

Title: Random feature

URL: https://en.wikipedia.org/wiki/Random_feature

PageID: 77992820

Categories: Category:Kernel methods for machine learning, Category:Machine learning, Category:Monte Carlo methods

Source: Wikipedia (CC BY-SA 4.0).

Random features (RF) are a technique used in machine learning to approximate kernel methods , introduced by Ali Rahimi and Ben Recht in their 2007 paper " Random Features for Large-Scale Kernel Machines " , [1] and extended by. [2] [3] RF uses a Monte Carlo approximation to kernel functions by randomly sampled feature maps . It is used for datasets that are too large for traditional kernel methods like support vector machine , kernel ridge regression , and gaussian process .

Mathematics

Kernel method

Given a feature map $\phi : \mathbb{R}^d \rightarrow V$, where V is a Hilbert space (more specifically, a reproducing kernel Hilbert space), the kernel trick replaces inner products in feature space $\langle \phi(x_i), \phi(x_j) \rangle_V$ by a kernel function $k(x_i, x_j) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$. Kernel methods replaces linear operations in high-dimensional space by operations on the kernel matrix: $K_X := [k(x_i, x_j)]_{i,j \in 1:N}$ where N is the number of data points.

Random kernel method

The problem with kernel methods is that the kernel matrix K_X has size $N \times N$. This becomes computationally infeasible when N reaches the order of a million. The random kernel method replaces the kernel function k by an inner product in low-dimensional feature space \mathbb{R}^D : $k(x, y) \approx \langle z(x), z(y) \rangle$ where z is a randomly sampled feature map $z : \mathbb{R}^d \rightarrow \mathbb{R}^D$.

This converts kernel linear regression into linear regression in feature space, kernel SVM into SVM in feature space, etc. Since we have $K_X \approx Z_X^T Z_X$ where $Z_X = [z(x_1), \dots, z(x_N)]$, these methods no longer involve matrices of size $O(N^2)$, but only random feature matrices of size $O(DN)$.

Random Fourier feature

Radial basis function kernel

The radial basis function (RBF) kernel on two samples $x_i, x_j \in \mathbb{R}^d$ is defined as [4]

where $\|x_i - x_j\|^2$ is the squared Euclidean distance and σ is a free parameter defining the shape of the kernel. It can be approximated by a random Fourier feature map $z : \mathbb{R}^d \rightarrow \mathbb{R}^D$: $z(x) := \frac{1}{\sqrt{D}} [\cos(\omega_1 \cdot x), \sin(\omega_1 \cdot x), \dots, \cos(\omega_D \cdot x), \sin(\omega_D \cdot x)]^T$ where $\omega_1, \dots, \omega_D$ are IID samples from the multidimensional normal distribution $N(0, \sigma^{-2}I)$.

Theorem — -

(Unbiased estimation) $E[\langle z(x), z(y) \rangle] = e^{-(x-y)^2 / (2\sigma^2)}$.

(Variance bound) $\text{Var}[\langle z(x), z(y) \rangle] = O(D^{-1})$

(Convergence) As $D \rightarrow \infty$, the approximation converges in probability to the true kernel.

(Unbiased estimation) By independence of $\omega_1, \dots, \omega_D$, it suffices to prove the case of $D = 1$. By the trigonometric identity $\cos(a-b) = \cos(a)\cos(b) + \sin(a)\sin(b)$, $\langle z(x), z(y) \rangle = \frac{1}{D} \sum_{i=1}^D \cos(\langle \omega_i, x-y \rangle)$. Apply the spherical symmetry of normal distribution, then evaluate the integral: $\int_{-\infty}^{\infty} \cos(kx) e^{-x^2/2} dx = e^{-k^2/2}$.

(Variance bound) Since $\omega_1, \dots, \omega_D$ are IID, it suffices to prove that the variance of $\cos(\langle \omega_1, x-y \rangle)$ is finite, which is true since it is bounded within $[-1, +1]$.

(Convergence) By Chebyshev's inequality.

Since \cos, \sin are bounded, there is a stronger convergence guarantee by Hoeffding's inequality. [1]: Claim 1

Random Fourier features

By Bochner's theorem, the above construction can be generalized to arbitrary positive definite shift-invariant kernel $k(x, y) = k(x-y)$.

Define its Fourier transform $p(\omega) = \frac{1}{2\pi} \int_{\mathbb{R}^d} e^{-j\langle \omega, \Delta \rangle} k(\Delta) d\Delta$ then $\omega_1, \dots, \omega_D$ are sampled IID from the probability distribution with probability density p . This applies for other kernels like the Laplace kernel and the Cauchy kernel.

Neural network interpretation

Given a random Fourier feature map z , training the feature on a dataset by featurized linear regression is equivalent to fitting complex parameters $\theta_1, \dots, \theta_D \in \mathbb{C}$ such that $f_\theta(x) = \text{Re}(\sum_k \theta_k e^{j\langle \omega_k, x \rangle})$ which is a neural network with a single hidden layer, with activation function $t \mapsto e^{it}$, zero bias, and the parameters in the first layer frozen.

In the overparameterized case, when $2D \geq N$, the network linearly interpolates the dataset $\{(x_i, y_i)\}_{i=1}^N$, and the network parameters is the least-norm solution: $\hat{\theta} = \arg \min_{\theta \in \mathbb{C}^D} \|\theta\|_2$ such that $f_\theta(x_k) = y_k \forall k \in 1:N$. At the limit of $D \rightarrow \infty$, the L2 norm $\|\hat{\theta}\|_2 \rightarrow \sqrt{f_K \cdot H}$ where f_K is the interpolating function obtained by the kernel regression with the original kernel, and $\|\cdot\|_H$ is the norm in the reproducing kernel Hilbert space for the kernel. [5]

Other examples

Random binning features

A random binning features map partitions the input space using randomly shifted grids at randomly chosen resolutions and assigns to an input point a binary bit string that corresponds to the bins in which it falls. The grids are constructed so that the probability that two points $x_i, x_j \in \mathbb{R}^d$ are assigned to the same bin is proportional to $K(x_i, x_j)$.

$K(x_{\{i\}}, x_{\{j\}})$. The inner product between a pair of transformed points is proportional to the number of times the two points are binned together, and is therefore an unbiased estimate of $K(x_{\{i\}}, x_{\{j\}})$. Since this mapping is not smooth and uses the proximity between input points, Random Binning Features works well for approximating kernels that depend only on the L_1 distance between datapoints.

Orthogonal random features

Orthogonal random features [6] uses a random orthogonal matrix instead of a random Fourier matrix.

Historical context

In NIPS 2006, deep learning had just become competitive with linear models like PCA and linear SVMs for large datasets, and people speculated about whether it could compete with kernel SVMs. However, there was no way to train kernel SVM on large datasets. The two authors developed the random feature method to train those.

It was then found that the $O(1/D)$ variance bound did not match practice: the variance bound predicts that approximation to within 0.01 requires $D \sim 10^4$, but in practice required only $\sim 10^2$. Attempting to discover what caused this led to the subsequent two papers. [2] [3] [7]

See also

Kernel method

Support vector machine

Fourier transform

Monte Carlo method

References

External links

Random Walks - Random Fourier features