

Kernel Methods

Winter 2023

Clayton Scott

These notes introduce nonlinear feature maps, a powerful technique for converting a linear algorithm to a nonlinear one, and kernels, a structure that makes this technique implementable in high dimensional settings.

1 Nonlinear Feature Maps

One way to create a nonlinear method for regression or classification is to first transform the feature vector via a nonlinear map

$$\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$$

and then apply a linear method to the transformed data $\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_n)$. In regression, the nonlinear function estimate is:

$$f(x) = \mathbf{w}^T \Phi(\mathbf{x}) + b$$

where $\mathbf{w} \in \mathbb{R}^m, b \in \mathbb{R}$.

Example 1. The goal of this example is to determine the least squares cubic polynomial fit to training data $(x_1, y_1), \dots, (x_n, y_n)$ where $x_i \in \mathbb{R}, y_i \in \mathbb{R}$.

Let's write $f(x) = a + bx + cx^2 + dx^3 = \mathbf{w}^T \Phi(x) + a$ where

$$\mathbf{w} = \begin{bmatrix} b \\ c \\ d \end{bmatrix} \quad \Phi(x) = \begin{bmatrix} x \\ x^2 \\ x^3 \end{bmatrix}.$$

Note that $\Phi(x)$ is nonlinear.

Since

$$\sum_i (y_i - f(x_i))^2 = \sum_i (y_i - \mathbf{w}^T \Phi(x_i) - a)^2,$$

The minimizer is

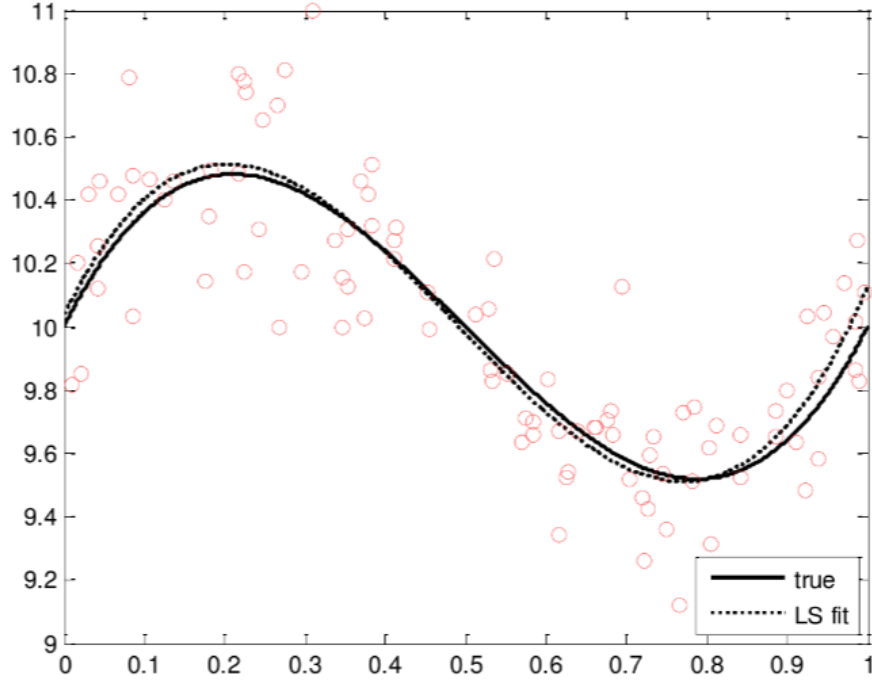
$$\begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

where

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & x_n^3 \end{bmatrix}$$

See Figure 1.

Figure 1: Least squares cubic polynomial fit.



The preceding example is a good motivation for regularization. As the polynomial degree increases, the matrix $\mathbf{X}^T \mathbf{X}$ becomes ill-conditioned, and regularization is necessary to avoid over-fitting.

In classification, the nonlinear classifier is expressed $f(\mathbf{x}) = \text{sign}\{\mathbf{w}^T \Phi(\mathbf{x}) + b\}$

Example 2. Consider the data shown in Figure 2. This data set can be perfectly classified by a circular classifier, that is, a classifier of the form

$$f(\mathbf{x}) = \text{sign}\{(x_1 - c_1)^2 + (x_2 - c_2)^2 - r^2\}$$

for a certain center

$$\mathbf{c} = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}$$

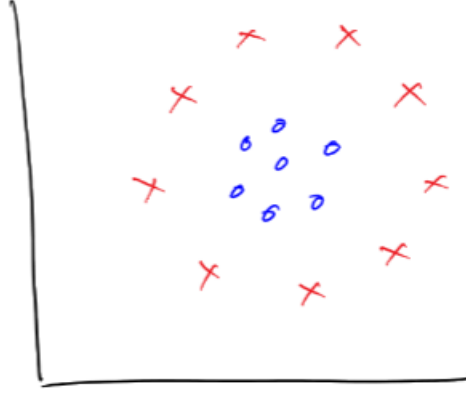
and radius r , and where the feature vector is

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^2.$$

This circular classifier can be viewed as a linear classifier in a transformed feature space. Consider

$$\Phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix}$$

Figure 2: Data separated by a circular classifier.



Then the circular classifier can be expressed

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \Phi(\mathbf{x}) + c)$$

where

$$\mathbf{w} = \begin{bmatrix} -2c_1 \\ -2c_2 \\ 0 \\ 1 \\ 1 \end{bmatrix}, \quad b = c_1^2 + c_2^2 - r^2.$$

In both of the preceding examples, it is also possible to eliminate the offset term by incorporating a constant entry of 1 into Φ . The point of these examples is that while linear predictors may not work well in the original feature space, they may be highly effective in a transformed feature space.

Many nonlinear machine learning methods are based on nonlinear feature transformations. Neural networks are a prominent example that we will return to. We now turn our attention to another class of methods known as kernel methods.

2 Kernels

One issue with nonlinear feature transformations is the way m explodes as d increases. $\Phi(\mathbf{x})$ cannot be computed explicitly, or stored, because it is too large. For example, if we extend the cubic polynomial to arbitrary degree p and dimension d , the total number of terms in the feature expansion is $\binom{p+d}{d}$. If the dimension is 100 and the degree is 10, this number exceeds 46 trillion!

Fortunately, the following two facts save us: (1) Many machine learning algorithms depend on $\Phi(\mathbf{x})$ only via *inner products* $\langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$, and (2) For certain Φ , the function

$$k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$$

can be computed efficiently even if m is huge (and possibly infinite)!

If a machine learning algorithm satisfies the first property, and Φ satisfies the second property, then we can use the function k to incorporate the nonlinear feature map Φ into the algorithm, without ever having

to explicitly compute $\Phi(\mathbf{x})$. Such functions k are called *kernels*,¹ and the algorithms obtained by using them are known as *kernel methods*. We will return to this idea below, and see several examples of kernel methods in subsequent lectures. First, we offer two equivalent formal definitions of kernels and offer some examples.

2.1 Inner Product Kernels

We begin with a simple example to illustrate the second property needed for kernel methods.

Example 3. When $d = 2$, the function $k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v})^2$ is associated to an $m = 3$ dimensional feature map as follows:

$$\begin{aligned} k(\mathbf{u}, \mathbf{v}) &= (\mathbf{u}^T \mathbf{v})^2 \\ &= \left(\begin{bmatrix} u_1 & u_2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \right)^2 \\ &= (u_1 v_1 + u_2 v_2)^2 \\ &= u_1^2 v_1^2 + 2u_1 u_2 v_1 v_2 + u_2^2 v_2^2 \\ &= \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle \end{aligned}$$

where

$$\Phi(\mathbf{u}) = \begin{bmatrix} u_1^2 \\ \sqrt{2}u_1 u_2 \\ u_2^2 \end{bmatrix}.$$

Extending to arbitrary d we have

$$\begin{aligned} k(\mathbf{u}, \mathbf{v}) &= (\mathbf{u}^T \mathbf{v})^2 \\ &= \left(\sum_{i=1}^d u_i v_i \right)^2 \\ &= \left(\sum_i u_i v_i \right) \left(\sum_j u_j v_j \right) \\ &= \sum_{i=1}^d \sum_{j=1}^d u_i u_j v_i v_j \\ &= \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle \end{aligned}$$

where $\Phi(\mathbf{u}) = [u_1^2 \cdots u_d^2 \sqrt{2}u_1 u_2 \cdots \sqrt{2}u_{d-1} u_d]^T$, and so $m = d + \frac{d(d-1)}{2}$.

Generalizing the above to arbitrary polynomial degree,

$$\begin{aligned} k(\mathbf{u}, \mathbf{v}) &= (\mathbf{u}^T \mathbf{v})^p \\ &= \sum_{(j_1 \dots j_d)} \binom{p}{j_1 \dots j_d} u_1^{j_1} \cdots u_d^{j_d} v_1^{j_1} \cdots v_d^{j_d} \\ &= \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle \end{aligned}$$

where $\Phi(\mathbf{u}) = [\cdots, \sqrt{\binom{p}{j_1 \dots j_d}} u_1^{j_1} \cdots u_d^{j_d}, \cdots]^T$. In words, the features associated to k are all monomials of degree p .

¹The term *kernel* is overloaded in machine learning, and can have a different meaning for different algorithms such as locally linear regression and kernel density estimation.

Another example comes from the function $k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v} + 1)^p$, whose feature map consists of all monomials (with appropriate scaling) of degree *at most* p . This is left as an exercise.

In these examples, the inner product is just the standard dot product on \mathbb{R}^d . More generally, we can allow the range of Φ to be any (real) inner product space, which is a vector space V on which we can define a real-valued function $\langle \mathbf{u}, \mathbf{v} \rangle$ (called an *inner product*), which must satisfy the following axioms:

1. $\langle \alpha_1 \mathbf{u}_1 + \alpha_2 \mathbf{u}_2, \mathbf{v} \rangle = \alpha_1 \langle \mathbf{u}_1, \mathbf{v} \rangle + \alpha_2 \langle \mathbf{u}_2, \mathbf{v} \rangle \quad \forall \alpha_1, \alpha_2 \in \mathbb{R}, \mathbf{u}_1, \mathbf{u}_2, \mathbf{v} \in V$
2. $\langle \mathbf{u}, \mathbf{v} \rangle = \langle \mathbf{v}, \mathbf{u} \rangle \quad \forall \mathbf{u}, \mathbf{v} \in V$
3. $\langle \mathbf{u}, \mathbf{u} \rangle \geq 0 \quad \forall \mathbf{u} \in V$, with equality iff $\mathbf{u} = 0$.

Example 4. If \mathbf{A} is a symmetric, positive definite matrix, then $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^T \mathbf{A} \mathbf{v}$ is a valid inner product on \mathbb{R}^d .

However, we need not be restricted to Euclidean inner product spaces.

We would like to know when it is possible to write a function k as

$$k(\mathbf{u}, \mathbf{v}) = \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle$$

for some inner product space V and $\Phi : \mathbb{R}^d \rightarrow V$. This leads to the following definition.

Definition 1. We say $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is an inner product kernel if \exists an inner product space V and a feature map $\Phi : \mathbb{R}^d \rightarrow V$ such that

$$k(\mathbf{u}, \mathbf{v}) = \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d$$

Note that Φ and V are generally not unique.

Example 5. Consider the function

$$k(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{u} - \mathbf{v}\|^2\right),$$

where $\sigma > 0$. Then k is an inner product kernel known as the *Gaussian* kernel. Proof of this fact is left as an exercise. The feature space associated to the Gaussian kernel turns out to be infinite dimensional and is difficult to describe. Nonetheless we can gain some intuition as to what kernel methods that use the Gaussian kernel do, and we will see examples of this later.

2.2 Symmetric, Positive Definite Kernels

One way to determine an inner product kernel is to construct Φ explicitly as we did in the examples above. We can also verify that k is an inner product kernel if it satisfies the following properties.

Definition 2. Let $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$. We say k is symmetric if $k(\mathbf{u}, \mathbf{v}) = k(\mathbf{v}, \mathbf{u}) \quad \forall \mathbf{u}, \mathbf{v}$. We say k is positive definite if the $n \times n$ matrix

$$[k(\mathbf{x}_i, \mathbf{x}_j)]_{ij} \tag{1}$$

is positive semi-definite² for all $n \geq 1$ and all $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$.

If k is both symmetric and positive definite, it is referred to as a symmetric, positive definite (SPD) kernel. The matrix in (1) is known as the *kernel matrix*.

Theorem 1. k is an SPD kernel $\iff k$ is an inner product kernel.

You are asked to prove one direction of the above equivalence as an exercise. The other direction is more advanced and will not be covered in this course.

²This is not a typo. The kernel is said to be positive definite when the associated kernel matrix is always positive semi-definite.

2.3 Examples of kernels

A summary of the above kernels is given below:

1. Homogeneous polynomial kernel

$$k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v})^p,$$

2. Inhomogeneous polynomial

$$k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v} + c)^p, \quad c > 0$$

3. Gaussian kernel

$$k(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{u} - \mathbf{v}\|^2\right), \quad \sigma > 0$$

Another example is the *neural tangent kernel*, which has strong connections to neural networks.

Kernels have also been defined on non-Euclidean feature spaces. For example, kernels have been constructed for strings, trees, graphs, images, functions, and probability distributions. This makes it possible to apply any kernelizable algorithm to data of all of these types. This is because in kernel methods, we only need to be able to compute $k(\mathbf{x}, \mathbf{x}')$, and never have to work directly with the raw features.

3 Kernel Methods

A machine learning algorithm is said to be *kernelizable* if it is possible to formulate the algorithm such that all training examples \mathbf{x}_i and all test instances \mathbf{x} occur in inner products of the form $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$ or $\langle \mathbf{x}_i, \mathbf{x} \rangle$. Many algorithms are kernelizable, but in most cases they must be rewritten in an equivalent form to make this obvious.

To *kernelize* a kernelizable algorithm is to select a kernel k , and replace every occurrence of $\langle \mathbf{u}, \mathbf{v} \rangle$, where $\mathbf{u}, \mathbf{v} \in \{\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}\}$ and \mathbf{x} denotes a test point at which a learned model is evaluated, by $k(\mathbf{u}, \mathbf{v})$. By the inner product characterization of kernels, this is equivalent to replacing every occurrence of $\langle \mathbf{u}, \mathbf{v} \rangle$, where $\mathbf{u}, \mathbf{v} \in \{\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}\}$, by $\langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle$. Therefore, kernelizing an algorithm is equivalent to first applying the nonlinear feature map Φ to elements of the feature space, and then applying the original algorithm to the transformed features. As we saw above, this is an effective way of converting linear algorithms to nonlinear ones. Furthermore, since the algorithm is expressed entirely in terms of the kernel, the feature map does not need to be computed.

Example 6. Consider the nearest centroid classifier. Given training data $(\mathbf{x}_i, y_i), i = 1, \dots, n, y_i \in \{-1, +1\}$, define the centroids

$$\mathbf{m}_- = \frac{1}{n_-} \sum_{i:y_i=-1} \mathbf{x}_i \quad \mathbf{m}_+ = \frac{1}{n_+} \sum_{i:y_i=+1} \mathbf{x}_i$$

where $n_- = |\{i : y_i = -1\}|$ and $n_+ = |\{i : y_i = +1\}|$. The nearest centroid classifier is

$$f(\mathbf{x}) = \text{sign}\{\|\mathbf{x} - \mathbf{m}_-\|^2 - \|\mathbf{x} - \mathbf{m}_+\|^2\}.$$

In words, we assign to a test point \mathbf{x} the label associated with the nearest of the two centroids. With a little algebra it can be shown that

$$f(\mathbf{x}) = \text{sign}\{\mathbf{w}^T \mathbf{x} + b\}$$

where $\mathbf{w} = \mathbf{m}_+ - \mathbf{m}_-$ and $b = \frac{\|\mathbf{m}_-\|^2 - \|\mathbf{m}_+\|^2}{2}$. To show that this classifier is kernelizable, we need to express $f(\mathbf{x})$ such that the training feature vectors \mathbf{x}_i and the test vector \mathbf{x} only appear in terms of the form $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$, $\langle \mathbf{x}_i, \mathbf{x} \rangle$, or $\langle \mathbf{x}, \mathbf{x} \rangle$. To see this simply substitute in the definitions of the centroids which yields

$$f(\mathbf{x}) = \text{sign} \left\{ \frac{1}{n_+} \sum_{i:y_i=1} \langle \mathbf{x}_i, \mathbf{x} \rangle - \frac{1}{n_-} \sum_{i:y_i=-1} \langle \mathbf{x}_i, \mathbf{x} \rangle + \frac{1}{2} \left(\frac{1}{n_-^2} \sum_{\substack{i:y_i=-1 \\ j:y_j=-1}} \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \frac{1}{n_+^2} \sum_{\substack{i:y_i=1 \\ j:y_j=1}} \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right) \right\}.$$

To kernelize the nearest centroid classifier, we need to choose a kernel k , such as the Gaussian kernel. Then the new classifier, which we might call the *kernel nearest centroid classifier*, is

$$f(\mathbf{x}) = \text{sign} \left\{ \frac{1}{n_+} \sum_{i:y_i=1} k(\mathbf{x}_i, \mathbf{x}) - \frac{1}{n_-} \sum_{i:y_i=-1} k(\mathbf{x}_i, \mathbf{x}) + \frac{1}{2} \left(\frac{1}{n_-^2} \sum_{\substack{i:y_i=-1 \\ j:y_j=-1}} k(\mathbf{x}_i, \mathbf{x}_j) - \frac{1}{n_+^2} \sum_{\substack{i:y_i=1 \\ j:y_j=1}} k(\mathbf{x}_i, \mathbf{x}_j) \right) \right\}.$$

Many kernels, such as the Gaussian kernel, can be viewed as measures of similarity between data points: The closer \mathbf{x} and \mathbf{x}' , the larger $k(\mathbf{x}, \mathbf{x}')$. If we think of $k(\mathbf{x}, \mathbf{x}')$ as a measure of similarity between \mathbf{x} and \mathbf{x}' then, ignoring the bias term, this classifier assigns the label of the class to which the test instance has higher average similarity. The bias term can be interpreted as favoring whichever class has higher average within-class similarity. In the case of the Gaussian kernel, normalized to integrate to one, we can actually say more. The first two terms in the final classifier are actually estimates of their respective class-conditional densities (kernel density estimates, to be precise). If we ignore the bias, the classifier assigns a test point to the class for which it has higher estimated density. Kernel density estimates will be studied in a later chapter.

Although the bias term looks cumbersome, it does not involve the test instance and can therefore be precomputed in $O(n^2d)$ operations. Then, for each test instance, the classifier can be evaluated in $O(nd)$ operations. This tends to be true for kernel methods in general: the kernelized version of an algorithm tends to have greater computational complexity than the original by one or two factors of n . This can still be much more efficient than computing $\Phi(\mathbf{x})$ explicitly when m is large.

Remark 1. The process of kernelizing a kernelizable algorithm is also known as *the kernel trick*.

Exercises

1. (☆) Use properties of Vandermonde matrices to prove that if the data points are unique, then in polynomial regression with $p < n$, the matrix $\mathbf{X}^T \mathbf{X}$ is invertible.
2. (★) If the smallest eigenvalue of a matrix is too small (relative to the largest eigenvalue), the matrix is said to be *ill-conditioned*. This means that while it may be invertible, computation of the inverse is prone to large numerical errors. In polynomial regression with p sufficiently large relative to n , the matrix $\mathbf{X}^T \mathbf{X}$ is ill-conditioned. Show that regularization overcomes this issue. In particular, show that adding a quadratic regularizer $\lambda \|\boldsymbol{\theta}\|^2$ to the least-squares objective makes the smallest eigenvalue (of the matrix that needs to be inverted) at least λ .
3. (★★) Let $k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^T \mathbf{v} + 1)^3$ where $\mathbf{u}, \mathbf{v} \in \mathbb{R}^2$. Find Φ such that $k(\mathbf{u}, \mathbf{v}) = \langle \Phi(\mathbf{u}), \Phi(\mathbf{v}) \rangle$.
4. (☆☆) Extend the previous example to arbitrary d and p .
5. (☆☆) Show that the number of monomials of degree p in d variables is $\binom{d+p-1}{p}$.
6. (☆) Use the previous result to show that the number of monomials of degree at most p in d variables is $\binom{d+p}{p}$.
7. (☆☆) Show that the nearest centroid classifier is equivalent to LDA in the situation where the class priors are equal and the covariance matrix is a multiple of the identity.
8. (★★) Kernelize the k -nearest neighbor classifier. Show that if $k = 1$ and the kernel is the Gaussian kernel, then the original and kernelized classifiers are the same.
9. (★★) Show that the sum of two kernels is a kernel.
10. (★★★) This problem asks you to kernelize logistic regression by kernelizing the Newton-Raphson algorithm for logistic regression. Since kernelizing corresponds to a high dimensional feature space, we will consider the regularized form of logistic regression. To simplify the problem somewhat, let's also assume that there is no offset term, so that the class probability model is

$$\eta(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})} = g(\boldsymbol{\theta}^T \mathbf{x}),$$

where $g(t) = 1/(1 + e^{-t})$. Let $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ be training data for binary classification, where $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \{0, 1\}$.

From an earlier exercise on logistic regression (and making a slight modification to account for the lack of an offset term), the Newton-Raphson update has the following form:

$$\boldsymbol{\theta}^{new} = \boldsymbol{\theta} + (\mathbf{X}^T \mathbf{W} \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{v} - \lambda \boldsymbol{\theta}),$$

where \mathbf{X} is an $n \times d$ matrix whose i th row is \mathbf{x}_i^T , \mathbf{W} is a diagonal matrix with i th diagonal entry equal to $g(\boldsymbol{\theta}^T \mathbf{x}_i)(1 - g(\boldsymbol{\theta}^T \mathbf{x}_i))$, \mathbf{v} is an $n \times 1$ column vector with i th entry given by $y_i - g(\boldsymbol{\theta}^T \mathbf{x}_i)$, and $\lambda > 0$ is the regularization parameter.

- (a) The matrix inversion lemma states that

$$(\mathbf{P} + \mathbf{Q} \mathbf{R} \mathbf{S})^{-1} = \mathbf{P}^{-1} - \mathbf{P}^{-1} \mathbf{Q} (\mathbf{R}^{-1} + \mathbf{S} \mathbf{P}^{-1} \mathbf{Q})^{-1} \mathbf{S} \mathbf{P}^{-1},$$

provided all of the matrix operations are well defined. Re-express the logistic regression iterative update step by applying the matrix inversion lemma to $(\mathbf{X}^T \mathbf{W} \mathbf{X} + \lambda \mathbf{I})^{-1}$.

We are free to initialize the algorithm in any way we choose. Let us initialize θ so that $\theta = \mathbf{X}^T \alpha$ for some vector $\alpha \in \mathbb{R}^n$. In other words, θ is a linear combination of data points. For example, we could initialize θ to be the first data point \mathbf{x}_1 , in which case $\alpha = [1 \ 0 \ 0 \ \cdots \ 0]^T$.

- (b) Now consider an arbitrary step in the iterative procedure. Show that if $\theta = \mathbf{X}^T \alpha$ for some α , then $\theta^{new} = \mathbf{X}^T \alpha^{new}$ for some α^{new} . Give an update formula for α^{new} in terms of α and the other quantities.
 - (c) Explain how to implement kernel logistic regression. That is, show how to incorporate inner product kernels into the iterative parameter update, and also how to evaluate the final classifier.
11. (☆☆☆) In this problem you are asked to show that the Gaussian kernel

$$k(\mathbf{x}, \mathbf{x}') = C \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}'\|^2\right), \quad C > 0, \sigma > 0,$$

is an inner product kernel (equivalently, a symmetric and positive definite kernel). You will be asked to verify various properties that let you construct kernels from other kernels, and you will see that both characterizations of kernels (IP/SPD) are useful in this regard.

- (a) Consider a sequence of symmetric, positive definite (SPD) kernels $\{k_n\}$ that converges pointwise to a function k , i.e., for all \mathbf{x}, \mathbf{x}' , $\lim_{n \rightarrow \infty} k_n(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}')$. Show that k is also a SPD kernel.
- (b) Prove that if $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$ are SPD kernels, then so is $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$.
Hint: Consider zero mean, independent random vectors V_1 and V_2 whose covariances are kernel matrices K_1 and K_2 . Determine another random vector whose covariance is the element-wise product of K_1 and K_2 .
- (c) Let $f(t) = \sum_{n=0}^{\infty} a_n t^n$ be a power series that converges for all $t \in \mathbb{R}$, and for which all $a_n \geq 0$. Argue that $k(\mathbf{x}, \mathbf{x}') = \sum_{n=0}^{\infty} a_n (k(\mathbf{x}, \mathbf{x}'))^n$ is a kernel.
- (d) Deduce that the exponential kernel $k(\mathbf{x}, \mathbf{x}') = \exp(\gamma \langle \mathbf{x}, \mathbf{x}' \rangle)$ is an inner product kernel, where $\gamma > 0$ and $\langle \cdot, \cdot \rangle$ denotes the dot product on \mathbb{R}^d .
- (e) Show that if k is an inner product kernel, then so is the normalized kernel

$$\tilde{k}(\mathbf{x}, \mathbf{x}') = \begin{cases} 0, & \text{if } k(\mathbf{x}, \mathbf{x}) = 0 \text{ or } k(\mathbf{x}', \mathbf{x}') = 0 \\ \frac{k(\mathbf{x}, \mathbf{x}')}{\sqrt{k(\mathbf{x}, \mathbf{x})k(\mathbf{x}', \mathbf{x}')}}, & \text{otherwise.} \end{cases}$$

- (f) Deduce that the Gaussian kernel is an inner product kernel.