

**Features**

Daniel Hsu

*Warning: notes are incomplete and unchecked for correctness***1 Feature expansions**

In many applications, no linear classifier over the “raw” set of features will perfectly separate the data. One recourse is to find additional features that are predictive of the label. This is called *feature engineering*, and is often a substantial part of the job of a machine learning practitioner. In some applications, it is possible to “throw in the kitchen sink”—i.e., include all possible features that might possibly be relevant. For instance, in document classification, one can include a feature for each possible word in the vocabulary that indicates whether that word is present in the given document (or counts the number of occurrences). One can also include a feature for each possible pair of consecutive words (“bi-grams”), each possible triple of consecutive words (“tri-grams”), and so on. In general, it is common to automatically generate features based on existing features  $\mathbf{x} \in \mathbb{R}^d$ , such as quadratic interaction features  $\mathbf{x} \mapsto (x_1x_2, x_1x_3, \dots, x_1x_d, x_2x_3, \dots, x_{d-1}x_d) \in \mathbb{R}^{\binom{d}{2}}$ , as well as higher-order interaction features.

The main drawback of these “kitchen sink” feature expansions is that it may be computationally expensive to work explicitly in the expanded feature space. Fortunately, there are some ways around this.

**1.1 The kernel trick**

Suppose  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$  is the mapping of the base feature vector  $\mathbf{x}$  to some expanded feature representation  $\phi(\mathbf{x})$  (where possibly  $D \gg d$ ). Imagine running the Online Perceptron algorithm in this expanded feature space  $\mathbb{R}^D$ . Naïvely, it takes  $\Theta(D)$  time to compute  $\langle \phi(\mathbf{x}), \mathbf{w} \rangle$  by explicitly working with the expanded feature representation. However, in some cases, this can be greatly improved. Recall that at any point in the algorithm, Online Perceptron’s weight vector is equal to

$$\mathbf{w} = \sum_{(\mathbf{x}, y) \in \mathcal{M}} y \phi(\mathbf{x})$$

where  $\mathcal{M}$  is the set of examples on which Online Perceptron has so far made mistakes (and hence used to update the weight vector). Therefore,

to compute a prediction on a new point  $\mathbf{x}'$  (with feature expansion  $\phi(\mathbf{x}')$ ), we need to compute the inner product

$$\langle \phi(\mathbf{x}'), \mathbf{w} \rangle = \sum_{(\mathbf{x}, y) \in \mathcal{M}} y \langle \phi(\mathbf{x}'), \phi(\mathbf{x}) \rangle.$$

This can be done in time  $\Theta(|\mathcal{M}| \cdot \tau_K)$ , where  $\tau_K$  is the time to compute the inner product  $\langle \phi(\mathbf{x}'), \phi(\mathbf{x}) \rangle$ . Naïvely,  $\tau_K$  appears to be  $\Theta(D)$ , so this may not appear to be an improvement. But fortunately, in many cases of interest,  $\tau_K$  is much smaller. This is called the *kernel trick*.

**All degree  $\leq 2$  interaction features.** Consider all  $D := 1 + 2d + \binom{d}{2}$  products of at most two features:

$$\begin{aligned} \phi(\mathbf{x}) := & \left( 1, \sqrt{2}x_1, \sqrt{2}x_2, \dots, \sqrt{2}x_d, \right. \\ & \left. x_1^2, x_2^2, \dots, x_d^2, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \dots, \sqrt{2}x_1x_d, \sqrt{2}x_2x_3, \dots, \sqrt{2}x_{d-1}x_d \right). \end{aligned}$$

There is some seemingly arbitrary scaling by  $\sqrt{2}$  in some parts, but this is not a critical problem for linear classifiers. We can compute  $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$  in  $\tau_K = O(d)$  time:

$$\begin{aligned} (1 + \langle \mathbf{x}, \mathbf{x}' \rangle)^2 &= 1 + 2\langle \mathbf{x}, \mathbf{x}' \rangle + \langle \mathbf{x}, \mathbf{x}' \rangle^2 \\ &= 1 + \sum_{i=1}^d (\sqrt{2}x_i)(\sqrt{2}x'_i) + \sum_{i=1}^d x_i^2 x_i'^2 + \sum_{1 \leq i < j \leq d} (\sqrt{2}x_i x_j)(\sqrt{2}x'_i x'_j) \\ &= \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle. \end{aligned}$$

**Products of all feature subsets.** Consider all  $D := 2^d$  products of subsets of features:

$$\phi(\mathbf{x}) := \left( \prod_{i \in S} x_i : S \subseteq [d] \right).$$

We can compute  $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$  in  $\tau_K = O(d)$  time:

$$\begin{aligned} \prod_{i=1}^d (1 + x_i x'_i) &= \sum_{S \subseteq [d]} \prod_{i \in S} (x_i x'_i) \\ &= \sum_{S \subseteq [d]} \prod_{i \in S} x_i \prod_{i \in S} x'_i \\ &= \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle. \end{aligned}$$

**An infinite dimensional feature expansion.** For any  $\sigma > 0$ , there is some infinite feature expansion  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^\infty$  (involving Hermite polynomials of all orders) such that

$$\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right),$$

which can be computed in  $\tau_K = O(d)$  time. This inner product is called the *Gaussian kernel with bandwidth  $\sigma$* .

**Kernels.** A positive definite *kernel function*  $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a symmetric function that satisfies the following property: for any  $n \geq 1$ , any  $c_1, c_2, \dots, c_n \in \mathbb{R}$ , and any  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$ ,

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0.$$

That is, for any set of data points  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ , the symmetric matrix whose  $(i, j)$ -th entry is  $K(\mathbf{x}_i, \mathbf{x}_j)$  is positive semi-definite. For any positive definite kernel  $K$ , there is some feature mapping  $\phi$  such that  $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = K(\mathbf{x}, \mathbf{x}')$  for any pair  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ . The mapping is not necessarily to  $\mathbb{R}^D$  for any finite  $D$ , but rather to an inner product space called a *Reproducing Kernel Hilbert Space* (RKHS). Also, kernel functions need not be defined over Euclidean spaces  $\mathbb{R}^d$ ; they can in fact be (and have been) defined over general sets (e.g., strings, trees, genetic sequences).

## 2 Kernel approximation

A drawback of using the kernel trick with Perceptron is that the representation of the weight vector  $\mathbf{w} \in \mathbb{R}^D$  may grow with the number of the training examples. One way to alleviate this drawback is to limit the number of data points used in the representation of the weight vector  $\mathbf{w} \in \mathbb{R}^D$ .

An alternative is to approximate the kernel function directly by finding a mapping  $\mathbf{z}: \mathbb{R}^d \rightarrow \mathbb{R}^m$  such that

$$\langle \mathbf{z}(\mathbf{x}), \mathbf{z}(\mathbf{x}') \rangle \approx K(\mathbf{x}, \mathbf{x}').$$

This can often be done with  $m \ll D$  (especially when  $D = \infty$ ) using *randomized feature expansions*.

**Shift-invariant kernels.** A very popular class of kernels are shift-invariant kernels: positive definite kernels of the form  $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x} - \mathbf{x}', \mathbf{0}) =: K(\mathbf{x} - \mathbf{x}')$ . This includes Gaussian kernels, Laplacian kernels, Cauchy kernels, and many others. The randomized approximation is grounded in *Bochner's theorem*, which states that for any continuous shift-invariant kernel  $K: \mathbb{R}^d \rightarrow \mathbb{R}$  on  $\mathbb{R}^d$ , then for some non-negative function  $\lambda: \mathbb{R}^d \rightarrow \mathbb{R}_+$ ,

$$K(\boldsymbol{\delta}) = \int_{\mathbb{R}^d} \exp(i\langle \boldsymbol{\delta}, \boldsymbol{\omega} \rangle) \lambda(\boldsymbol{\omega}) d\boldsymbol{\omega}$$

where  $i = \sqrt{-1}$ . This says that  $K$  is the inverse Fourier transform of a non-negative function  $\lambda$ . Given any continuous shift-invariant kernel  $K$ , we can obtain  $\lambda$  via the Fourier transform:

$$\lambda(\boldsymbol{\omega}) = \int_{\mathbb{R}^d} \exp(-i\langle \boldsymbol{\delta}, \boldsymbol{\omega} \rangle) K(\boldsymbol{\delta}) d\boldsymbol{\delta}.$$

By appropriately scaling the kernel, the function  $\lambda$  above will be a probability density.

For example, if  $K(\boldsymbol{\delta}) = \exp(-\|\boldsymbol{\delta}\|^2/(2\sigma^2))/(2\pi)^d$  is the (scaled) Gaussian kernel, then

$$\lambda(\boldsymbol{\omega}) = \frac{1}{(2\pi\sigma^{-2})^{d/2}} \exp\left(-\frac{\|\boldsymbol{\omega}\|^2}{2\sigma^{-2}}\right),$$

which is simply the multivariate Gaussian density with spherical covariance  $\sigma^{-2}\mathbf{I}$ . Therefore, we can write the (scaled) Gaussian kernel as the following expectation over a Gaussian random vector  $\boldsymbol{\omega} \sim \mathcal{N}(\mathbf{0}, \sigma^{-2}\mathbf{I})$ :

$$\begin{aligned} K(\mathbf{x} - \mathbf{x}') &= \mathbb{E}(\exp(i\langle \mathbf{x} - \mathbf{x}', \boldsymbol{\omega} \rangle)) \\ &= \mathbb{E}(\exp(i\langle \mathbf{x}, \boldsymbol{\omega} \rangle) \exp(-i\langle \mathbf{x}', \boldsymbol{\omega} \rangle)) \\ &= \mathbb{E}(\cos(\langle \mathbf{x}, \boldsymbol{\omega} \rangle) \cos(\langle \mathbf{x}', \boldsymbol{\omega} \rangle) + \sin(\langle \mathbf{x}, \boldsymbol{\omega} \rangle) \sin(\langle \mathbf{x}', \boldsymbol{\omega} \rangle)). \end{aligned}$$

The same holds for other appropriately scaled continuous shift-invariant kernels  $K$  when  $\boldsymbol{\omega} \sim \lambda$  for the corresponding probability density  $\lambda$ .

**Random Fourier features.** Let  $\lambda$  be the probability density obtained as the Fourier transform of an appropriately scaled continuous shift-invariant kernel  $K$ . Suppose we randomly draw  $\boldsymbol{\omega} \sim \lambda$ , and then define the feature map  $\mathbf{z}_{\boldsymbol{\omega}}: \mathbb{R}^d \rightarrow \mathbb{R}^2$  by

$$\mathbf{z}_{\boldsymbol{\omega}}(\mathbf{x}) := (\cos(\langle \mathbf{x}, \boldsymbol{\omega} \rangle), \sin(\langle \mathbf{x}, \boldsymbol{\omega} \rangle)).$$

Then for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ ,

$$\begin{aligned}\mathbb{E}\langle \mathbf{z}_\omega(\mathbf{x}), \mathbf{z}_\omega(\mathbf{x}') \rangle &= \mathbb{E}(\cos(\langle \mathbf{x}, \omega \rangle) \cos(\langle \mathbf{x}', \omega \rangle) + \sin(\langle \mathbf{x}, \omega \rangle) \sin(\langle \mathbf{x}', \omega \rangle)) \\ &= K(\mathbf{x} - \mathbf{x}').\end{aligned}$$

Now suppose we draw  $m$  independent random vectors  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m \sim \lambda$ , and define  $\mathbf{z}: \mathbb{R}^d \rightarrow \mathbb{R}^{2m}$  by

$$\mathbf{z}(\mathbf{x}) := \sqrt{\frac{1}{m}}(\cos(\langle \mathbf{x}, \omega_1 \rangle), \sin(\langle \mathbf{x}, \omega_1 \rangle), \dots, \cos(\langle \mathbf{x}, \omega_m \rangle), \sin(\langle \mathbf{x}, \omega_m \rangle)).$$

Then, again, for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ ,

$$\begin{aligned}\mathbb{E}\langle \mathbf{z}(\mathbf{x}), \mathbf{z}(\mathbf{x}') \rangle &= \mathbb{E}\left(\frac{1}{m} \sum_{i=1}^m \cos(\langle \mathbf{x}, \omega_i \rangle) \cos(\langle \mathbf{x}', \omega_i \rangle) + \sin(\langle \mathbf{x}, \omega_i \rangle) \sin(\langle \mathbf{x}', \omega_i \rangle)\right) \\ &= K(\mathbf{x} - \mathbf{x}').\end{aligned}$$

The variance of  $\langle \mathbf{z}(\mathbf{x}), \mathbf{z}(\mathbf{x}') \rangle$  is roughly  $O(1/m)$  since the  $\{\omega_i\}_{i \in [n]}$  are independent. Therefore if  $m = \Omega(1/\epsilon^2)$ , then  $|\langle \mathbf{z}(\mathbf{x}), \mathbf{z}(\mathbf{x}') \rangle - K(\mathbf{x} - \mathbf{x}')| \leq \epsilon$  with high probability. (This can be made precise using via Hoeffding's inequality.)

To recap, we have shown that using a relatively large number  $m$  of these *random Fourier features* based on an appropriate probability distribution, we are able to approximate the kernel function using standard inner products in  $\mathbb{R}^m$  (well,  $\mathbb{R}^{2m}$ ). Moreover, this mapping  $\mathbf{z}$  is constructed in a completely data-oblivious way.

### 3 Dimensionality reduction

In some cases, it is feasible to explicitly compute the feature expansion  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$  because it has a compact implicit representation. For instance, if the base feature vectors  $\mathbf{x} \in \mathbb{R}^d$  are sparse (i.e., they have only a small number of non-zero entries), then low-degree interaction features will also be (relatively) sparse. However, it may still be undesirable to maintain a high-dimensional weight vector in  $\mathbb{R}^D$ . A possible recourse is dimensionality reduction: construct a mapping  $\mathbf{A}: \mathbb{R}^D \rightarrow \mathbb{R}^m$  of high-dimensional feature vectors  $\phi(\mathbf{x}) \in \mathbb{R}^D$  to lower-dimensional vectors in  $\mathbb{R}^m$  (for some  $m \ll D$ ), and then learn a linear classifier in  $\mathbb{R}^m$ . As before, we would like the mapping  $\mathbf{A}$  to approximately preserve inner products:

$$\langle \mathbf{A}(\phi(\mathbf{x})), \mathbf{A}(\phi(\mathbf{x}')) \rangle \approx \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle.$$

Henceforth, we ignore the distinction between  $\mathbf{x}$  and  $\phi(\mathbf{x})$ , and simply discuss techniques for constructing mappings  $\mathbf{A}$  from  $\mathbb{R}^d$  to  $\mathbb{R}^m$ . In fact, the mappings we construct will be *linear maps*, and hence will be representable as matrices  $\mathbf{A} \in \mathbb{R}^{m \times d}$ .

**Gaussian random matrix.** Again, we shall use randomization to construct the mapping  $\mathbf{A} \in \mathbb{R}^{m \times d}$  in a data-oblivious way. Fix a pair of vectors  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ . We would like the mapping  $\mathbf{A}$  to approximately preserve their inner product:  $\langle \mathbf{Ax}, \mathbf{Ax}' \rangle \approx \langle \mathbf{x}, \mathbf{x}' \rangle$ . It turns out that because of the structure of Euclidean space, it suffices to preserve the (squared Euclidean) lengths of  $\mathbf{x} + \mathbf{x}'$  and  $\mathbf{x} - \mathbf{x}'$ . This is because for any pair of vectors  $\mathbf{u}, \mathbf{v}$  in Euclidean space,

$$\langle \mathbf{u}, \mathbf{v} \rangle = \frac{1}{4} (\|\mathbf{u} + \mathbf{v}\|^2 - \|\mathbf{u} - \mathbf{v}\|^2).$$

The upshot is that if  $\|\mathbf{A}(\mathbf{x} + \mathbf{x}')\|^2 \approx \|\mathbf{x} + \mathbf{x}'\|^2$  and  $\|\mathbf{A}(\mathbf{x} - \mathbf{x}')\|^2 \approx \|\mathbf{x} - \mathbf{x}'\|^2$ , then  $\langle \mathbf{Ax}, \mathbf{Ax}' \rangle \approx \langle \mathbf{x}, \mathbf{x}' \rangle$ .

Our first approach is to let  $\mathbf{A} \in \mathbb{R}^{m \times d}$  be a matrix of i.i.d. Gaussian random variables, each with variance  $1/m$ .

**Theorem 1.** *Let  $\mathbf{A} := (1/\sqrt{m})\mathbf{Z}$ , where  $\mathbf{Z}$  is a  $k \times d$  matrix of i.i.d.  $N(0, 1)$  random variables. Fix any  $\mathbf{v} \in \mathbb{R}^d$  and  $\delta \in (0, 1)$ . With probability at least  $1 - \delta$ ,*

$$\|\mathbf{v}\|^2 \left( 1 - 2\sqrt{\frac{\ln(2/\delta)}{m}} \right) \leq \|\mathbf{Av}\|^2 \leq \|\mathbf{v}\|^2 \left( 1 + 2\sqrt{\frac{\ln(2/\delta)}{m}} + \frac{2\ln(2/\delta)}{m} \right).$$

*Proof.* Fix  $\mathbf{v} \in \mathbb{R}^d$  and  $\delta \in (0, 1)$ . Without loss of generality, we may assume  $\mathbf{v}$  is a unit vector. Let  $Y_i := \sum_{j=1}^d Z_{i,j} v_j$  for each  $i \in [m]$ , so

$$\|\mathbf{Av}\|^2 = \frac{1}{m} \|\mathbf{Zv}\|^2 = \frac{1}{m} \sum_{i=1}^m \left( \sum_{j=1}^d Z_{i,j} v_j \right)^2 = \frac{1}{m} \sum_{i=1}^m Y_i^2.$$

Since  $Y_i$  is a linear combination of standard Gaussian random variables  $Z_{i,1}, Z_{i,2}, \dots, Z_{i,d}$ , it is itself distributed as a Gaussian random variable, with mean zero and variance  $\sum_{j=1}^d v_j^2 = 1$ , i.e.,  $Y_i \sim N(0, 1)$ . The independence of the  $Z_{i,j}$  implies that the  $Y_i$  are independent. The sum-of-squares of  $m$  independent standard Gaussian random variables  $\sum_{i=1}^m Y_i^2$  is distributed according to the  $\chi^2$  distribution with  $m$  degrees of freedom. In expectation,

$\sum_{i=1}^m Y_i^2$  is  $m$ . Moreover,

$$\Pr\left(\sum_{i=1}^m Y_i^2 > m + 2\sqrt{mt} + 2t\right) \leq e^{-t}$$

and

$$\Pr\left(\sum_{i=1}^m Y_i^2 < m - 2\sqrt{mt}\right) \leq e^{-t}$$

for any  $t > 0$ . Letting  $t := \ln(2/\delta)$  and combining these probability tail bounds with a union bound completes the proof.  $\square$

Suppose  $\|\mathbf{x}\|$  and  $\|\mathbf{x}'\|$  are both  $O(1)$ . Theorem 1 implies that if  $m = \Omega(1/\epsilon^2)$ , then  $|\langle \mathbf{A}\mathbf{x}, \mathbf{A}\mathbf{x}' \rangle - \langle \mathbf{x}, \mathbf{x}' \rangle| \leq O(\epsilon)$  with high probability.

**Sparse random matrix.** It is possible to improve on Theorem 1 in terms of computation. Observe that if  $\mathbf{x}$  has  $k$  non-zero entries, then computing  $\mathbf{A}\mathbf{x}$  requires  $O(km)$  time. When  $\mathbf{A}$  is a Gaussian random matrix, there is not really any structure that can be exploited to speed-up this computation.

It turns out it is possible to construct  $\mathbf{A}$  more cleverly so that it has only  $O(\sqrt{m \log d})$  non-zero entries per column. This way, computing  $\mathbf{A}\mathbf{x}$  can be done in  $O(k\sqrt{m \log d})$  time.

The construction is achieved using binary code vectors  $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_d \in \{0, 1\}^m$ , each with exactly  $s$  non-zero entries, that collectively satisfy

$$\langle \mathbf{c}_i, \mathbf{c}_j \rangle \leq \frac{2s^2}{m}, \quad \text{for all } 1 \leq i < j \leq d.$$

It can be shown that such codes exist when  $s = \Omega(\sqrt{m \log d})$ . In fact, a construction of these codes can be achieved using a (good) hash function  $h: [s] \times [d] \rightarrow [m/s]$  as follows. The  $j$ -th code vector  $\mathbf{c}_j \in \{0, 1\}^m$  is viewed as the concatenation of  $s$  vectors  $\mathbf{c}_{1,j}, \mathbf{c}_{2,j}, \dots, \mathbf{c}_{s,j} \in \{0, 1\}^{m/s}$ ; the  $i$ -th vector  $\mathbf{c}_{i,j}$  has a single non-zero entry in the  $h(i, j)$ -th position.

The linear mapping  $\mathbf{A}$  is constructed as the following:

$$\mathbf{A} := \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \\ \vdots \\ \mathbf{A}_s \end{bmatrix} \in \{-1, 0, 1\}^{m \times d}$$

where, for each  $i \in [s]$ ,

$$\mathbf{A}_i := \frac{1}{\sqrt{s}} \left[ \sigma_{i,1} \mathbf{c}_{i,1} \mid \sigma_{i,2} \mathbf{c}_{i,2} \mid \dots \mid \sigma_{i,d} \mathbf{c}_{i,d} \right] \in \{-1, 0, 1\}^{(m/s) \times d},$$

and the  $\{\sigma_{i,j}\}_{i \in [s], j \in [d]}$  are i.i.d. random variables with

$$\Pr(\sigma_{i,j} = 1) = \Pr(\sigma_{i,j} = -1) = \frac{1}{2}.$$

(The  $\{\sigma_{i,j}\}$  are often called *Rademacher random variables* or *random signs*.) For all practical purposes, the random signs can also be thought of as being specified by a hash function  $\sigma: [d] \times [s] \rightarrow \{\pm 1\}$ , although the next theorem relies on random signs.

**Theorem 2.** *Let  $\mathbf{A} \in \{-1, 0, 1\}^{m \times d}$  be the random matrix constructed from the code vectors and random signs as described above, and let  $s$  be the number of non-zero entries in each code vector. Fix any  $\mathbf{v} \in \mathbb{R}^d$  and  $\delta \in (0, 1)$ . With probability at least  $1 - \delta$ ,*

$$|\|\mathbf{A}\mathbf{v}\|^2 - \|\mathbf{v}\|^2| \leq \|\mathbf{v}\|^2 \cdot O\left(\sqrt{\frac{\log(1/\delta)}{m}} + \frac{\log(1/\delta)}{s}\right).$$

Suppose  $\|\mathbf{x}\|$  and  $\|\mathbf{x}'\|$  are both  $O(1)$ . Theorem 2 implies that if  $m = \Omega(1/\epsilon^2)$  and  $s = \Omega(\sqrt{m \log d})$  (this latter condition is also used to guarantee the existence of the required code vectors), then  $|\langle \mathbf{A}\mathbf{x}, \mathbf{A}\mathbf{x}' \rangle - \langle \mathbf{x}, \mathbf{x}' \rangle| \leq O(\epsilon)$  with high probability.

Let  $h: [s] \times [d] \rightarrow [m/s]$  and  $\sigma: [s] \times [d] \rightarrow \{\pm 1\}$  be the hash functions representing  $\mathbf{A}$ . Computing  $\mathbf{A}\mathbf{x}$  is done as follows:

Let  $\mathbf{y} := \mathbf{0} \in \mathbb{R}^m$ .

**for**  $i = 1, 2, \dots, s$  **do**

**for** each non-zero entry  $x_j$  in  $\mathbf{x}$  (with  $j \in [d]$  and  $x_j \in \mathbb{R} \setminus \{0\}$ ) **do**

        Let  $y_{(i-1)s+h(i,j)} := y_{(i-1)s+h(i,j)} + \sigma(i,j)x_j$

**end for**

**end for**

If  $\mathbf{x}$  has  $k$  non-zero entries, then the computation required is  $O(ks) = O(k\sqrt{m \log d})$  (compared to  $O(km)$  for the Gaussian random matrix). Computing inner products  $\langle \mathbf{A}\mathbf{x}, \mathbf{w} \rangle$  with a weight vector  $\mathbf{w} \in \mathbb{R}^m$  is similar. It is also possible to show non-trivial guarantees for  $s = 1$  and also when the random signs  $\sigma$  are omitted.