


Problem 1: What is the connection between soft-margin SVM and logistic regression?

The soft-margin SVM is defined via the minimization problem

$$\text{minimize } f_0(w, b, \epsilon) = \frac{1}{2} w^T w + C \sum_{i=1}^N \epsilon_i$$

$$\text{subject to } y_i (w^T x_i + b) - 1 + \epsilon_i \geq 0 \quad i = 1, \dots, N$$

$$\epsilon_i \geq 0 \quad i = 1, \dots, N$$

Due to the complementary slackness conditions, we have

$$\alpha_i (y_i (w^T x_i + b) - 1 + \epsilon_i) = 0$$

Thus, for points that lie inside or beyond the margin (i.e. we have $\alpha_i > 0$ and $y_i (w^T x_i + b) < 1$), it must hold that $\epsilon_i = 1 - y_i (w^T x_i + b)$. Otherwise, ϵ_i is minimized and therefore 0. In other words,

$$\epsilon_i = \begin{cases} 1 - (w^T x_i + b) & \text{if } y_i (w^T x_i + b) < 1 \\ 0 & \text{otherwise} \end{cases} = \max(0, 1 - y_i (w^T x_i + b)) = \text{hinge}(y_i (w^T x_i + b))$$

Dividing by C, the error (or loss) function is

$$E(w, b, C) = \underbrace{\frac{1}{2C}}_J w^T w + \sum_{i=1}^N \text{hinge}(y_i (w^T x_i + b)) \quad (1)$$

Applying logistic regression to the same classification problem, we have

$$p(y_i = 1 | x_i, w) = \sigma(w^T x_i + b),$$

$$p(y_i = -1 | x_i, w) = \sigma(-(w^T x_i + b)) = 1 - \sigma(w^T x_i + b),$$

where the last equality holds due to the definition of the sigmoid function,

$\sigma(z) = \frac{1}{1+e^{-z}}$. Keep in mind that in this case $y_i \in \{-1, 1\}$ and not $\{0, 1\}$ as we use originally.

$$\text{Altogether, we have } p(y | x, w) = \prod_{i=1}^N \sigma(y_i (w^T x_i + b))$$

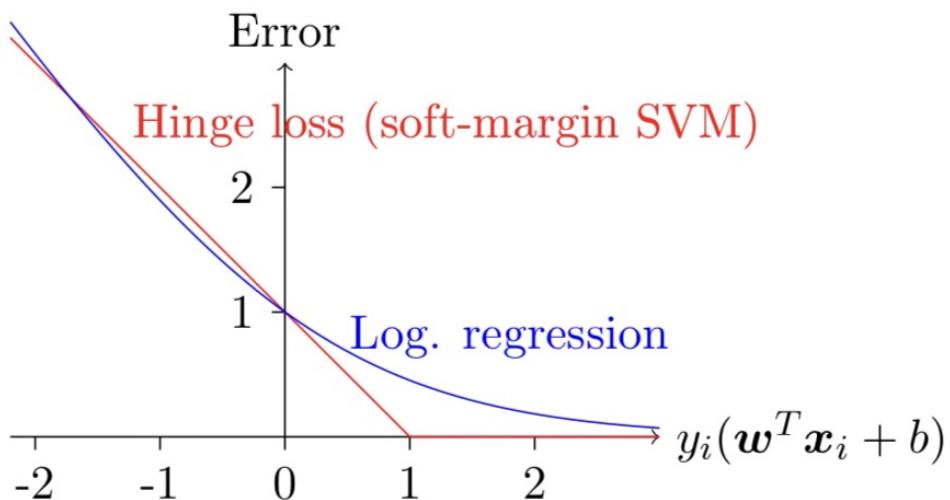
We use the negative log-likelihood as the error function, i.e.

$$J(w, b) = -\ln(p(y|X, w)) = -\sum_{i=1}^N \ln\left(\frac{1}{1+e^{-y_i(w^T x_i + b)}}\right) = \sum_{i=1}^N \ln(1+e^{-y_i(w^T x_i + b)})$$

Additionally, we introduce L2 regularization term and get

$$J(w, b, \lambda) = \lambda w^T w + \sum_{i=1}^N \ln(1+e^{-y_i(w^T x_i + b)}) \quad (2)$$

Hence, we see the close relationship between the soft-margin SVM (1) and logistic regression (2). While soft-margin SVM uses the hinge function for its loss, logistic regression uses $\ln(1+e^{-x})$. For better comparison, we can plot these two (with logistic regression rescaled by $\frac{1}{\ln 2}$):



Problem 2: Consider a soft-margin SVM fitted to a linearly separable dataset \mathcal{D} using the Hinge loss formulation of the optimization task.

$$\underset{\mathbf{w}, b}{\text{minimize}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$$

a) Is it guaranteed that all training samples in \mathcal{D} will be assigned the correct label by the model?

No, as it might happen that misclassifying a few points that are very close to the decision boundary would lead to a significantly increasing the margin.

In general, larger values of C make this behavior less likely.

b) Prove that if for some $C_0 \geq 0$ the resulting model classifies all training samples correctly and all samples lie outside of the margin then it will also be the case if we train the model with any larger $C > C_0$.

Denote by $h(\mathbf{w}, b) = \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$ the overall Hinge loss and let $C > C_0 \geq 0$ be as in the task formulation. Note that the hyperplane defined by \mathbf{w} and b assigns the correct labels to all samples such that all of them also lie outside of the margin if and only if $h(\mathbf{w}, b) = 0$.

We also define the corresponding solutions as

$$(\mathbf{w}^*, b^*) = \underset{\mathbf{w}, b}{\arg \min} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C_0 h(\mathbf{w}, b) \text{ and}$$

$$(\mathbf{v}^*, d^*) = \underset{\mathbf{w}, b}{\arg \min} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C h(\mathbf{w}, b)$$

Then the following inequality holds due to optimality of these points:

$$\frac{1}{2} \|\mathbf{w}^*\|^2 + C_0 h(\mathbf{w}^*, b^*) \leq \frac{1}{2} \|\mathbf{v}^*\|^2 + C h(\mathbf{v}^*, d^*).$$

$$\frac{1}{2} \|\mathbf{v}^*\|^2 + C h(\mathbf{v}^*, d^*) \leq \frac{1}{2} \|\mathbf{w}^*\|^2 + C h(\mathbf{w}^*, b^*)$$

Summing the left and right hand side of these inequalities we get

$$\frac{1}{2} \|w^*\|^2 + \frac{1}{2} \|v^*\|^2 + \mathcal{L}(h(w^*, b^*), h(v^*, d^*)) \leq \frac{1}{2} \|v^*\|^2 + \frac{1}{2} \|w^*\|^2 + \mathcal{L}(h(v^*, d^*), h(w^*, b^*))$$



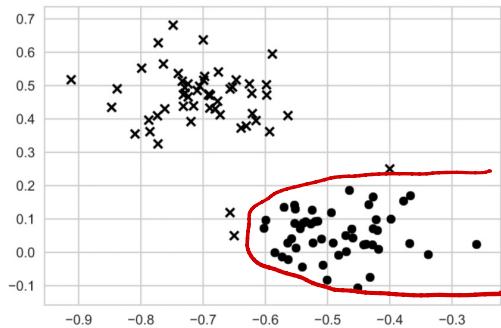
$$((C - C_0)(h(v^*, d^*) - h(w^*, b^*))) \leq 0 \quad \text{and since } C - C_0 > 0$$



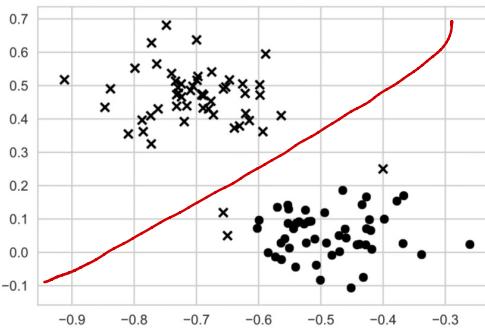
$$h(v^*, d^*) \leq h(w^*, b^*)$$

That means the overall Hinge loss decreases for the optimal solution if we solve the SVM fitting problem with a larger value of C . Since we could separate the data perfectly (with no points within the margin) with the solution (w^*, b^*) for C we have $h(w^*, b^*) = 0$. As shown above the Hinge Loss can not get worse for larger C , hence also $h(v^*, d^*)$ has to be 0 and the solution (v^*, d^*) also corresponds to a decision boundary with no error on training data and no points lying inside the margin.

Problem 3: Sketch the decision boundary of an SVM with a quadratic kernel (polynomial with degree 2) for the data in the figure below, for two specified values of the penalty parameter C . The two classes are denoted as \bullet 's and \times 's.



(a) $C = 10^{10}$



(b) $C = 10^{-10}$

- a) With such a large penalty, SVM will try to correctly classify all of the instances in the training set
- b) Given the small penalty, we can allow few missclassified instances, and obtain a larger margin between the two classes.

2 Kernels

Problem 4: Consider the Gaussian kernel

$$k_G(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{2\sigma^2}\right), \text{ with } \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d.$$

- a) Suppose you have found a feature map $\theta : \mathbb{R}^d \rightarrow \mathbb{R}^{d_1}$ that transforms your data into a feature space in which a SVM with a Gaussian kernel works well. However, computing $\theta(\mathbf{x})$ is computationally expensive and luckily you discover an efficient method to compute the scalar product

$$k(\mathbf{x}_1, \mathbf{x}_2) = \theta(\mathbf{x}_1)^T \theta(\mathbf{x}_2)$$

in your feature space without having to compute $\theta(\mathbf{x}_1)$ and $\theta(\mathbf{x}_2)$ explicitly. Show how you can use the scalar product $k(\mathbf{x}_1, \mathbf{x}_2)$ to efficiently compute the Gaussian kernel $k_G(\theta(\mathbf{x}_1), \theta(\mathbf{x}_2))$ in your feature space.

By expanding the quadratic term and applying the definition of $K(x_1, x_2)$ we get

$$\begin{aligned} K_G(\theta(\mathbf{x}_1), \theta(\mathbf{x}_2)) &= \exp\left(-\frac{\|\theta(\mathbf{x}_1) - \theta(\mathbf{x}_2)\|^2}{2\sigma^2}\right) \\ &= \exp\left(-\frac{(\theta(\mathbf{x}_1) - \theta(\mathbf{x}_2))^T (\theta(\mathbf{x}_1) - \theta(\mathbf{x}_2))}{2\sigma^2}\right) \\ &= \exp\left(-\frac{\theta(\mathbf{x}_1)^T \theta(\mathbf{x}_1) - 2\theta(\mathbf{x}_1)^T \theta(\mathbf{x}_2) + \theta(\mathbf{x}_2)^T \theta(\mathbf{x}_2)}{2\sigma^2}\right) \\ &= \exp\left(-\frac{K(\mathbf{x}_1, \mathbf{x}_1) - 2K(\mathbf{x}_1, \mathbf{x}_2) + K(\mathbf{x}_2, \mathbf{x}_2)}{2\sigma^2}\right) \end{aligned}$$

Thus we can compute $K_G(\theta(\mathbf{x}_1), \theta(\mathbf{x}_2))$ from $K(\mathbf{x}_1, \mathbf{x}_2)$ using the above equation

- b) One of the nice things about kernels is that new kernels can be constructed out of already given ones. Use the five kernel construction rules from the lecture to prove that k_G is a kernel.

Hint: Use the Taylor expansion of the exponential function to prove that $\exp \circ k_1$ is a kernel if k_1 is a kernel. Also, consider $k_2(\phi(x_1), \phi(x_2))$ with the linear kernel $k_2(x_1, x_2) = x_1^T x_2$ and a feature map ϕ with only one feature.

First we prove that $\exp(K(x_1, x_2))$ is a Kernel if $K(x_1, x_2)$ is a Kernel. The Taylor expansion of the exponential function is

$$\exp(K(x_1, x_2)) = 1 + \sum_{n=1}^{\infty} \frac{1}{n!} K(x_1, x_2)^n$$

The power $K(x_1, x_2)^n$ is a Kernel by iterated application of rule 3 ($K_1(x_1, x_2) K_2(x_1, x_2)$ is a Kernel). The product $(1/n!) K(x_1, x_2)^n$ is a Kernel by rule 2 (if $K_1(x_1, x_2)$ is a Kernel for $d > 0$) because $(1/n!)$ is always positive. The sum $\sum_{n=1}^N (1/n!) K(x_1, x_2)^n$ is a Kernel for arbitrary $N \in \mathbb{N}$ by iterated application of rule 1 ($K_1(x_1, x_2) + K_2(x_1, x_2)$ is a Kernel). Therefore, the pointwise limit of kernels

$$\sum_{n=1}^{\infty} \frac{1}{n!} K(x_1, x_2)^n = \lim_{N \rightarrow \infty} \sum_{n=1}^N \frac{1}{n!} K(x_1, x_2)^n$$

is again a Kernel. The constant 1 is a Kernel by rule 4: $K_3(\phi(x_1), \phi(x_2))$ with $K_3(x_1, x_2) = x_1^T x_2$ and $\phi(z) = 1$. Thus $1 + \sum_{n=1}^{\infty} \frac{1}{n!} K(x_1, x_2)^n$ is a Kernel by rule 1.

We expand the argument of the exponential function.

$$\exp\left(-\frac{(x_1 - x_2)^2}{2\sigma^2}\right) = \exp\left(-\frac{x_1^T x_1}{2\sigma^2}\right) \exp\left(-\frac{x_2^T x_2}{2\sigma^2}\right) \exp\left(\frac{x_1^T x_2}{2\sigma^2}\right)$$

Consider the last term first. The scalar product $x_1^T x_2$ is linear Kernel by rule 2 the product $x_1^T x_2 / \sigma^2$ is a Kernel because σ^2 is positive. As proved above the term $\exp(x_1^T x_2 / \sigma^2)$ is then a Kernel as well.

The product of the first 2 terms $\exp(-x_1^T x_1 / 2\sigma^2) \exp(-x_2^T x_2 / 2\sigma^2)$ is a Kernel by rule 4 with $K_3(x_1, x_2) = x_1^T x_2$ and the feature map $\phi(z) = \exp(-z^T z / 2\sigma^2)$.

Finally, by rule 3 the product of the first two terms with the third term is a Kernel.

- c) Can any finite set of points be linearly separated in the feature space of the Gaussian kernel if σ can be chosen freely?

Consider the limit $\sigma \rightarrow 0$. Then

$$K_G(x_1, x_2; \sigma) \xrightarrow{\sigma \rightarrow 0} K(x_1, x_2) = \begin{cases} 1 & \text{if } x_1 = x_2 \\ 0 & \text{if } x_1 \neq x_2 \end{cases}$$

All training samples are correctly classified if

$$y_i (w^T \phi_G(x_i) + b) > 0 \quad \text{for all } i$$

We will construct a classifier with the desired property by using the dual representation of $w = \sum_j y_j \alpha_j \phi_G(x_j)$. Substituting it into the above expression and replacing the scalar product $\phi_G(x_i)^T \phi_G(x_j)$ by the Kernel function $K_G(x_i, x_j; \sigma)$ we get the same condition in the following form

$$y_i \left(\sum_j y_j \alpha_j K_G(x_i, x_j; \sigma) + b \right) > 0$$

Using the limit Kernel function K in particular we can see that the left hand side will converge to $\xrightarrow{\sigma \rightarrow 0} y_i^2 \alpha_i + y_i b$

for small σ . By choosing $b=0$ we see that the resulting condition is fulfilled for all training samples since $y_i^2 = 1$ for all i and we can simply set all $\alpha_i > 0$. Hence all finite sets of points can be linearly separated using the Gaussian kernel if the variance σ is chosen small enough.

Problem 5: Let $\mathcal{M} = \cup_{n \in \mathbb{N}} \cup_{m \in \mathbb{N}} \mathbb{R}^{n \times m}$ denote the set of all real-valued matrices of arbitrary size. Prove that the function $k : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}$ is a valid kernel, where

$$k(\mathbf{X}, \mathbf{Y}) = \min(\text{rank}(\mathbf{X}), \text{rank}(\mathbf{Y})).$$

We will prove it by constructing the corresponding feature map. Consider the function $\phi(\mathbf{x}) = (\underbrace{1, 1, 1, \dots, 1}_{\text{rank}(\mathbf{x}) \text{ ones}}, \underbrace{0, 0, 0, 0, \dots}_{\text{zeros everywhere else}})$

$$\text{Then } \phi(\mathbf{x})^\top \phi(\mathbf{y}) = \sum_{j=1}^{\infty} \phi_j(\mathbf{x}) \phi_j(\mathbf{y}) = \min(\text{rank}(\mathbf{x}), \text{rank}(\mathbf{y})) = k(\mathbf{x}, \mathbf{y})$$

the inner product. Hence, $k(\mathbf{x}, \mathbf{y})$ is a valid kernel.

3 SVM

Problem 6: Explain the similarities and differences between the SVM and perceptron algorithms. How do they perform classification? In what way do they differ?

Both are supervised classification methods separate two classes by a hyperplane. The difference is that SVM also tries to maximize the margin between the hyperplane and data, while perceptron only cares about separation.

Problem 7: Recall that the dual function in the setting of the SVM training task (Slide 17) can be written as

$$g(\alpha) = \frac{1}{2} \alpha^T Q \alpha + \alpha^T \mathbf{1}_N.$$

- (a) Write down the matrix Q using the vector of labels \mathbf{y} and feature matrix \mathbf{X} . Denote the element-wise product between two matrices (in case you want to use it) by \odot (also known as Hadamard product or Schur product).

Dual function g is $\rightarrow g(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \mathbf{y}_i \mathbf{y}_j \mathbf{x}_i^T \mathbf{x}_j$

can be written as $\rightarrow g(\alpha) = \frac{1}{2} \alpha^T Q \alpha + \alpha^T \mathbf{1}_N = \sum_{i=1}^N \alpha_i + \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i Q_{ij} \alpha_j$

meaning that $Q_{ij} = -\mathbf{y}_i \mathbf{y}_j \mathbf{x}_i^T \mathbf{x}_j$. From that we see that Q can be decomposed as (-1 times) the element-wise product of two matrices A and B with $A_{ij} = \mathbf{y}_i \mathbf{y}_j$ and $B_{ij} = \mathbf{x}_i^T \mathbf{x}_j$ so that

$$A = \begin{pmatrix} \mathbf{y}_2 \mathbf{y}_1 & \mathbf{y}_1 \mathbf{y}_2 & \cdots & \mathbf{y}_1 \mathbf{y}_N \\ \mathbf{y}_2 \mathbf{y}_2 & \mathbf{y}_2 \mathbf{y}_2 & \cdots & \mathbf{y}_2 \mathbf{y}_N \\ \vdots & \vdots & & \vdots \\ \mathbf{y}_N \mathbf{y}_1 & \mathbf{y}_N \mathbf{y}_2 & \cdots & \mathbf{y}_N \mathbf{y}_N \end{pmatrix} = \mathbf{y} \mathbf{y}^T \quad \text{and} \quad B = \begin{pmatrix} \mathbf{x}_1^T \mathbf{x}_1 & \mathbf{x}_1^T \mathbf{x}_2 & \cdots & \mathbf{x}_1^T \mathbf{x}_N \\ \mathbf{x}_2^T \mathbf{x}_1 & \mathbf{x}_2^T \mathbf{x}_2 & \cdots & \mathbf{x}_2^T \mathbf{x}_N \\ \vdots & \vdots & & \vdots \\ \mathbf{x}_N^T \mathbf{x}_1 & \mathbf{x}_N^T \mathbf{x}_2 & \cdots & \mathbf{x}_N^T \mathbf{x}_N \end{pmatrix}$$

Finally, we get $Q = -\mathbf{y} \mathbf{y}^T \odot \mathbf{X} \mathbf{X}^T$

- (b) Prove that we can search for a *local* maximizer of g to find its *global* maximum (don't forget to prove properties of Q that you decide to use in this task).

We will first show that $A \odot B = \mathbf{y} \mathbf{y}^T \odot \mathbf{X} \mathbf{X}^T \in \mathbb{R}^{N \times N}$ is a positive semi-definite matrix by definition. For that consider an arbitrary $\alpha \in \mathbb{R}^N$ and

$$\begin{aligned} \alpha^T (A \odot B) \alpha &= \sum_{i=1}^N \sum_{j=1}^N \alpha_i \mathbf{y}_i \mathbf{y}_j \mathbf{x}_i^T \mathbf{x}_j \alpha_j \\ &= \sum_{i=1}^N \sum_{j=1}^N (\alpha_i \mathbf{y}_i \mathbf{x}_i^T) (\alpha_j \mathbf{y}_j \mathbf{x}_j) = \left(\sum_{i=1}^N \alpha_i \mathbf{y}_i \mathbf{x}_i^T \right) \left(\sum_{j=1}^N \alpha_j \mathbf{y}_j \mathbf{x}_j^T \right) = \left\| \sum_{i=1}^N \alpha_i \mathbf{y}_i \mathbf{x}_i \right\|^2 \geq 0 \end{aligned}$$

As $A \odot B$ is PSD, it follows that $Q = -A \odot B$ is negative semi-definite. This means that the function g is concave, thus every local maximum is a global maximum.

Problem 8: Consider training a standard hard-margin SVM on a linearly separable training set of N samples. Let s denote the number of support vectors we would obtain if we would train on the entire dataset. Furthermore, let ε denote the leave-one-out cross validation (LOOCV) misclassification rate. Prove that the following relation holds:

$$\varepsilon \leq \frac{s}{N}.$$

Intuitively, the result follows from the following claim (which we prove below): if x_i is not a support vector when training on the entire training set, then the optimal w^* and b^* do not change when leaving x_i out of the training set.

Since the originally data are linearly separable and since we are using a hard margin classifier, the hypothesis given by the original w^* and b^* will not make an error on x_i , and hence, no error will be made in the i -th step of the LOOCV. Equivalently, the only possible errors in the LOOCV procedure are made on x_i 's which are support vectors when training on the entire training set, and hence $\varepsilon \leq \frac{s}{N}$

Formal prove:

Let (w_0^*, b_0^*) and α_0^* denote the optimal primal and dual solutions for the SVM when training on the entire D . Also let, $D_i = D \setminus \{(x_i, y_i)\}$ be the set of training examples when omitting the i -th example, and let $(w_{D_i}^*, b_{D_i}^*)$ and $\alpha_{D_i}^*$ be the optimal primal and dual variables of the optimization problem when training on D_i .

$\alpha_{D_i}^*$ consists of $n-1$ variables, namely $\alpha_{D_i,1}, \dots, \alpha_{D_i,i+1}, \dots, \alpha_{D_i,N}$. Now consider setting the dual variables as follows $\alpha_{D_i,j} = \alpha_{D,j}^*$ for $j \neq i$. Note that, if x_i is not a support vector when training on D then $\alpha_{D_i,i}^* = 0$. We can verify that (w_0^*, b_0^*) and $\alpha_{D_i}^*$ satisfy the KKT conditions for the SVM optimization problem when training on D_i . From this, and the fact that w_0^*, b_0^* are unique since the objective function is strictly convex we can conclude that w^*, b^* do not change when omitting $\{(x_i, y_i)\}$ as desired.

4 Kernels

Problem 10: Show that for $N \in \mathbb{N}$ and $a_i \geq 0$ for $i = 0, \dots, N$ the following function k is a valid kernel.

$$k(\mathbf{x}_1, \mathbf{x}_2) = \sum_{i=1}^N a_i (\mathbf{x}_1^T \mathbf{x}_2)^i + a_0, \text{ with } \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d.$$

The function $K_1(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{x}_2$ is a kernel because it is the scalar product of the input vectors. The constant function $K_c(\mathbf{x}_1, \mathbf{x}_2) = a_0 \geq 0$ is a kernel because we can define the feature map $\phi(x) = \sqrt{a_0}$ and obtain this kernel by calculating the scalar product in feature space $\phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) = \sqrt{a_0}^2 = a_0$. Overall, K is a kernel since it is build up of sums and products of kernels and multiplication with non-negative scalars.

Problem 11: Find the feature transformation $\phi(x)$ corresponding to the kernel

$$k(x_1, x_2) = \frac{1}{1 - x_1 x_2}, \text{ with } x_1, x_2 \in (0, 1).$$

Hint: Consider an infinite-dimensional feature space and infinite series.

We use the geometric series to transform K :

$$K(x_1, x_2) = \frac{1}{1 - x_1 x_2} = \sum_{i=0}^{\infty} x_1^i x_2^i = \phi(x_1)^\top \phi(x_2).$$

with the feature transformation

$$\phi(x) = [1, x, x^2, x^3, x^4, \dots]^\top$$

Geometric series:

$$S = \frac{1}{1-r} = 1 + r + r^2 + \dots \quad \text{if } 0 < r < 1$$