

Machine Learning Exercise Sheet 05

Linear Classification

Exercise sheets consist of two parts: homework and in-class exercises. You solve the homework exercises on your own or with your registered group and upload it to Moodle for a possible grade bonus. The in-class exercises will be solved and explained during the tutorial. You do not have to upload any solutions of the in-class exercises.

In-class Exercises

Multi-Class Classification

Problem 1: Consider a generative classification model for C classes defined by class probabilities $p(y = c) = \pi_c$ and generative conditional densities $p(\mathbf{x} | y = c, \theta_c)$ where $\mathbf{x} \in \mathbb{R}^D$ is the input feature vector and $\theta = \{\theta_c\}_{c=1}^C$ are further model parameters. Suppose we are given a training set $\mathcal{D} = \{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}_{n=1}^N$ where $\mathbf{y}^{(n)}$ is a binary target vector of length C that uses the 1-of- C (one-hot) encoding scheme, so that it has components $y_c^{(n)} = \delta_{ck}$ if pattern n is from class $y = k$. Assuming that the data points are i.i.d., show that the maximum-likelihood solution for the class probabilities π is given by

$$\pi_c = \frac{N_c}{N}$$

where N_c is the number of data points assigned to class c .

The data likelihood given the parameters $\{\pi_c, \theta_c\}_{c=1}^C$ is

$$p(\mathcal{D} | \{\pi_c, \theta_c\}_{c=1}^C) = \prod_{n=1}^N \prod_{c=1}^C (p(\mathbf{x}^{(n)} | \theta_c) \pi_c)^{y_c^{(n)}}$$

and so the data log-likelihood is given by

$$\log p(\mathcal{D} | \{\pi_c, \theta_c\}_{c=1}^C) = \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \pi_c + \text{const w.r.t. } \pi_c.$$

In order to maximize the log likelihood with respect to π_c we need to preserve the constraint $\sum_c \pi_c = 1$. For this we use the method of Lagrange multipliers where we introduce λ as an unconstrained additional parameter and find a local extremum of the unconstrained function

$$\sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \pi_c - \lambda \left(\sum_{c=1}^C \pi_c - 1 \right).$$

instead. See wikipedia article on Lagrange multipliers for an intuition of why this works. This function is a sum of concave terms in π_c as well as λ and is therefore itself concave in these variables.

We can find the extremum by finding the root of the derivatives. Setting the derivative with respect to π_c equal to zero, we obtain

$$\pi_c = \frac{1}{\lambda} \sum_{n=1}^N y_c^{(n)} = \frac{N_c}{\lambda}.$$

Setting the derivative with respect to λ equal to zero, we obtain the original constraint

$$\sum_{c=1}^C \pi_c = 1$$

where we can now plug in the previous result $\pi_c = \frac{N_c}{\lambda}$ and obtain $\lambda = \sum_c N_c = N$. Plugging this in turn into the expression for π_c we obtain

$$\pi_c = \frac{N_c}{N}$$

which we wanted to show.

Linear Discriminant Analysis

Problem 2: Using the same classification model as in the previous question, now suppose that the class-conditional densities are given by Gaussian distributions with a *shared* covariance matrix, so that

$$p(\mathbf{x}^{(n)} | y = c, \boldsymbol{\theta}) = p(\mathbf{x}^{(n)} | \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}^{(n)} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}).$$

Show that the maximum likelihood estimate for the mean of the Gaussian distribution for class c is given by

$$\boldsymbol{\mu}_c = \frac{1}{N_c} \sum_{\substack{n=1 \\ y^{(n)}=c}}^N \mathbf{x}^{(n)}$$

which represents the mean of the observations assigned to class c .

Similarly, show that the maximum likelihood estimate for the shared covariance matrix is given by

$$\boldsymbol{\Sigma} = \sum_{c=1}^C \frac{N_c}{N} \mathbf{S}_c \quad \text{where} \quad \mathbf{S}_c = \frac{1}{N_c} \sum_{\substack{n=1 \\ y^{(n)}=c}}^N (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)(\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)^T.$$

Thus $\boldsymbol{\Sigma}$ is given by a weighted average of the sample covariances of the data associated with each class, in which the weighting coefficients N_c/N are the prior probabilities of the classes.

We begin by writing out the data log-likelihood.

$$\begin{aligned} & \log p(\mathcal{D} | \{\pi_c, \boldsymbol{\theta}_c\}_{c=1}^C) \\ &= \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \pi_c \cdot p(\mathbf{x}^{(n)} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}) \end{aligned}$$

Then we plug in the definition of the multivariate Gaussian

$$= \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \left((2\pi)^{-\frac{D}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp \left(-\frac{1}{2} (x^{(n)} - \mu_c)^T \Sigma^{-1} (x^{(n)} - \mu_c) \right) \right) + y^{(n)} \log \pi_c$$

and simplify.

$$= -\frac{1}{2} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \left(D \log 2\pi + \log \det(\Sigma) + (x^{(n)} - \mu_c)^T \Sigma^{-1} (x^{(n)} - \mu_c) - 2 \log \pi_c \right)$$

This expression is concave in μ_c , so we can obtain the maximizer by finding the root of the derivative. With the help of the matrix cookbook, we identify the derivative with respect to μ_c as

$$\sum_{n=1}^N y_c^{(n)} \Sigma^{-1} (x^{(n)} - \mu_c)$$

which we can set to 0 and solve for μ_c to obtain

$$\mu_c = \frac{1}{\sum_{n=1}^N y_c^{(n)}} \sum_{n=1}^N y_c^{(n)} x^{(n)} = \frac{1}{N} \sum_{\substack{n=1 \\ y^{(n)}=c}}^N x^{(n)}.$$

To find the optimal Σ , we need the trace trick

$$a = \text{Tr}(a) \text{ for all } a \in \mathbb{R}^{n \times n} \text{ and } \text{Tr}(AB) = \text{Tr}(BCA).$$

With this we can rewrite

$$(x^{(n)} - \mu_c)^T \Sigma^{-1} (x^{(n)} - \mu_c) = \text{Tr} \left(\Sigma^{-1} (x^{(n)} - \mu_c)(x^{(n)} - \mu_c)^T \right)$$

and use the matrix-trace derivative rule $\frac{\partial}{\partial A} \text{Tr}(AB) = B^T$ to find the derivative of the data log-likelihood with respect to Σ . Because the log-likelihood contains both Σ and Σ^{-1} , we convert one into the other with $\log \det A = -\log \det A^{-1}$ to obtain

$$-\frac{1}{2} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \left(-\log \det \Sigma^{-1} + \text{Tr} \left(\Sigma^{-1} (x^{(n)} - \mu_c)(x^{(n)} - \mu_c)^T \right) \right) + \text{const w.r.t. } \Sigma.$$

Finally, we use rule (57) from the matrix cookbook $\frac{\partial \log |\det X|}{\partial X} = (X^{-1})^T$ and compute the derivative of the log-likelihood with respect to Σ^{-1} as

$$-\frac{1}{2} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \left(-\Sigma^T + (x^{(n)} - \mu_c)(x^{(n)} - \mu_c)^T \right).$$

We find the root with respect to Σ and find

$$\Sigma = \frac{1}{\sum_{n=1}^N \sum_{c=1}^C y_c^{(n)}} \left(\sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} (x^{(n)} - \mu_c)(x^{(n)} - \mu_c)^T \right)^T = \frac{1}{N} \sum_{c=1}^C \sum_{\substack{n=1 \\ y^{(n)}=c}}^N (x^{(n)} - \mu_c)(x^{(n)} - \mu_c)^T$$

which we can immediately break apart into the representation in the instructions.

Homework

Linear classification

Problem 3: We want to create a generative binary classification model for classifying *non-negative* one-dimensional data. This means, that the labels are binary ($y \in \{0, 1\}$) and the samples are $x \in [0, \infty)$.

We assume uniform class probabilities

$$p(y = 0) = p(y = 1) = \frac{1}{2}.$$

As our samples x are non-negative, we use exponential distributions (and not Gaussians) as class conditionals:

$$p(x \mid y = 0) = \text{Expo}(x \mid \lambda_0) \quad \text{and} \quad p(x \mid y = 1) = \text{Expo}(x \mid \lambda_1),$$

where $\lambda_0 \neq \lambda_1$. Assume, that the parameters λ_0 and λ_1 are known and fixed.

- a) Suppose you are given an observation x . What is the name of the posterior distribution $p(y \mid x)$? You only need to provide the name of the distribution (e.g. “normal”, “gamma” etc.), not estimate its parameter.

Bernoulli.

Remark: y can only take values in $\{0, 1\}$, so obviously Bernoulli is the only possible answer.

- b) What values of x are classified as class 1? (As usual, we assume that the classification decision is $\hat{y} = \arg \max_k p(y = k \mid x)$).

Sample x is classified as class 1 if $p(y = 1 \mid x) > p(y = 0 \mid x)$. This is the same as saying

$$\frac{p(y = 1 \mid x)}{p(y = 0 \mid x)} \stackrel{!}{>} 1 \quad \text{or equivalently} \quad \log \frac{p(y = 1 \mid x)}{p(y = 0 \mid x)} \stackrel{!}{>} 0.$$

We begin by simplifying the left hand side.

$$\begin{aligned} \log \frac{p(y = 1 \mid x)}{p(y = 0 \mid x)} &= \log \frac{p(x \mid y = 1) p(y = 1)}{p(x \mid y = 0) p(y = 0)} \\ &= \log \frac{p(x \mid y = 1)}{p(x \mid y = 0)} \\ &= \log \frac{\lambda_1 \exp(-\lambda_1 x)}{\lambda_0 \exp(-\lambda_0 x)} \\ &= \log \frac{\lambda_1}{\lambda_0} + \lambda_0 x - \lambda_1 x = \log \frac{\lambda_1}{\lambda_0} + (\lambda_0 - \lambda_1)x \end{aligned}$$

To figure out which x are classified as class 1, we need to solve for x .

$$\log \frac{\lambda_1}{\lambda_0} + (\lambda_0 - \lambda_1)x > \log 1 \quad \Leftrightarrow \quad (\lambda_0 - \lambda_1)x > -\log \frac{\lambda_1}{\lambda_0} = \log \lambda_0 - \log \lambda_1$$

We have to be careful, because if $(\lambda_0 - \lambda_1) < 0$, dividing by it will flip the inequality sign. Hence the answer is

$$\begin{cases} x \in \left(\frac{\log \lambda_0 - \log \lambda_1}{\lambda_0 - \lambda_1}, \infty \right) & \text{if } \lambda_0 > \lambda_1 \\ x \in \left[0, \frac{\log \lambda_0 - \log \lambda_1}{\lambda_0 - \lambda_1} \right) & \text{otherwise.} \end{cases}$$

Problem 4: Let $\mathcal{D} = \{(x_i, y_i)\}$ be a linearly separable dataset for 2-class classification, i.e. there exists a vector \mathbf{w} such that $\text{sign}(\mathbf{w}^T \mathbf{x})$ separates the classes. Show that the maximum likelihood parameter \mathbf{w} of a logistic regression model has $\|\mathbf{w}\| \rightarrow \infty$. Assume that \mathbf{w} contains the bias term.

How can we modify the training process to prefer a \mathbf{w} of finite magnitude?

In logistic regression, we model the posterior distribution as

$$y_i | \mathbf{x} \sim \text{Bernoulli}(\sigma(\mathbf{w}^T \mathbf{x}_i)) \quad \text{where } \sigma(a) = \frac{1}{1 + \exp(-a)}.$$

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

We fit the logistic regression model by choosing the parameter \mathbf{w} that maximizes the data log-likelihood or alternatively minimizes the negative log-likelihood which expands to

$$E(\mathbf{w}) = -\log p(\mathbf{y} | \mathbf{w}, \mathbf{X}) = -\sum_{i=1}^N y_i \log \sigma(\mathbf{w}^T \mathbf{x}_i) + (1 - y_i) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_i)).$$

We assumed that the data-set is linearly separable, so by definition there is a $\tilde{\mathbf{w}}$ such that

$$\tilde{\mathbf{w}}^T \mathbf{x}_i > 0 \text{ if } y_i = 1 \quad \text{and} \quad \tilde{\mathbf{w}}^T \mathbf{x}_i < 0 \text{ if } y_i = 0.$$

Scaling this separator $\tilde{\mathbf{w}}$ by a factor $\lambda \gg 0$ makes the negative log-likelihood smaller and smaller. To see this, we compute the limit

$$\lim_{\lambda \rightarrow \infty} E(\lambda \tilde{\mathbf{w}}) = -\left(\sum_{\substack{i=1 \\ y_i=1}}^N \log \lim_{\lambda \rightarrow \infty} \sigma(\lambda \overbrace{\tilde{\mathbf{w}}^T \mathbf{x}_i}^{>0}) + \sum_{\substack{i=1 \\ y_i=0}}^N \log \left(1 - \lim_{\lambda \rightarrow \infty} \sigma(\lambda \overbrace{\tilde{\mathbf{w}}^T \mathbf{x}_i}^{<0}) \right) \right) = 0$$

which equals the smallest achievable value (E is the negative log of a probability, so $E(\mathbf{w}) \in [0, \infty)$ and thus $E(\mathbf{w}) \geq 0$).

We can see that E is a convex function because \log is concave and σ is convex if $a < 0$ and concave if $a > 0$. So $\log \sigma(a)$ is concave if $a > 0$ and $\log(1 - \sigma(a))$ is concave if $a < 0$. It follows that E is a convex function because E is the negative sum of concave functions.

A convex function has a unique minimum *if* it attains its minimum value. We know that E tends towards its minimum as $\lambda \rightarrow \infty$, so E cannot have a finite minimizer and all its minima are only achieved in the limit. It follows that any solution to the loss minimization problem has infinite norm.

Because E is convex and tends towards a limit of 0 in some directions, we can move the minimum into the space of finite vectors by adding any convex term that achieves its minimum such as $\mathbf{w}^T \mathbf{w}$ or similar forms of weight regularization.

Problem 5: Show that the softmax function is equivalent to a sigmoid in the 2-class case.

$$\begin{aligned} \frac{\exp(\mathbf{w}_1^T \mathbf{x})}{\exp(\mathbf{w}_1^T \mathbf{x}) + \exp(\mathbf{w}_0^T \mathbf{x})} &= \frac{1}{1 + \exp(\mathbf{w}_0^T \mathbf{x}) / \exp(\mathbf{w}_1^T \mathbf{x})} \\ &= \frac{1}{1 + \exp((\mathbf{w}_0 - \mathbf{w}_1)^T \mathbf{x})} \\ &= \frac{1}{1 + \exp(-\hat{\mathbf{w}}^T \mathbf{x})} \end{aligned}$$

where $\hat{\mathbf{w}} = \mathbf{w}_1 - \mathbf{w}_0$.

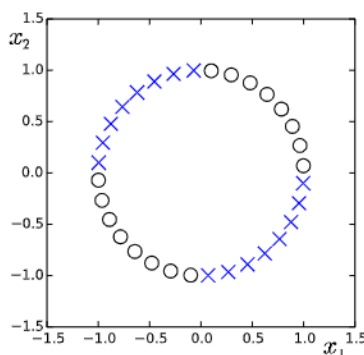
One conclusion we can draw from this is that if we have C parameter vectors \mathbf{w}_c for C classes, the logistic regression model is unidentifiable. This means that adding a constant $\boldsymbol{\tau} \in \mathbb{R}^D$ to each vector $\mathbf{w}_c := \mathbf{w}_c + \boldsymbol{\tau}$ would lead to the same logistic regression model. We can fix this issue by adding a constraint $\mathbf{w}_1 = \mathbf{0}$, which is what is done implicitly when we use sigmoid (instead of 2-class softmax) in binary classification.

Problem 6: Show that the derivative of the sigmoid function $\sigma(a) = (1 + e^{-a})^{-1}$ can be written as

$$\frac{\partial \sigma(a)}{\partial a} = \sigma(a) (1 - \sigma(a)).$$

$$\frac{\partial \sigma(a)}{\partial a} = -\frac{1}{(1 + e^{-a})^2} \cdot e^{-a} \cdot (-1) = \frac{1}{1 + e^{-a}} \frac{e^{-a}}{1 + e^{-a}} = \sigma(a) \frac{1 + e^{-a} - 1}{1 + e^{-a}} = \sigma(a) (1 - \sigma(a))$$

Problem 7: Give a basis function $\phi(x_1, x_2)$ that makes the data in the example below linearly separable (crosses in one class, circles in the other).



One example is $\phi(x) = x_1 x_2$ which makes the data separable by the hyperplane $w = (1)$ because the circles will be mapped to the positive real numbers while the crosses go to the negative numbers, i.e. $w^T x > 0$ if x is a circle and $w^T x < 0$ otherwise.

Naive Bayes

Problem 8: In 2-class classification the decision boundary Γ is the set of points where both classes are assigned equal probability,

$$\Gamma = \{x \mid p(y = 1 \mid x) = p(y = 0 \mid x)\}.$$

Show that Naive Bayes with Gaussian class likelihoods produces a quadratic decision boundary in the 2-class case, i.e. that Γ can be written with a quadratic equation of x ,

$$\Gamma = \{x \mid x^T A x + b^T x + c = 0\},$$

for some A , b and c .

As a reminder, in Naive Bayes we assume class prior probabilities

$$p(y = 0) = \pi_0 \quad \text{and} \quad p(y = 1) = \pi_1$$

and class likelihoods

$$p(x \mid y = c) = \mathcal{N}(x \mid \mu_c, \Sigma_c)$$

with per-class means μ_c and *diagonal* (because of the feature independence) covariances Σ_c .

Because $p(y = 1 \mid x) + p(y = 0 \mid x) = 1$ and we want them to be equal, we can assume that $p(y = 0 \mid x) > 0$ and rewrite the defining equation as

$$\frac{p(y = 1 \mid x)}{p(y = 0 \mid x)} = 1.$$

Now apply the logarithm to both sides and simplify.

$$\begin{aligned}
 \log \frac{p(y=1|x)}{p(y=0|x)} &= \log \left(\frac{p(x|y=1)p(y=1)}{p(x)} \cdot \frac{p(x)}{p(x|y=0)p(y=0)} \right) \\
 &= \log(p(x|y=1)p(y=1)) - \log(p(x|y=0)p(y=0)) \\
 &= \log \mathcal{N}(x|\mu_1, S_1) - \log \mathcal{N}(x|\mu_0, S_0) + \log \frac{\pi_1}{\pi_0} \\
 &= -\frac{1}{2} \log(2\pi)^D |S_1| - \frac{1}{2} (x - \mu_1)^T S_1^{-1} (x - \mu_1) \\
 &\quad + \frac{1}{2} \log(2\pi)^D |S_0| + \frac{1}{2} (x - \mu_0)^T S_0^{-1} (x - \mu_0) + \log \frac{\pi_1}{\pi_0} \\
 &= -\frac{1}{2} x^T S_1^{-1} x + x^T S_1^{-1} \mu_1 - \frac{1}{2} \mu_1^T S_1^{-1} \mu_1 \\
 &\quad + \frac{1}{2} x^T S_0^{-1} x - x^T S_0^{-1} \mu_0 + \frac{1}{2} \mu_0^T S_0^{-1} \mu_0 + \frac{1}{2} \log \frac{|S_0|}{|S_1|} + \log \frac{\pi_1}{\pi_0} \\
 &= \frac{1}{2} x^T [S_0^{-1} - S_1^{-1}] x + x^T [S_1^{-1} \mu_1 - S_0^{-1} \mu_0] \\
 &\quad + \frac{1}{2} \mu_1^T S_1^{-1} \mu_1 - \frac{1}{2} \mu_0^T S_0^{-1} \mu_0 + \log \frac{\pi_1}{\pi_0} + \frac{1}{2} \log \frac{|S_0|}{|S_1|}
 \end{aligned}$$

Assignment Project Exam Help

This shows that Γ is quadratic and can alternatively be written as

$$\Gamma = \{x | x^T A x + b x + c = 0\}$$

where

$$A = \frac{1}{2} S_0^{-1} - S_1^{-1}, \quad b = S_1^{-1} \mu_1 - S_0^{-1} \mu_0$$

$$c = -\frac{1}{2} \mu_1^T S_1^{-1} \mu_1 + \frac{1}{2} \mu_0^T S_0^{-1} \mu_0 + \log \frac{\pi_1}{\pi_0} + \frac{1}{2} \log \frac{|S_0|}{|S_1|}.$$

If both classes had the same covariance matrix ($S_0 = S_1$), A would be the zero matrix and we would obtain a linear decision boundary as we did in the lecture (also, $\log \frac{|S_0|}{|S_1|} = 0$).