

## Lecture 12 - Random Features Continued

In the last lecture we introduced the idea of random (Fourier) features via Bochner's theorem

$$K(\underline{x}, \underline{x}') = k(\underline{x} - \underline{x}') = \mathbb{E}_{\underline{\omega} \sim \Lambda} \exp(i \underline{\omega}^T \underline{x}) \exp(-i \underline{\omega}^T \underline{x}')$$

$$\varphi_j: \mathbb{R}^d \rightarrow \mathbb{C} \quad \approx \frac{1}{N} \sum_{j=1}^N \underbrace{\exp(i \underline{\omega}_j^T \underline{x})}_{\varphi_j(\underline{x})} \underbrace{\exp(-i \underline{\omega}_j^T \underline{x}')}_{\overline{\varphi_j(\underline{x}')}}$$

In fact since both  $k$  &  $\Lambda$  are real valued we infer that we can choose our feature maps to be real valued

$$\begin{aligned} \exp(i \underline{\omega}^T (\underline{x} - \underline{x}')) &= \cos(\underline{\omega}^T (\underline{x} - \underline{x}')) + i \sin(\underline{\omega}^T (\underline{x} - \underline{x}')) \\ &= \cos(\underline{\omega}^T \underline{x}) \cos(\underline{\omega}^T \underline{x}') + \sin(\underline{\omega}^T \underline{x}) \sin(\underline{\omega}^T \underline{x}') \end{aligned}$$

Thus, we can define

$$\varphi_j: \mathbb{R}^d \rightarrow \mathbb{R}^2 \quad \varphi_j(\underline{x}) = (\cos(\underline{\omega}_j^T \underline{x}), \sin(\underline{\omega}_j^T \underline{x}))$$

①

& take the approx

$$\begin{cases} K(\underline{x}, \underline{x}') \approx \frac{1}{N} \sum_{j=1}^N \varphi_j^T(\underline{x}) \varphi_j(\underline{x}') \\ \underline{\omega}_j \stackrel{\text{iid}}{\sim} \Lambda \end{cases}$$

These are not the only possible ways to construct random features. Another option (Preferred by Rahimi & Recht) is to take

$$\varphi_j: \mathbb{R}^d \rightarrow \mathbb{R} \quad \varphi_j(\underline{x}) = \cos(\sqrt{2} \underline{\omega}_j^T \underline{x} + b_j), \\ \underline{\omega}_j \stackrel{\text{iid}}{\sim} \Lambda \quad \underline{b} \sim \mathcal{U}[0, 2\pi]$$

Observe: Here we are exploiting the lack of uniqueness of feature maps' i.e. in general there are many maps  $\varphi: \mathcal{X} \rightarrow \mathcal{F}$  st.

$$K(\underline{x}, \underline{x}') = \langle \varphi(\underline{x}), \varphi(\underline{x}') \rangle_{\mathcal{F}}$$

The question remains, if we use random features then how do we estimate the RKHS norms? There is a simple calculation that reveals this (although it is not a complete proof).

Suppose  $N \geq 1$

$$K(\underline{x}, \underline{x}') = \frac{1}{\sqrt{N}} \sum_{j=1}^N \varphi_j(\underline{x}) \varphi_j(\underline{x}')$$

$$= \sum_{j=1}^N \left( \frac{1}{\sqrt{N}} \varphi_j(\underline{x}) \right) \left( \frac{1}{\sqrt{N}} \varphi_j(\underline{x}') \right)$$

Now consider  $f(\underline{x}) = \sum_{j=1}^N \alpha_j \left( \frac{1}{\sqrt{N}} \varphi_j(\underline{x}) \right)$

& try to satisfy the reproducing property

$$\langle f, K(\underline{x}, \cdot) \rangle = \left\langle \sum_j \alpha_j \left( \frac{1}{\sqrt{N}} \varphi_j \right), \sum_j \underbrace{\left( \frac{1}{\sqrt{N}} \varphi_j(\underline{x}) \right)}_{\text{same coeff.}} \left( \frac{1}{\sqrt{N}} \varphi_j \right) \right\rangle$$

$$= f(\underline{x}) = \sum_{j=1}^N \alpha_j \left( \frac{1}{\sqrt{N}} \varphi_j(\underline{x}) \right)$$

from this we see

$$f = \sum_{j=1}^N \alpha_j \varphi_j, \quad g = \sum_{j=1}^N \beta_j \varphi_j$$

$$\text{Then } \langle f, g \rangle_N = \underline{\alpha}^T \underline{\beta} \text{ \& so } \|f\|_N = \|\underline{\alpha}\|_2$$

Be careful, the  $\|\cdot\|_N$  norm is the RKHS norm corresponding to  $K_N$  & not  $K = \lim_{N \rightarrow \infty} K_N$ !

## An example in Classification

Consider  $\mathcal{X} \subseteq \mathbb{R}^d$  & a function

$$l: \mathcal{X} \rightarrow \{-1, +1\}. \quad \text{eg } \mathcal{X} \subseteq \mathbb{R}^{784} \text{ for MNIST!}$$

The function  $l$  assigns a label to any point  $\underline{x} \in \mathcal{X}$ . Now suppose we are given a training

data  $\mathcal{X} = \{\underline{x}_1, \dots, \underline{x}_n\} \subset \mathcal{X}$  along with

labels  $\underline{y} = \{y_1, \dots, y_n\} \subseteq \{-1, +1\}^n$  &  $y_j \approx l(\underline{x}_j)$

Goal: given a new point  $\underline{x}_{n+1} \in \mathcal{X}$   
Predict  $l(\underline{x}_{n+1})$ .

(4)

we will do this using the so called probit method. First, we need a model for the  $y_j$ .

We assume

$$(*) \quad y_j = \text{Sign}(f(x_j) + \varepsilon_j)$$

where  $\varepsilon_j \stackrel{\text{iid}}{\sim} \psi$ , for eg  $\psi(t) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp(-\frac{t^2}{2\sigma^2})$

&  $\text{sign}(t) = \begin{cases} +1 & \text{if } t \geq 0 \\ -1 & \text{if } t < 0 \end{cases}$ . The function  $f: \mathcal{X} \rightarrow \mathbb{R}$  is called a latent variable / function. It is a convenient tool in our model.

Based on  $(*)$  we can formulate a log-likelihood / loss function for  $f$ .

Suppose  $\psi$  is symmetric. Then

$$\begin{aligned} \Pr(y_j = +1 | f) &= \Pr(f(x_j) + \varepsilon_j \geq 0) \\ &= \Pr(\varepsilon_j \geq -f(x_j)) = \int_{-f(x_j)}^{\infty} \psi(t) dt \\ &= \int_{-\infty}^{f(x_j)} \psi(t) dt = \Psi(f(x_j)) = \Psi(y_j f(x_j)) \end{aligned}$$

where  $\Psi$  is the CDF of  $\psi$ .

⑤

Repeating a similar calculation as above shows

$$\begin{aligned}\Pr(y_j = -1 | f) &= \Pr(f(x_j) + \varepsilon_j < 0) \\ &= \Psi(-f(x_j)) = \Psi(y_j f(x_j))\end{aligned}$$

In summary  $\Pr(y_j | f) = \Psi(y_j f(x_j))$

Since the  $\varepsilon_j$  are iid we have

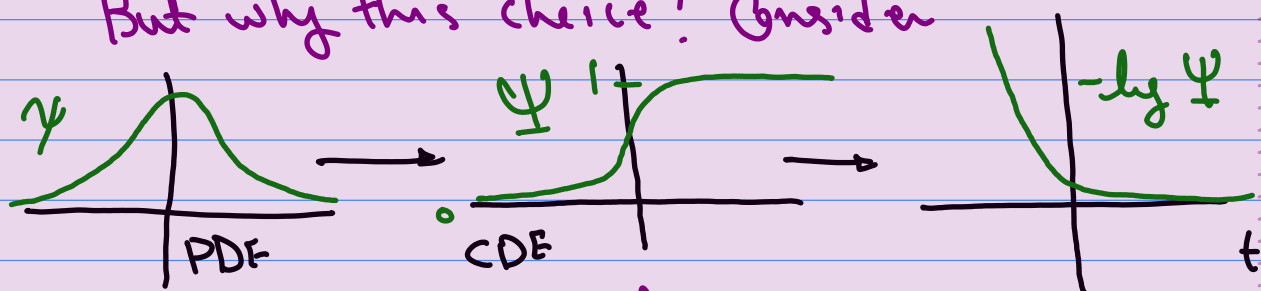
$$\Pr(\underline{y} | f) = \prod_{j=1}^n \Psi(y_j f(x_j))$$

This is simply the likelihood of  $\underline{y}$  given  $f$ . We therefore use the negative log of this function as our misfit term, i.e.

$$L(f) = -\log \Pr(\underline{y} | f) = -\sum_{j=1}^n \log \Psi(y_j f(x_j))$$

This loss is often referred to as the probit loss or model.

But why this choice? Consider



(6)

This loss is small if  $y_j f(x_j) \gg 0$ , i.e.  $y_j$  has the same sign as  $f(x_j)$ . Conversely it blows up if the signs

Disagree!

To this end we can formulate an opt. prob.  
to find  $f: \mathcal{X} \rightarrow \mathbb{R}$

$$f^* = \arg \min_{f \in \mathcal{H}_K} \sum_{j=1}^n -\log \Psi(y_j; f(x_j)) + \lambda \|f\|_K^2$$

so an application of the rep. th. yields

$$f^* = K(\cdot, X) \left( K(X, X) + \sigma^2 \bar{I} \right)^{-1} \underline{z}^*$$

↑ nugget if needed

$$\underline{z}^* = \arg \min_{\underline{z} \in \mathbb{R}^n} \sum_{j=1}^n -\log \Psi(y_j; z_j) + \lambda \underline{z}^T \left( K(X, X) + \sigma^2 \bar{I} \right)^{-1} \underline{z}$$

Alternatively, using random features, ie,

$$K(\underline{x}, \underline{x}') = \sum_{j=1}^N \left( \frac{1}{\sqrt{N}} \varphi_j(\underline{x}) \right) \left( \frac{1}{\sqrt{N}} \varphi_j(\underline{x}') \right)$$

$$f^* = \sum_{j=1}^N \omega_j^* \left( \frac{1}{\sqrt{N}} \varphi_j \right)$$

$$\underline{\omega}^* = \arg \min_{\underline{\omega} \in \mathbb{R}^N} \sum_{j=1}^n -\log \Psi(y_j; f(\underline{x}_j)) + \lambda \|\underline{\omega}\|_2^2$$

$$\text{s.t. } f(\underline{x}) = \sum_{j=1}^N \omega_j \left( \frac{1}{\sqrt{N}} \varphi_j \right)$$

If  $\Psi$  is taken to be log concave, then  $\Psi$  is also log-concave (e.g. Gaussian, logistic, Laplace, etc.) Then both problems become convex! & have unique solutions! They can be solved efficiently!



