

CS 688 – Pattern Recognition Lecture 4

Linear Models for Classification

- Probabilistic generative models
- Probabilistic discriminative models

Probabilistic Generative Models

- We now turn to a probabilistic approach to classification
- We'll see how models with linear decision boundaries arise from simple assumptions about the distribution of the data

Generative Approach

- Solve the inference problem of estimating the class-conditional densities $p(\mathbf{x} | C_k)$ for each class C_k
- Infer the prior class probabilities $p(C_k)$
- Use Bayes' theorem to find the class posterior probabilities:

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{p(\mathbf{x})}$$

where
$$p(\mathbf{x}) = \sum_k p(\mathbf{x} | C_k)p(C_k)$$

- Use decision theory to determine class membership for each new input \mathbf{x}

Probabilistic Generative Models

Two class case:

$$p(C_1 | \mathbf{x}) = \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_1)p(C_1) + p(\mathbf{x} | C_2)p(C_2)}$$

$$= \frac{1}{1 + e^{-a}} = \sigma(a)$$

$$a = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)}$$

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

logistic sigmoid function

Probabilistic Generative Models

K > 2 classes:

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{\sum_j p(\mathbf{x} | C_j)p(C_j)} = \frac{e^{a_k}}{\sum_j e^{a_j}}$$

$$a_k = \ln p(\mathbf{x} | C_k)p(C_k)$$

$$\sigma(a) = \frac{e^{a_k}}{\sum_j e^{a_j}} \quad \text{softmax function}$$

Probabilistic Generative Models

Lets assume the class-conditional densities are Gaussian with the same covariance matrix:

$$p(\mathbf{x} | C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} e^{\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_k)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_k)\right)}$$

Two class case first. We can show the following result:

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

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

$$w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}$$

Probabilistic Generative Models

We have shown:

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

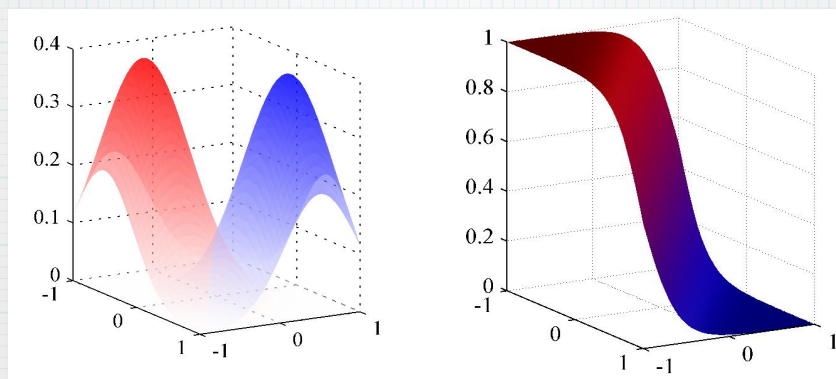
$$\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

$$w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}$$

Decision boundary: $p(C_1 | \mathbf{x}) = p(C_2 | \mathbf{x}) = 0.5$

$$\Rightarrow \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + w_0)}} = 0.5 \Rightarrow \mathbf{w}^T \mathbf{x} + w_0 = 0$$

Probabilistic Generative Models



Probabilistic Generative Models

K > 2 classes:

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{\sum_j p(\mathbf{x} | C_j)p(C_j)} = \frac{e^{a_k}}{\sum_j e^{a_j}}$$

$$a_k = \ln p(\mathbf{x} | C_k)p(C_k)$$

We can show the following result:

$$p(C_k | \mathbf{x}) = \frac{e^{\mathbf{w}_k^T \mathbf{x} + w_{k0}}}{\sum_j e^{\mathbf{w}_j^T \mathbf{x} + w_{j0}}}$$

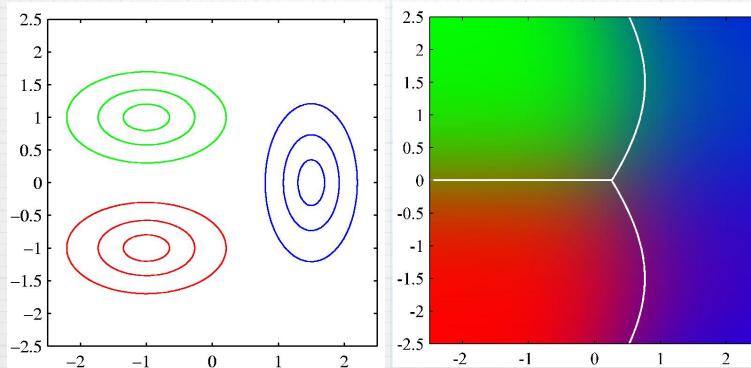
$$\mathbf{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k \quad w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(C_k)$$

Probabilistic Generative Models

The cancellation of the quadratic terms is due to the assumption of shared covariances.

If we allow each class conditional density to have its own covariance matrix, the cancellations no longer occur, and we obtain a quadratic function of \mathbf{x} , i.e. a **quadratic discriminant**.

Probabilistic Generative Models



Maximum likelihood solution

We have a parametric functional form for the class-conditional densities:

$$p(\mathbf{x} | C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} e^{\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_k)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_k)\right)}$$

We can estimate the parameters and the prior class probabilities using maximum likelihood.

➤ Two class case with shared covariance matrix.

➤ Training data:

$$\{\mathbf{x}_n, t_n\} \quad n = 1, \dots, N$$

$t_n = 1$ denotes class C_1 $t_n = 0$ denotes class C_2

Priors : $p(C_1) = \pi$, $p(C_2) = (1 - \pi)$

Maximum likelihood solution

$$\{\mathbf{x}_n, t_n\} \quad n = 1, \dots, N$$

$t_n = 1$ denotes class C_1 $t_n = 0$ denotes class C_2

Priors : $p(C_1) = \pi, \quad p(C_2) = (1 - \pi)$

For a data point \mathbf{x}_n from class C_1 we have $t_n = 1$ and therefore

$$p(\mathbf{x}_n, C_1) = p(C_1)p(\mathbf{x}_n | C_1) = \pi N(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})$$

For a data point \mathbf{x}_n from class C_2 we have $t_n = 0$ and therefore

$$p(\mathbf{x}_n, C_2) = p(C_2)p(\mathbf{x}_n | C_2) = (1 - \pi)N(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})$$

Maximum likelihood solution

$$\{\mathbf{x}_n, t_n\} \quad n = 1, \dots, N$$

$t_n = 1$ denotes class C_1 $t_n = 0$ denotes class C_2

Priors : $p(C_1) = \pi, \quad p(C_2) = (1 - \pi)$

$$p(\mathbf{x}_n, C_1) = p(C_1)p(\mathbf{x}_n | C_1) = \pi N(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})$$

$$p(\mathbf{x}_n, C_2) = p(C_2)p(\mathbf{x}_n | C_2) = (1 - \pi)N(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})$$

Assuming observations are drawn independently, we can write the likelihood function as follows:

$$p(\mathbf{t} | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{n=1}^N p(t_n | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) =$$

$$\prod_{n=1}^N [p(\mathbf{x}_n, C_1)]^{t_n} [p(\mathbf{x}_n, C_2)]^{1-t_n} =$$

$$\prod_{n=1}^N [\pi N(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1 - \pi)N(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n} \quad \mathbf{t} = (t_1, \dots, t_N)^T$$

Maximum likelihood solution

We want to find the values of the parameters that maximize the likelihood function, i.e., fit a model that best describes the observed data.

$$p(\mathbf{t} | \pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N [\pi N(\mathbf{x}_n | \mu_1, \Sigma)]^{t_n} [(1 - \pi) N(\mathbf{x}_n | \mu_2, \Sigma)]^{1-t_n}$$

As usual, we consider the log of the likelihood:

$$\ln p(\mathbf{t} | \pi, \mu_1, \mu_2, \Sigma) = \sum_{n=1}^N [t_n \ln \pi + t_n \ln N(\mathbf{x}_n | \mu_1, \Sigma) + (1 - t_n) \ln(1 - \pi) + (1 - t_n) \ln N(\mathbf{x}_n | \mu_2, \Sigma)]$$

Maximum likelihood solution

$$\ln p(\mathbf{t} | \pi, \mu_1, \mu_2, \Sigma) = \sum_{n=1}^N [t_n \ln \pi + t_n \ln N(\mathbf{x}_n | \mu_1, \Sigma) + (1 - t_n) \ln(1 - \pi) + (1 - t_n) \ln N(\mathbf{x}_n | \mu_2, \Sigma)]$$

We first maximize the log likelihood with respect to π . The terms that depend on π are $\sum_{n=1}^N [t_n \ln \pi + (1 - t_n) \ln(1 - \pi)]$

$$\begin{aligned} \frac{\partial}{\partial \pi} \sum_{n=1}^N [t_n \ln \pi + (1 - t_n) \ln(1 - \pi)] &= \frac{1}{\pi} \sum_{n=1}^N t_n - \frac{1}{1 - \pi} \sum_{n=1}^N (1 - t_n) \\ &= \frac{1}{\pi} \sum_{n=1}^N t_n - \frac{1}{1 - \pi} N + \frac{1}{1 - \pi} \sum_{n=1}^N t_n = \frac{1}{\pi(1 - \pi)} \sum_{n=1}^N t_n - \frac{1}{1 - \pi} N = 0 \\ \Rightarrow \pi &= \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2} \end{aligned}$$

Maximum likelihood solution

$$\pi_{ML} = \frac{N_1}{N_1 + N_2}$$

Thus, the maximum likelihood estimate of π is the fraction of points in class C_1

The result can be generalized to the multiclass case: the maximum likelihood estimate of $p(C_k)$ is given by the fraction of points in the training set that belong to C_k

Maximum likelihood solution

$$\ln p(\mathbf{t} | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \sum_{n=1}^N [t_n \ln \pi + t_n \ln N(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) + (1-t_n) \ln(1-\pi) + (1-t_n) \ln N(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]$$

We now maximize the log likelihood with respect to $\boldsymbol{\mu}_1$. The terms that depend on $\boldsymbol{\mu}_1$ are $\sum_{n=1}^N [t_n \ln N(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]$

$$\frac{\partial}{\partial \boldsymbol{\mu}_1} \sum_{n=1}^N [t_n \ln N(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})] = \frac{\partial}{\partial \boldsymbol{\mu}_1} \left[-\frac{1}{2} \sum_{n=1}^N [t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1)] + \text{const} \right]$$

$$= \sum_{n=1}^N [t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}] = 0 \Rightarrow \sum_{n=1}^N [t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)] = 0 \Rightarrow \sum_{n=1}^N t_n \mathbf{x}_n = N_1 \boldsymbol{\mu}_1$$

$$\Rightarrow \boldsymbol{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n$$

$$[N(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} e^{\left(-\frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1)\right)}]$$

Maximum likelihood solution

$$\boldsymbol{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n$$

Thus, the maximum likelihood estimate of $\boldsymbol{\mu}_1$ is the sample mean of all the input vectors \mathbf{x}_n assigned to class C_1

By maximizing the log likelihood with respect to $\boldsymbol{\mu}_2$ we obtain a similar result for $\boldsymbol{\mu}_2$

$$\boldsymbol{\mu}_2 = \frac{1}{N_2} \sum_{n=1}^N (1-t_n) \mathbf{x}_n$$

Maximum likelihood solution

Maximizing the log likelihood with respect to $\boldsymbol{\Sigma}$ we obtain the maximum likelihood estimate $\boldsymbol{\Sigma}_{ML}$

$$\boldsymbol{\Sigma}_{ML} = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2$$

$$S_1 = \frac{1}{N_1} \sum_{n \in C_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)(\mathbf{x}_n - \boldsymbol{\mu}_1)^T \quad S_2 = \frac{1}{N_2} \sum_{n \in C_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)(\mathbf{x}_n - \boldsymbol{\mu}_2)^T$$

➤ Thus: the maximum likelihood estimate of the covariance is given by the weighted average of the sample covariance matrices associated with each of the classes.

➤ This results extend to K classes.

Probabilistic Discriminative Models

Two-class case: $p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$

Multiclass case: $p(C_k | \mathbf{x}) = \frac{e^{\mathbf{w}_k^T \mathbf{x} + w_{k0}}}{\sum_j e^{\mathbf{w}_j^T \mathbf{x} + w_{j0}}}$

Discriminative approach: use the functional form of the generalized linear model for the posterior probabilities and determine its parameters directly using maximum likelihood.

Probabilistic Discriminative Models

Advantages:

- Fewer parameters to be determined
- Improved predictive performance, especially when the class-conditional density assumptions give a poor approximation of the true distributions.

Probabilistic Discriminative Models

Two-class case:

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) = y(\mathbf{x})$$

$$p(C_2 | \mathbf{x}) = 1 - p(C_1 | \mathbf{x})$$

In the terminology of statistics, this model is known as logistic regression.

Assuming $\mathbf{x} \in \mathfrak{R}^M$ how many parameters do we need to estimate?

$$M + 1$$

Probabilistic Discriminative Models

How many parameters did we estimate to fit Gaussian class-conditional densities (generative approach)?

$$p(C_1) \Rightarrow 1$$

$$2 \text{ mean vectors} \Rightarrow 2M$$

$$\Sigma \Rightarrow M + \frac{M^2 - M}{2} = \frac{M^2 + M}{2}$$

$$\text{total} = 1 + 2M + \frac{M^2 + M}{2} = O(M^2)$$

Logistic Regression

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) = y(\mathbf{x})$$

We use maximum likelihood to determine the parameters of the logistic regression model.

$$\{\mathbf{x}_n, t_n\} \quad n = 1, \dots, N$$

$t_n = 1$ denotes class C_1 $t_n = 0$ denotes class C_2

We want to find the values of \mathbf{w} that maximize the posterior probabilities associated to the observed data

Likelihood function :

$$L(\mathbf{w}) = \prod_{n=1}^N P(C_1 | \mathbf{x}_n)^{t_n} (1 - P(C_1 | \mathbf{x}_n))^{1-t_n} = \prod_{n=1}^N y(\mathbf{x}_n)^{t_n} (1 - y(\mathbf{x}_n))^{1-t_n}$$

Logistic Regression

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) = y(\mathbf{x})$$

$$L(\mathbf{w}) = \prod_{n=1}^N P(C_1 | \mathbf{x}_n)^{t_n} (1 - P(C_1 | \mathbf{x}_n))^{1-t_n} = \prod_{n=1}^N y(\mathbf{x}_n)^{t_n} (1 - y(\mathbf{x}_n))^{1-t_n}$$

We consider the negative logarithm of the likelihood:

$$E(\mathbf{w}) = -\ln L(\mathbf{w}) = -\ln \prod_{n=1}^N y(\mathbf{x}_n)^{t_n} (1 - y(\mathbf{x}_n))^{1-t_n} =$$

$$-\sum_{n=1}^N (t_n \ln y(\mathbf{x}_n) + (1 - t_n) \ln(1 - y(\mathbf{x}_n)))$$

$$\arg \min_{\mathbf{w}} E(\mathbf{w})$$

Logistic Regression

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) = y(\mathbf{x})$$

We compute the derivative of the error function with respect to \mathbf{w} (gradient):

$$\frac{\partial}{\partial \mathbf{w}} E(\mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} \left[- \sum_{n=1}^N (t_n \ln y(\mathbf{x}_n) + (1-t_n) \ln(1-y(\mathbf{x}_n))) \right]$$

We need to compute the derivative of the logistic sigmoid function:

$$\begin{aligned} \frac{\partial}{\partial a} \sigma(a) &= \frac{\partial}{\partial a} \frac{1}{1+e^{-a}} = \frac{e^{-a}}{(1+e^{-a})^2} = \frac{1}{1+e^{-a}} \frac{e^{-a}}{1+e^{-a}} = \\ &= \frac{1}{1+e^{-a}} \left(1 - \frac{1}{1+e^{-a}} \right) = \sigma(a)(1-\sigma(a)) \end{aligned}$$

Logistic Regression

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) = y(\mathbf{x})$$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} E(\mathbf{w}) &= \frac{\partial}{\partial \mathbf{w}} \left[- \sum_{n=1}^N (t_n \ln y(\mathbf{x}_n) + (1-t_n) \ln(1-y(\mathbf{x}_n))) \right] = \\ &= - \sum_{n=1}^N \left(\frac{t_n}{y_n} y_n (1-y_n) \mathbf{x}_n + \frac{(1-t_n)}{(1-y_n)} (-y_n)(1-y_n) \mathbf{x}_n \right) = \\ &= - \sum_{n=1}^N (t_n (1-y_n) \mathbf{x}_n - (1-t_n) y_n \mathbf{x}_n) = - \sum_{n=1}^N (t_n - t_n y_n - y_n + t_n y_n) \mathbf{x}_n = \\ &= - \sum_{n=1}^N (t_n - y_n) \mathbf{x}_n = \sum_{n=1}^N (y_n - t_n) \mathbf{x}_n \\ \nabla E(\mathbf{w}) &= \sum_{n=1}^N (y_n - t_n) \mathbf{x}_n \end{aligned}$$

Logistic Regression

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

- The gradient of E at \mathbf{w} gives the direction of the steepest increase of E at \mathbf{w} . We need to minimize E . Thus we need to update \mathbf{w} so that we move along the opposite direction of the gradient: $\mathbf{w}^{t+1} - \mathbf{w}^t \propto -\nabla E(\mathbf{w})$

This technique is called **gradient descent**

- It can be shown that E is a concave function of \mathbf{w} . Thus, it has a unique minimum.
- An efficient iterative technique exists to find the optimal \mathbf{w} parameters (*Newton-Raphson optimization*).

Batch vs. on-line learning

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

- The computation of the above gradient requires the processing of the entire training set (*batch technique*)

$$\mathbf{w}^{t+1} - \mathbf{w}^t \propto -\nabla E(\mathbf{w})$$

- If the data set is large, the above technique can be costly;
- For real time applications in which data become available as continuous streams, we may want to update the parameters as data points are presented to us (*on-line technique*).

On-line learning

- After the presentation of each data point n , we compute the contribution of that data point to the gradient (stochastic gradient):

$$\nabla E_n(\mathbf{w}) = (y_n - t_n)\mathbf{x}_n$$

- The on-line updating rule for the parameters becomes:

$$\mathbf{w}^{t+1} - \mathbf{w}^t \propto -\nabla E_n(\mathbf{w})$$

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla E_n(\mathbf{w}) = \mathbf{w}^t - \eta (y_n - t_n)\mathbf{x}_n$$

$\eta > 0$ is called learning rate.

It's value needs to be chosen carefully to ensure convergence

Multiclass Logistic Regression

Multiclass case:
$$p(C_k | \mathbf{x}) = \frac{e^{\mathbf{w}_k^T \mathbf{x} + w_{k0}}}{\sum_j e^{\mathbf{w}_j^T \mathbf{x} + w_{j0}}} = y_k(\mathbf{x})$$

We use maximum likelihood to determine the parameters of the logistic regression model.

$$\{\mathbf{x}_n, \mathbf{t}_n\} \quad n = 1, \dots, N$$

$$\mathbf{t}_n = (0, \dots, 1, \dots, 0) \text{ denotes class } C_k$$

We want to find the values of $\mathbf{w}_1, \dots, \mathbf{w}_K$ that maximize the posterior probabilities associated to the observed data

Likelihood function :

$$L(\mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K P(C_k | \mathbf{x}_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_k(\mathbf{x}_n)^{t_{nk}}$$

Multiclass Logistic Regression

$$p(C_k | \mathbf{x}) = \frac{e^{\mathbf{w}_k^T \mathbf{x} + w_{k0}}}{\sum_j e^{\mathbf{w}_j^T \mathbf{x} + w_{j0}}} = y_k(\mathbf{x})$$

$$L(\mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K y_k(\mathbf{x}_n)^{t_{nk}}$$

We consider the negative logarithm of the likelihood:

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln L(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_k(\mathbf{x}_n)$$

$$\arg \min_{\mathbf{w}_j} E(\mathbf{w}_j)$$

Multiclass Logistic Regression

$$p(C_k | \mathbf{x}) = \frac{e^{\mathbf{w}_k^T \mathbf{x} + w_{k0}}}{\sum_j e^{\mathbf{w}_j^T \mathbf{x} + w_{j0}}} = y_k(\mathbf{x})$$

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_k(\mathbf{x}_n)$$

We compute the gradient of the error function with respect to one of the parameter vectors:

$$\frac{\partial}{\partial \mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \frac{\partial}{\partial \mathbf{w}_j} \left[-\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_k(\mathbf{x}_n) \right]$$

Multiclass Logistic Regression

$$p(C_k | \mathbf{x}) = \frac{e^{\mathbf{w}_k^T \mathbf{x} + w_{k0}}}{\sum_j e^{\mathbf{w}_j^T \mathbf{x} + w_{j0}}} = y_k(\mathbf{x})$$

$$\frac{\partial}{\partial \mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \frac{\partial}{\partial \mathbf{w}_j} \left[- \sum_{n=1}^N \sum_{k=1}^K t_{nk} y_k(\mathbf{x}_n) \right]$$

Thus, we need to compute the derivatives of the softmax function:

$$\frac{\partial}{\partial a_k} y_k = \frac{\partial}{\partial a_k} \frac{e^{a_k}}{\sum_j e^{a_j}} = \frac{e^{a_k} \sum_j e^{a_j} - e^{a_k} e^{a_k}}{\left(\sum_j e^{a_j} \right)^2} = y_k - y_k^2 = y_k(1 - y_k)$$

Multiclass Logistic Regression

$$p(C_k | \mathbf{x}) = \frac{e^{\mathbf{w}_k^T \mathbf{x} + w_{k0}}}{\sum_j e^{\mathbf{w}_j^T \mathbf{x} + w_{j0}}} = y_k(\mathbf{x})$$

$$\frac{\partial}{\partial \mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \frac{\partial}{\partial \mathbf{w}_j} \left[- \sum_{n=1}^N \sum_{k=1}^K t_{nk} y_k(\mathbf{x}_n) \right]$$

Thus, we need to compute the derivatives of the softmax function:

$$\text{for } j \neq k, \quad \frac{\partial}{\partial a_j} y_k = \frac{\partial}{\partial a_j} \frac{e^{a_k}}{\sum_j e^{a_j}} = \frac{-e^{a_k} e^{a_j}}{\left(\sum_j e^{a_j} \right)^2} = -y_k y_j$$

Multiclass Logistic Regression

$$\frac{\partial}{\partial a_k} y_k = \frac{\partial}{\partial a_k} \frac{e^{a_k}}{\sum_j e^{a_j}} = \frac{e^{a_k} \sum_j e^{a_j} - e^{a_k} e^{a_k}}{\left(\sum_j e^{a_j}\right)^2} = y_k - y_k^2 = y_k(1 - y_k)$$

$$\text{for } j \neq k, \frac{\partial}{\partial a_j} y_k = \frac{\partial}{\partial a_j} \frac{e^{a_k}}{\sum_j e^{a_j}} = \frac{-e^{a_k} e^{a_j}}{\left(\sum_j e^{a_j}\right)^2} = -y_k y_j$$

Compact expression:

$$\frac{\partial}{\partial a_j} y_k = y_k (I_{kj} - y_j)$$

where I_{kj} are the elements of the identity matrix

Multiclass Logistic Regression

$$p(C_k | \mathbf{x}) = \frac{e^{\mathbf{w}_k^T \mathbf{x} + w_{k0}}}{\sum_j e^{\mathbf{w}_j^T \mathbf{x} + w_{j0}}} = y_k(\mathbf{x})$$

$$\frac{\partial}{\partial a_j} y_k = y_k (I_{kj} - y_j)$$

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \frac{\partial}{\partial \mathbf{w}_j} \left[- \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_k(\mathbf{x}_n) \right] =$$

$$- \sum_{n=1}^N \sum_{k=1}^K t_{nk} \frac{1}{y_{nk}} y_{nk} (I_{kj} - y_{nj}) \mathbf{x}_n = \sum_{n=1}^N \sum_{k=1}^K (t_{nk} y_{nj} - t_{nk} I_{kj}) \mathbf{x}_n =$$

$$\sum_{n=1}^N (y_{nj} - t_{nj}) \mathbf{x}_n$$

Multiclass Logistic Regression

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

- It can be shown that E is a concave function of \mathbf{w} . Thus, it has a unique minimum.
- For a batch solution, we can use the *Newton-Raphson optimization* technique.
- On-line solution (stochastic gradient descent):

$$\mathbf{w}_j^{t+1} = \mathbf{w}_j^t - \eta \nabla E_n(\mathbf{w}) = \mathbf{w}_j^t - \eta (y_{nj} - t_{nj}) \mathbf{x}_n$$