

# Chapter 3

## Kernels, Beginning of neural networks

### Reading

1. Bishop 6.1-6.3
2. Goodfellow 6.1-6.4
3. “Random features for large-scale kernel machines” by [Rahimi and Recht \(2008\)](#).

### 3.1 Digging deeper into the perceptron

#### 3.1.1 Convergence rate

How many iterations does a perceptron need to fit on a given dataset? We will assume that the training data are bounded, i.e.,  $\|x^i\| \leq R$  for some  $R$  and for all  $i \in \{1, \dots, n\}$ . Let us also assume that the training dataset is indeed linearly separable, i.e., a solution  $w^*$  exists for the perceptron weights with training error exactly zero. This means

$$y^i w^{*\top} x^i > 0 \quad \forall i.$$

We will also assume that this classifier *separates the data well*. Note that the distance of each input  $x^i$  from the decision boundary (i.e., all  $x$  such that  $w^{*\top} x = 0$ ) is given by the component of  $x^i$  in the direction of  $w^*$  if the label is  $y^* = +1$  and in the direction  $-w^*$  if the label is negative. In other words,

$$\frac{y^i w^{*\top} x^i}{\|w^*\|} = \rho^i$$

gives the distance to the decision boundary. The quantity on the right hand side is called the *margin*, it is simply the distance of the sample  $i$  from the

17 decision boundary. If  $w^*$  is the classifier with the largest average margin,

$$\rho = \min_{i \in \{1, \dots, n\}} \rho^i$$

18 is a good measure of how hard a particular machine learning problem is.

19 You can now try to prove that after each update of the perceptron the inner  
20 product of the current weights with the try solution  $\langle w_t, w^* \rangle$  increases at least  
21 linearly and that the squared norm  $\|w_t\|^2$  increases at most linearly in the  
22 number of updates  $t$ . Together the two will give you a result that after  $t$  weight  
23 updates

$$t \leq \frac{R^2}{\rho^2} \quad (3.1)$$

24 all training data are classified correctly. Notice a few things about this expres-  
25 sion.

- 26 1. The quantity  $\frac{R^2}{\rho^2}$  is dimension independent; that the number of steps  
27 reach a given accuracy is independent of the dimension of the data will  
28 be a property shared by optimization algorithms in general.
- 29 2. There are no constant factors, this is also the worst case number of  
30 updates; this is rare.
- 31 3. The number of updates scales with the hardness of the problem; if the  
32 margin  $\rho$  was small, we need lots of updates to drive the training error  
33 to zero.

### 34 3.1.2 Dual representation

35 Let us see how the parameters of the perceptron look after training on the entire  
36 dataset. At each iteration, the weights are updated in the direction  $(x^t, y^t)$   
37 or they are not updated at all. Therefore, if  $\alpha^i$  is the number of times the  
38 perceptron sampled the datum  $(x^i, y^i)$  during the course of its training and got  
39 it wrong, we can write the weights of the perceptron as a linear combination

$$w^* = \sum_{i=1}^n \alpha^i y^i x^i. \quad (3.2)$$

40 where  $\alpha^i \in \{0, 1, \dots\}$ . The perceptron therefore using the classifier

$$\begin{aligned} f(x, w) &= \text{sign}(\hat{y}) \\ \text{where } \hat{y} &= \left( \sum_{i=1}^n \alpha^i y^i x^i \right)^\top x \\ &= \sum_{i=1}^n \alpha^i y^i x^{i\top} x. \end{aligned} \quad (3.3)$$

41 Remember this special form: the inner product of the new input  $x$  with  
42 all the other inputs  $x^i$  in the training dataset is combined linearly to get the  
43 prediction. The weights of this linear combination are the dual variables which  
44 is a measure of how many tries it took the perceptron to fit that sample during  
45 training.  
46

▲ As you see in (3.3), computing the prediction for a new input  $x$  involves, either remembering all the weights  $w$  at the end of training, or storing all the  $\{\alpha^i\}_{i=1, \dots, n}$  along with the training dataset. The latter is called the dual representation of a perceptron and the scalars  $\{\alpha^i\}$  are called the dual parameters.

## 3.2 Creating nonlinear classifiers from linear ones

Linear classifiers such as the perceptron, or the support vector machine (SVM) can be extended to nonlinear ones. The trick is essentially the same that we saw when we fit polynomials (polynomials are nonlinear) using the formula for linear regression. We are interested in mapping input data  $x$  to some different space, this is usually a higher-dimensional space called the *feature space*.

$$x \mapsto \phi(x).$$

The quantity  $\phi(x)$  is called a feature vector.

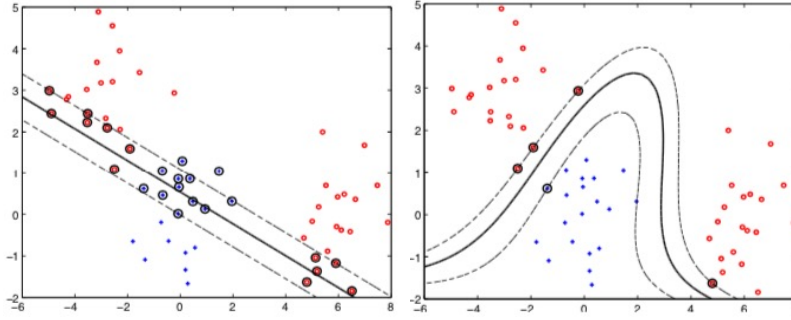


Figure 3.1

For example, in the polynomial regression case for scalar input data  $x \in \mathbb{R}$  we used

$$\phi(x) := \begin{bmatrix} 1, \sqrt{2}x, x^2 \end{bmatrix}^\top$$

to get a quadratic feature space. The role of  $\sqrt{2}$  will become clear shortly. Certainly this trick of polynomial expansion also works for higher dimensional input

$$\phi(x) := \begin{bmatrix} 1, x_1, x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2 \end{bmatrix}^\top.$$

Having fixed a feature vector  $\phi(x)$ , we can now fit a linear perceptron on the input data  $\{\phi(x^i), y^i\}$ . This involves updating the weights at each iteration as

$$w_{t+1} = \begin{cases} w_t + y^t \phi(x^t) & \text{if } \text{sign}(w_t^\top \phi(x^t)) \neq y^t \\ w_t & \text{else.} \end{cases} \quad (3.4)$$

At the end of such training, the perceptron is

$$w^* = \sum_{i=1}^n \alpha^i y^i \phi(x^i)$$

and predictions are made by first mapping the new input to our feature space

$$f(x; w) = \text{sign} \left( \sum_{i=1}^n \alpha^i y^i \phi(x^i)^\top \phi(x) \right). \quad (3.5)$$

Notice that we now have a linear combination of the *features* not the data directly.

❓ The concept of a feature space seems like a panacea. If we have complex data, we simply map it to some high-dimensional feature and fit a linear function to these features. However, the “curse of dimensionality” coined by Richard Bellman states that to fit a function in  $\mathbb{R}^d$  the number of data needs to be exponential in  $d$ . It therefore stands to reason that we need a lot more data to fit a classifier in feature space than in the original input space. Why would we still be interested in the feature space then?

### 3.3 Kernels

Observe the expression of the classifier in (3.5). Each time we make predictions on the new input, we need to compute  $n$  inner products of the form

$$\phi(x^i)^\top \phi(x).$$

If the feature dimension is high, we need to enumerate the large number of feature dimensions if we are using the weights of the perceptron, or these inner products if we are using the dual variables. Observe however that even if the feature vector is large, we can compactly evaluate the inner product

$$\begin{aligned}\phi(x) &= [1, \sqrt{2}x, x^2] \\ \phi(x') &= [1, \sqrt{2}x', x'^2] \\ \phi(x)^\top \phi(x') &= 1 + 2xx' + (xx')^2 = (1 + xx')^2.\end{aligned}$$

for input  $x \in \mathbb{R}$ . Kernels are a formalization of exactly this idea. A kernel

$$k : X \times X \rightarrow \mathbb{R}.$$

is a symmetric, positive semi-definite function of its two arguments for which it holds that

$$k(x, x') = \phi(x)^\top \phi(x')$$

for some feature  $\phi$ . Few examples of kernels are

$$\begin{aligned}k(x, x') &= (x^\top x' + c)^2, \\ k(x, x') &= \exp(-\|x - x'\|^2 / (2\sigma^2)).\end{aligned}$$

#### 3.3.1 Kernel perceptron

We can now give the kernel version of the perceptron algorithm. The idea is to simply replace any inner product in the algorithm that looks like  $\phi(x)^\top \phi(x')$  by the kernel  $k(x, x')$ .

##### Kernel perceptron

Initialize dual variables  $\alpha^i = 0$  for all  $i \in \{1, \dots, n\}$ . Perform the following steps for iterations  $t = 1, 2, \dots$

1. At the  $t^{\text{th}}$  iteration, sample a data point with index  $\omega_t$  from  $D_{\text{train}}$  uniformly randomly, call it  $(x^{\omega_t}, y^{\omega_t})$ .
2. If there is a mistake, i.e., if

$$\begin{aligned}0 &\geq y^{\omega_t} \left( \sum_{i=1}^n \alpha^i y^i \phi(x^i)^\top \phi(x^{\omega_t}) \right) \\ &= y^{\omega_t} \left( \sum_{i=1}^n \alpha^i y^i k(x^i, x^{\omega_t}) \right),\end{aligned}$$

then update

$$\alpha^{\omega_t} \leftarrow \alpha^{\omega_t} + 1.$$

**▲** Feature spaces can become large very quickly. What is the dimensionality of  $\phi(x)$  for a tenth-order polynomial with a three-dimensional input data?

Notice that we do not ever compute  $\phi(x)$  so it does not matter what the dimensionality of the feature vector is. It can even be infinite, e.g., for the radial basis function kernel. Observe also that we do not maintain weights  $w$ . We instead maintain the dual variables  $\{\alpha^1, \dots, \alpha^n\}$  while running the algorithm.

Note that the kernel perceptron computes the kernel over *all* data samples in the training set at each iteration. It is expensive and seems wasteful. The Gram matrix denoted by  $G \in \mathbb{R}^{n \times n}$

$$G_{ij} = k(x^i, x^j) \quad (3.6)$$

helps address this problem by computing the kernel on all pairs in the training dataset. With this in hand, we can modify step 2 in the kernel perceptron using

$$y^t \left( \sum_{i=1}^n \alpha^i y^i k(x^i, x^t) \right) = y^t (\alpha \odot Y)^T G e_t.$$

where  $e_t = [0, \dots, 0, 1, 0, \dots]$  with a 1 on the  $t^{\text{th}}$  element,  $\alpha = [\alpha^1, \dots, \alpha^n]$  denotes the vector of all the dual variables,  $Y = [y^1, \dots, y^n]$  is a vector of all the labels, and the notation  $\alpha \odot Y = [\alpha^1 y^1, \dots, \alpha^n y^n]$  denotes the element-wise (Hadamard) product. This expression now only involves a matrix-vector multiplication, which is much easier than computing the kernel at each iteration. Gram matrices can become very big. If the number of samples is  $n = 10^6$ , not an unusual number today, the Gram matrix has  $10^{12}$  elements. The big failing of kernel methods is that they require a large amount of memory at training time. Nystrom methods compute low-rank approximations of the Gram matrix which makes operations with kernels easier.

### 3.3.2 Mercer's theorem

This theorem shows that given any kernel that satisfies some regularity properties can be rewritten as an inner product.

**Theorem 1 (Mercer's Theorem).** For any symmetric function  $k : X \times X \rightarrow \mathbb{R}$  which is square integrable in  $X \times X$  and satisfies

$$\int_{X \times X} k(x, x') f(x) f(x') dx dx' \geq 0 \quad (3.7)$$

for all square integrable functions  $f \in L_2(X)$ , there exist functions  $\phi_i : X \rightarrow \mathbb{R}$  and numbers  $\lambda_i \geq 0$  where

$$k(x, x') = \sum_i \lambda_i \phi_i^\top(x) \phi_i(x')$$

for all  $x, x' \in X$ . The condition in (3.7) is called Mercer's condition. You will also have seen it written as *for any finite set of inputs  $\{x^1, \dots, x^n\}$  and any choice of real-valued coefficients  $c_1, \dots, c_n$  a valid kernel should satisfy*

$$\sum_{i,j} c_i c_j k(x^i, x^j) \geq 0.$$

There can be an infinite number of coefficients  $\lambda_i$  in the summation.

❓ Kernels look great, you can fit perceptrons in powerful feature spaces using essentially the same algorithm. How expensive is each iteration of the perceptron?

▲ When ML algorithms are implemented in a system, there exist tradeoffs between the feature-space version and the Gram matrix version of linear classifiers. The former is preferable if the number of samples in the dataset is large, while the latter is used when the dimensionality of features is large.

❓ Logistic regression with a loss function

$$\ell_{\text{logistic}}(w) = \log(1 + e^{-yw^\top x})$$

is also a linear classifier. Write down how you will fit a logistic regression using stochastic gradient descent; this is similar to the perceptron algorithm. Write down the feature-space version of the algorithm and a kernelized logistic regression that uses the Gram matrix.

💡 A function  $f : X \rightarrow \mathbb{R}$  is square integrable iff

$$\int_{x \in X} |f(x)|^2 dx < \infty.$$

**Remark 2 (Checking if a function is a valid kernel).** Note that Mercer's condition states that the Gram matrix of any dataset is positive semi-definite:

$$u^\top G u \geq 0 \quad \text{for all } u \in \mathbb{R}^n.$$

This is easy to show.

$$\begin{aligned} u^\top G u &= \sum_{ij} u_i u_j G_{ij} \\ &= \sum_{ij} u_i u_j \phi(x_i)^\top \phi(x_j) \\ &= \left( \sum_i u_i \phi(x_i) \right)^\top \left( \sum_j u_j \phi(x_j) \right) \\ &= \left\| \sum_i u_i \phi(x_i) \right\|^2 \\ &\geq 0. \end{aligned}$$

The integral in Theorem 1 in Mercer's condition is really just the continuous analogue of the vector-matrix-vector multiplication above. So if you have a function that you would like to use as a kernel, checking its validity is easy by showing that the Gram matrix is positive semi-definite.

Kernels are powerful because they do not require you to think of the feature and parameter spaces. For instance, we may wish to design a machine learning algorithm for spam detection that takes in a variable length of feature vector depending on the particular input. If  $x[i]$  is the  $i^{\text{th}}$  character of a string, a good feature vector to use is to consider the set of all length  $k$  subsequences. The number of components in this feature vector is exponential. However, as you can imagine, given two strings  $x, x'$

```
this string is interesting
txws sbhtqg is iyubqtnhpqg
```

you can write a Python function to check their similarities with respect to some rules *you define*. Mercer's theorem is useful here because it says that so long as your function satisfies the basic properties of a kernel function, there exists some feature space which your function implicitly constructs.

### 3.4 Learning the feature vector

The central idea behind deep learning is to learn the feature vectors  $\phi$  instead of choosing them a priori.

How do we choose what set of feature vectors to learn from? For instance, we can pick all polynomials; we can pick all possible Gabor filters that you saw in HW 1; we can also pick all possible string kernels.

### 3.4.1 Random features

Suppose that we have a finite-dimensional feature  $\phi(x) \in \mathbb{R}^p$ . We saw in the perceptron that

$$f(x; w) = \text{sign} \left( \sum_i w_i \phi_i(x) \right)$$

where  $\phi(x) = [\phi_1(x), \dots, \phi_p(x)]$  and  $w = [w_1, \dots, w_p]$  are the feature and weight vectors respectively. We will set

$$\phi(x) = \sigma(S^\top x), \quad (3.8)$$

where  $S \in \mathbb{R}^{d \times p}$  is a matrix. The function  $\sigma(\cdot)$  is a nonlinear function of its argument and acts on all elements of the argument element-wise

$$\sigma(z) = [\sigma(z_1), \dots, \sigma(z_p)]^\top.$$

We will abuse notation that denote both the vector version of  $\sigma$  and the element-wise version of  $\sigma$  using the same Greek letter. Notice that this is a special type of feature vector (or a special type of kernel), it is a linear combination of the input elements. What matrix  $S$  should we pick to combine these input elements? The paper by Rahimi and Recht (2008) proposed the idea that for shift-invariant kernels (which have the property  $k(x, x') = k(x - x')$  one may use a matrix with random elements as our  $S$

$$S^\top = \begin{bmatrix} \omega_1 \\ \vdots \\ \omega_p \end{bmatrix}$$

where  $\omega_i \in \mathbb{R}^d$  are random variables drawn from, say, a Gaussian distribution and the function

$$\sigma(z) = \cos(z)$$

is a cosine function. Using a random matrix is a cheap trick, it lets us create a lot of features quickly without worrying about their quality. Our classifier is now

$$f(x; w) = \text{sign}(w^\top \sigma(S^\top x)) \quad (3.9)$$

and we can solve the optimization problem

$$w^* = \underset{w}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n \ell_{\text{hinge}}(y^i, \hat{y}^i) \quad (3.10)$$

with  $\hat{y}^i = w^\top \sigma(S^\top x^i)$  and fit the weights  $w$  using SGD as before.



Figure 3.2

As an example consider the heatmap of Gabor-like kernel  $k(x, x')$  in Figure 3.2 on the left. We can think of the decomposition

$$\begin{aligned} \text{left-most picture} &= k(x, x') = \phi(x)^\top \phi(x') \\ &= w_1 \sigma(\omega_1^\top x) + w_2 \sigma(\omega_2^\top x) + w_k \sigma(\omega_k^\top x) + \dots \\ &= \text{right-most picture} \end{aligned}$$

In other words, the random elements of the matrix  $S$ , namely  $\omega_k$  can combine together *linearly* using the trained weights  $w_k$  to give us a kernel that looks like a useful kernel on the left. A large random matrix  $S$  allows us to learn may such kernels and combine their output linearly.

### 3.4.2 Learning the feature matrix as well

Random features do not work easily for all kinds of data. For instance, if you have an image of size  $100 \times 100$ , and you are trying to find a fruit



Features:

$\mathbb{1}_{x_{ij} = \text{red}} \quad \forall i, j$

$\mathbb{1}_{\text{circle of radius 5 at } x_{ij} \text{ has red pixels}} \quad \forall i, j$

In order to get a decent training error, you would need a lot of these random features

❓ What kind of data do you think random features will work well for?

we can design random features of the form

$$\phi_{ij,kl} = \mathbb{1}_{\{\text{mostly red color in a box formed by pixels } (ij) \text{ and } (kl)\}}.$$

We will need lots and lots of such features before we can design an object detector that works well for this image. In other words, random features do not solve the problem that you need to be clever about picking your feature space/kernel.

We can now simply motivate deep learning as learning the matrix  $S$  in (3.9) in addition to the coefficients  $w$ . The classifier now is

$$f(x; w, S) = \text{sign}(w^\top \sigma(S^\top x)) \quad (3.11)$$

but we now solve the optimization problem

$$w^*, S^* = \underset{w, S}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n \ell_{\text{hinge}}(y^i, \hat{y}^i) \quad (3.12)$$

with  $\hat{y}^i = w^\top \sigma(S^\top x^i)$  as before. We have hereby seen our first *deep network*. The classifier in (3.11) is a two-layer neural network.

Moving from the problem in (3.10) to this new problem in (3.12) is a very big change.

1. **Nonlinearity.** The classifier in (3.11) is not linear anymore. It is a nonlinear function of its parameters  $w, S$  (both of which we will call weights).



- 
- 181 2. **High-dimensionality.** We added a lot more weights to the classifier,  
182 the original classifier had  $w \in \mathbb{R}^p$  parameters to learn while the new  
183 one also has  $S \in \mathbb{R}^{d \times p}$  more weights. The curse of dimensionality  
184 suggests that we will need a lot more data to fit the new classifier.
- 185 3. **Non-convex optimization.** The optimization problem in (3.12) much  
186 harder than the one in (3.10). The latter is a convex function (we  
187 will discuss this soon) which are easy to minimize. The former is  
188 a non-convex function in its parameters  $w, S$  because they interact  
189 multiplicatively, such functions are harder to minimize. We could write  
190 down the solution of the perceptron using the final values of the dual  
191 variables. We cannot do this for a two-layer neural network.

# 192 Bibliography

- 193 Rahimi, A. and Recht, B. (2008). Random features for large-scale kernel  
194 machines. In *Advances in neural information processing systems*, pages  
195 1177–1184.