to using linear combination of (nonlinearly) transformed vectors $\phi(\mathbf{x})$ i.e., $y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$.
- Appropriately chosen linear models of **basis functions** $\phi(\mathbf{x})$ can be equivalent to nonlinear models in \mathbf{x}



Three approaches to classification

- 1 Finding a function $y(\mathbf{x})$ called **discriminant function** which maps a new input \mathbf{x}_* onto a class label.
- 2 **Generative models:** Determining **class-conditional densities** $p(\mathbf{x}|\mathcal{C}_k)$ or joint distributions $p(\mathbf{x}, \mathcal{C}_k)$, followed by the estimation of posterior densities:

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)}$$

Finally, decision theory is used to determine class membership for new \mathbf{x} .

- 3 **Discriminative models:** Obtaining **posterior class probabilities** $p(\mathcal{C}_k|\mathbf{x})$ directly, followed by discrimination.

Probabilistic Generative Models

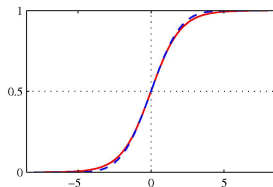
Consider a two-class problem; the posterior probability for \mathcal{C}_1 is

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\ &= \frac{1}{1 + \exp(-\alpha)} = \sigma(\alpha) \end{aligned} \quad (17)$$

where:

$$\alpha = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \quad (18)$$

and $\sigma(\alpha)$ is the *logistic sigmoid* function.



Logistic sigmoid function

Properties:

- **squashing function**: maps the whole real axis into a finite interval.
- $\sigma(-\alpha) = 1 - \sigma(\alpha)$
- $d\sigma/d\alpha = \sigma(1 - \sigma)$
- Inverse (known as *logit function*): $\alpha = \ln \left(\frac{\sigma}{1-\sigma} \right)$
- α represents the **log of the ratio of conditional probabilities for the two classes**:

$$\alpha(\mathbf{x}) = \ln \left(\frac{p(\mathcal{C}_1|\mathbf{x})}{p(\mathcal{C}_2|\mathbf{x})} \right) \quad (19)$$

also known as the *log odds*.

Softmax function

For $K > 2$ we obtain the **normalized exponential**:

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_{j=1}^K p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)} = \frac{\exp(\alpha_k)}{\sum_{j=1}^K \exp(\alpha_j)} \quad (20)$$

where:

$$\alpha_k = \ln p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k) \quad (21)$$

The normalized exponential (20) is also called **softmax** function:

- when $\alpha_k \gg \alpha_j$ for all $j \neq k$, then $p(\mathcal{C}_k|\mathbf{x}) \simeq 1$ and $p(\mathcal{C}_j|\mathbf{x}) \simeq 0$

Probabilistic Generative Models

Assuming that i) $K = 2$; ii) the class-conditional densities are Gaussian; and iii) there is a common covariance matrix $\Sigma_k = \Sigma$, $k = 1, 2$, the class-conditional density of \mathbf{x} is

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) \right\} \quad (22)$$

where D is the dimensionality of \mathbf{x} .

Substituting to the above expression for $p(\mathcal{C}_1|\mathbf{x})$:

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) \quad (23)$$

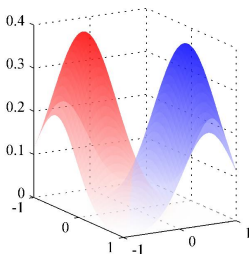
where:

$$\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \quad (24)$$

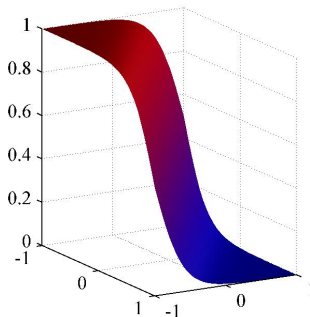
$$w_0 = -\frac{1}{2} \left(\boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 \right) + \ln \left(\frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)} \right) \quad (25)$$

Probabilistic Generative Models

- The resulting **decision boundaries** are linear in \mathbf{x} . **Why?**
- The prior probabilities $p(\mathcal{C}_k)$ enter through the bias parameter w_0 only, having the effect of making **parallel shifts of decision boundaries**.



Left: class-conditional densities for two classes (red and blue)
Right: Class posterior probability $p(\mathcal{C}_1|\mathbf{x})$



Probabilistic Generative Models

For $K > 2$ and again using a common covariance matrix $\Sigma_k = \Sigma$, $k = 1, \dots, K$:

$$\alpha_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (26)$$

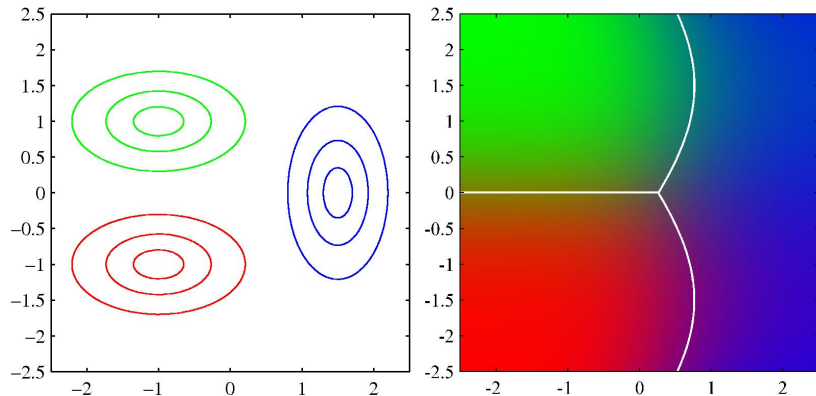
where:

$$\mathbf{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k \quad (27)$$

$$w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(\mathcal{C}_k). \quad (28)$$

In the above, using different Σ_k for each class will lead to quadratic functions ([discriminants](#)) of \mathbf{x} .

Probabilistic Generative Models



$$\Sigma_1 = \Sigma_2$$
$$\Sigma_3 \neq \Sigma_j, j = 1, 2.$$

Probabilistic Generative Models: Maximum Likelihood

- **Given:** Functional form for $p(\mathbf{x}, \mathcal{C}_k)$ together with training data $\{\mathbf{x}_n, t_n\}$ where $n = 1, \dots, N$. Here, $t_n = 1$ for class \mathcal{C}_1) and $t_n = 0$ for class \mathcal{C}_2).
- **Aim:** Estimating parameters of the $p(\mathbf{x}, \mathcal{C}_k)$ along with prior class probabilities $p(\mathcal{C}_k)$ via ML
- **Likelihood Function:**

$$\begin{aligned} p(\mathbf{t}, \mathbf{X} | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) &= \prod_{n=1}^N (p(\mathbf{x}_n, \mathcal{C}_1))^{t_n} (p(\mathbf{x}_n, \mathcal{C}_2))^{1-t_n} \\ &= \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma)]^{t_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma)]^{1-t_n} \end{aligned}$$

where $\mathbf{t} = (t_1, \dots, t_N)^T$.

Probabilistic Generative Models: Maximum Likelihood

- The ML estimate of μ_k , $k = 1, 2$ and Σ are:

$$\mu_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n \quad (29)$$

- The ML estimate of Σ is:

$$\Sigma = \mathbf{S} = \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2 \quad (30)$$

where

$$\mathbf{S}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mu_1)(\mathbf{x}_n - \mu_1)^T \quad (31)$$

$$\mathbf{S}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mu_1)(\mathbf{x}_n - \mu_1)^T \quad (32)$$

- The ML estimate of the prior class probabilities are

$$p(\mathcal{C}_1) = \pi = \frac{N_1}{N_1 + N_2}$$

Generative vs Discriminative Models

Until now we followed a **generative approach** to estimate posterior class probabilities:

$$\{p(\mathbf{x}|\mathcal{C}_k), p(\mathcal{C}_k)\} \Rightarrow p(\mathcal{C}_k|\mathbf{x}). \quad (33)$$

where we used the fact that

$$p(\mathcal{C}_k|\mathbf{x}) = \sigma(\mathbf{w}_k^T \mathbf{x} + w_0)$$

where: where:

$$\begin{aligned} \mathbf{w}_k &= \Sigma^{-1} \boldsymbol{\mu}_k \\ w_{k0} &= -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(\mathcal{C}_k). \end{aligned}$$

Linear discriminant functions:

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

Generative vs Discriminative Models

When we train **discriminative models**:

- we use a function form of $p(\mathcal{C}_k|\mathbf{x}) = \sigma(\mathbf{w}_k^T \mathbf{x} + w_0)$ explicitly without considering class conditional densities $p(\mathbf{x}|\mathcal{C}_k)$
- we optimize/estimate the parameters of the above function directly using training data.
- Pros and Cons?

Fixed basis functions

Sometimes it is beneficial to use a (fixed) nonlinear transformation of the data $\mathbf{x} \rightarrow \phi$ before applying the linear discriminative model.

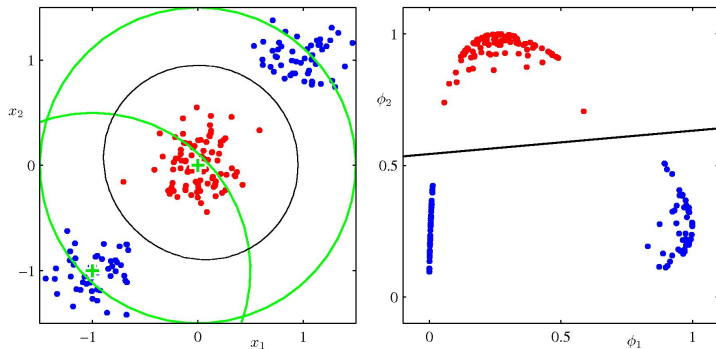


Illustration of the role of nonlinear basis in linear classification models. Use of 2 Gaussian basis functions centered at the green crosses converts the nonlinearly separable data (in \mathbf{x}) to linearly separable data in ϕ .

Logistic Regression

- Considering a 2-class classification problem, the posterior probability of class 1 can be written as a logistic sigmoid acting on a linear function of ϕ :

$$\begin{aligned}p(\mathcal{C}_1|\phi(\mathbf{x}_n)) &= y(\phi(\mathbf{x}_n)) = \sigma(\mathbf{w}^T \phi(\mathbf{x}_n)) \\p(\mathcal{C}_2|\phi(\mathbf{x}_n)) &= 1 - p(\mathcal{C}_1|\phi(\mathbf{x}_n))\end{aligned}$$

- Goal:** For a set of data $\{\phi(\mathbf{x}_n), t_n\}$, $n = 1, \dots, N$ with $t_n \in \{0, 1\}$ and $\mathbf{t} = [t_1, \dots, t_N]^T$, the goal is to estimate the optimal parameters \mathbf{w} using ML.
- Likelihood function:**

$$\begin{aligned}p(\mathbf{t}|\mathbf{w}) &= \prod_{n=1}^N (p(\mathcal{C}_1|\phi_n))^{t_n} (p(\mathcal{C}_2|\phi_n))^{1-t_n} \\&= \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}\end{aligned}$$

where $y_n = p(\mathcal{C}_1|\phi(\mathbf{x}_n))$

Logistic Regression

- We define the **error function** as the negative log-likelihood, which is called *cross-entropy* error function (suitable for minimization):

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \left(t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \right) \quad (34)$$

where $y_n = \sigma(\alpha_n)$ and $\alpha_n = \mathbf{w}^T \phi(\mathbf{x}_n)$.

- We optimize the parameters \mathbf{w} by applying **stochastic gradient descent**. But why not obtain \mathbf{w} directly as in LS?
- The gradient of error function w.r.t \mathbf{w} is

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi(\mathbf{x}_n) \quad (35)$$

Iterative Reweighted Least Squares

Reminder: Newton-Raphson method defines the update rule:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \mathbf{H}^{-1} \nabla E(\mathbf{w}^{(\tau)}) \quad (36)$$

where \mathbf{H} is the *Hessian matrix* having elements $H_{ij} = \frac{\partial^2 E_n(\mathbf{w})}{\partial w_i \partial w_j}$.

Using the cross-entropy error function:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi(\mathbf{x}_n) = \Phi^T (\mathbf{y} - \mathbf{t}) \quad (37)$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1 - y_n) \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T = \Phi^T \mathbf{R} \Phi, \quad (38)$$

where $\mathbf{R} = \text{diag}(y_n(1 - y_n))$ is a $N \times N$ square diagonal depending on \mathbf{w} .

Iterative Reweighted Least Squares

Using the property $0 < y_n < 1$:

- \mathbf{H} is positive definite: for any \mathbf{u} the value $\mathbf{u}^T \mathbf{H} \mathbf{u} > 0$
- the error function $E(\mathbf{w})$ is a *convex* function of \mathbf{w}
- the error function has a unique (global) minimum.

The **iterative update rule** (due to dependence of \mathbf{R} on \mathbf{w}) is:

$$\begin{aligned}\mathbf{w}^{(\tau+1)} &= \mathbf{w}^{(\tau)} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \left(\Phi^T \mathbf{R} \Phi \mathbf{w}^{(\tau)} - \Phi^T (\mathbf{y} - \mathbf{t}) \right) \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}\end{aligned}\tag{39}$$

where $\mathbf{z} \in \mathbb{R}^N$:

$$\mathbf{z} = \Phi \mathbf{w}^{(\tau)} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t}).\tag{40}$$

Multiclass logistic regression

Reminder:

$$p(\mathcal{C}_k | \phi(\mathbf{x}_n)) = y_k(\phi(\mathbf{x}_n)) = \frac{\exp(\alpha_k)}{\sum_{j=1}^K \exp(\alpha_j)} \quad (41)$$

where $\alpha_k = \mathbf{w}_k \phi(\mathbf{x}_n)$ and $\partial y_k / \partial \alpha_j = y_k(l_{kj} - y_j)$ with l_{kj} being the kj -th element of the identity matrix.

Using 1-of- K coding for the target vectors \mathbf{t}_n forming $\mathbf{T} \in \mathbb{R}^{N \times K}$, the **likelihood function** is:

$$p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K p(\mathcal{C}_k | \phi(\mathbf{x}_n))^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}} \quad (42)$$

with $y_{nk} = y_k(\phi(\mathbf{x}_n))$, and t_{nk} is the $\{n, k\}$ -th element of \mathbf{T} .

Multiclass logistic regression

The negative log-likelihood, which is called *cross-entropy* error function, (suitable for minimization) is:

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} \quad (43)$$

Derivatives of $E(\mathbf{w})$:

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi(\mathbf{x}_n) \quad (44)$$

$$\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T. \quad (45)$$

We can use $\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K)$ for applying SGD-based optimization, or both $\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K)$ and $\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K)$ to apply IRLS-based optimization.