

On the Error of Random Fourier Features

Dougal J. Sutherland and Jeff Schneider

{dsutherl,schneide}@cs.cmu.edu

Carnegie Mellon University



Random Fourier features

Random Fourier features (Rahimi and Recht, 2007) scale shift-invariant kernels to large numbers of inputs by using linear models on $z(x)$, where $z : \mathbb{R}^d \rightarrow \mathbb{R}^D$ has $k(x, y) \approx z(x)^\top z(y)$.

Let $\Delta := x - y$, and $k(x, y) = k(\Delta)$, $k(0) = 1$ be a continuous PSD kernel. Its Fourier transform $P(\omega)$ is a probability measure (Bochner's theorem).

One embedding is:

$$\phi(x) := \sqrt{\frac{2}{D}} \begin{bmatrix} \sin(\omega_1^\top x) \\ \cos(\omega_1^\top x) \\ \vdots \\ \sin(\omega_{D/2}^\top x) \\ \cos(\omega_{D/2}^\top x) \end{bmatrix}, \quad \omega_i \stackrel{iid}{\sim} P(\omega).$$

$\phi(x)^\top \phi(y)$ is an average of $D/2$ terms $\cos(\omega_i^\top \Delta)$; note $\mathbb{E} \cos(\omega^\top \Delta) = \mathfrak{R} \int e^{\omega^\top \Delta \sqrt{-1}} dP(\omega) = k(\Delta)$.

Another has more samples from $P(\omega)$, but with additional non-shift-invariant noise:

$$\psi(x) := \sqrt{\frac{2}{D}} \begin{bmatrix} \cos(\omega_1^\top x + b_1) \\ \vdots \\ \cos(\omega_D^\top x + b_D) \end{bmatrix}, \quad \begin{array}{l} \omega_i \stackrel{iid}{\sim} P(\omega) \\ b_i \stackrel{iid}{\sim} \text{Unif}_{[0, 2\pi]} \end{array}$$

$\psi(x)^\top \psi(y)$ is the mean of D terms of the form $\cos(\omega_i^\top \Delta) + \cos(\omega_i^\top (x + y) + 2b_i)$.

Our contribution

We show the ϕ embedding has lower variance than ψ for the Gaussian kernel, and improve the theoretical understanding of both embeddings' errors.

Prevalence of the embeddings

- The original publication discussed both ϕ and ψ . Online revisions only mention ψ (but give a bound only for ϕ). Later work used only ψ .
- Of the first 100 citations on Google Scholar, 15 used ψ , 14 used ϕ , 28 did not specify.
- All three library implementations we found (scikit-learn, Shogun, and JSAT) use ψ .

Variance

Using trig identities, we can show that

$$\begin{aligned} \text{Var } \phi(x)^\top \phi(y) &= \frac{1}{D} [1 + k(2\Delta) - 2k(\Delta)^2] \\ \text{Var } \psi(x)^\top \psi(y) &= \frac{1}{D} [1 + \frac{1}{2}k(2\Delta) - k(\Delta)^2]. \end{aligned}$$

So ϕ is lower-variance when

$$\text{Var } \cos(\omega^\top \Delta) = \frac{1}{2} + \frac{1}{2}k(2\Delta) - k(\Delta)^2 \leq \frac{1}{2}.$$

ϕ is better for Gaussian kernels

For $k(\Delta) := \exp(-\|\Delta\|^2/(2\sigma^2))$,

$$\text{Var } \cos(\omega) = \frac{1}{2} (1 - \exp(-\|\Delta\|^2/\sigma^2)) \leq \frac{1}{2}.$$

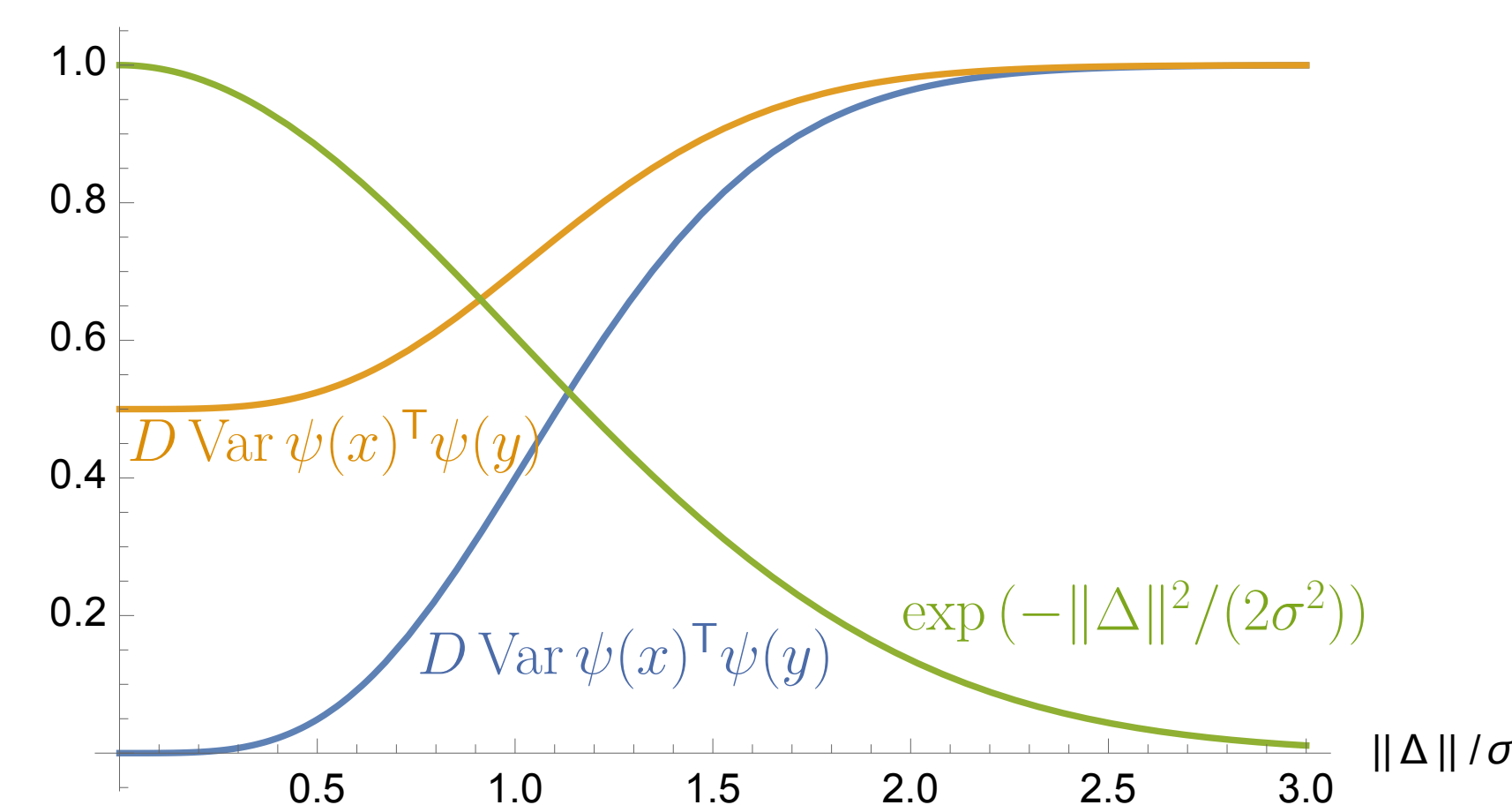


Figure 1: The variance per dimension for the Gaussian kernel. The difference in variance is higher for larger kernel values.

Improved uniform convergence

We can tighten the bound for ϕ and show one for ψ .

- Let ℓ be the diameter of the domain $\mathcal{X} \subset \mathbb{R}^d$.
- Let $\sigma_p^2 := \mathbb{E} \|\omega\|^2$, $\sigma_w^2 := \sup_{\Delta} [2 \text{Var } \cos(\omega^\top \Delta)]$.

Define $f_\phi(x, y) := \phi(x)^\top \phi(y) - k(x, y)$ to be the error for ϕ , and let $\alpha_\varepsilon := \min(1, \frac{1}{2}\sigma_w^2 + \frac{1}{3}\varepsilon)$; then

$$\Pr(\|f_\phi\|_\infty \geq \varepsilon) \leq \beta_d \left(\frac{\sigma_p \ell}{\varepsilon} \right)^{\frac{2}{1+\frac{1}{d}}} \exp\left(-\frac{D\varepsilon^2}{8(d+2)\alpha_\varepsilon}\right).$$

For ψ , define $f_\psi(x, y) := \psi(x)^\top \psi(y) - k(x, y)$ as well as $\alpha'_\varepsilon := \min(1, \frac{1}{8}(1 + \sigma_w^2) + \frac{1}{6}\varepsilon)$; then

$$\Pr(\|f_\psi\|_\infty \geq \varepsilon) \leq \beta'_d \left(\frac{\sigma_p \ell}{\varepsilon} \right)^{\frac{2}{1+\frac{1}{d}}} \exp\left(-\frac{D\varepsilon^2}{32(d+2)\alpha'_\varepsilon}\right).$$

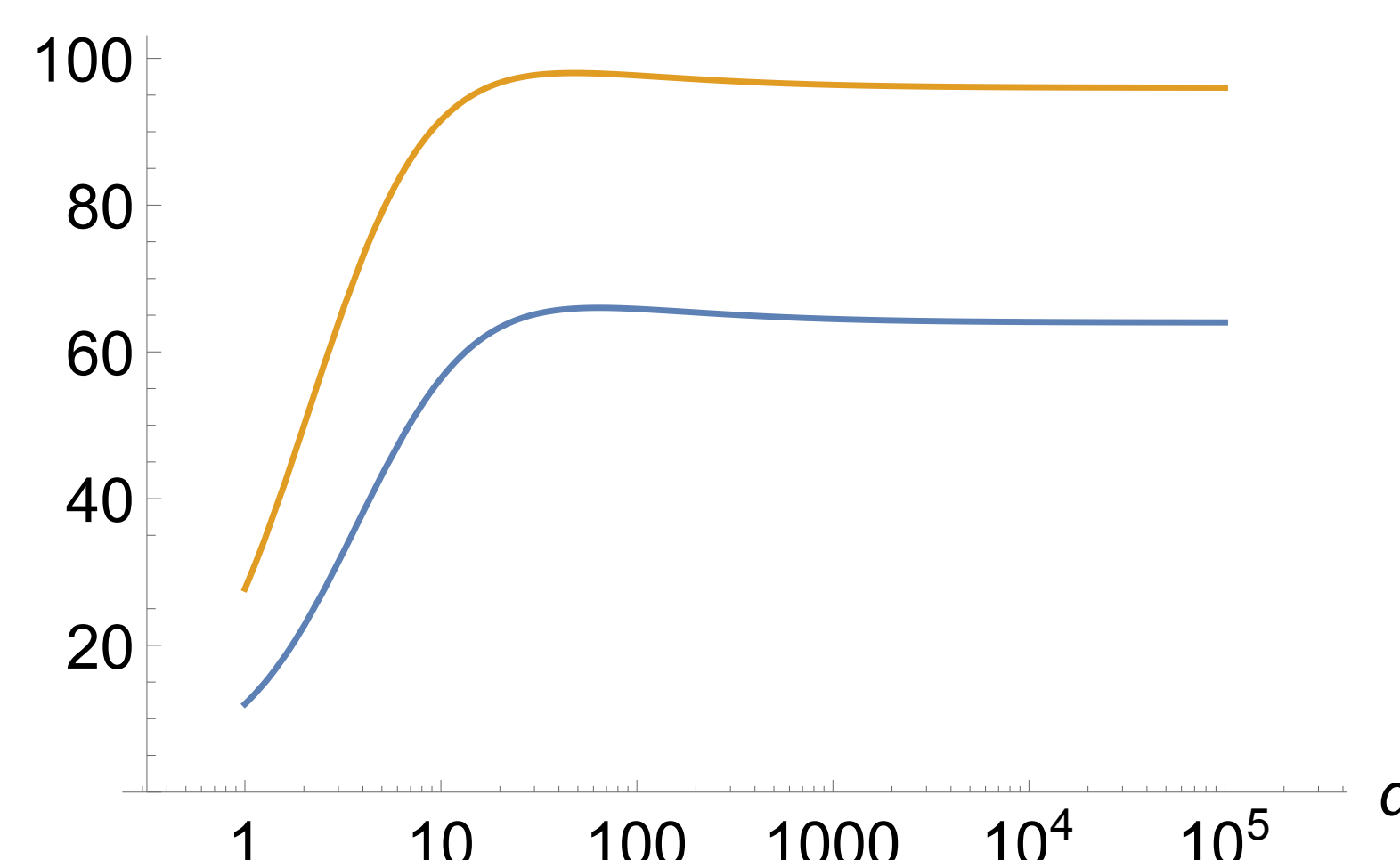


Figure 2: The coefficients β_d (blue, for ϕ) and β'_d (orange, for ψ).

The bound for ϕ is always tighter than that for ψ .

Expected max error

- Suppose $k(\Delta)$ is L -Lipschitz.
- Let $\gamma \approx 0.964$, $0.8 < \gamma' < 1.55$ depending on \mathcal{X} .

$$\mathbb{E} \|f_\phi\|_\infty \leq \frac{24\gamma\sqrt{d}\ell}{\sqrt{D}} \left(L + \mathbb{E} \max_{i=1, \dots, \frac{D}{2}} \|\omega_i\| \right),$$

$$\mathbb{E} \|f_\psi\|_\infty \leq \frac{48\gamma'\sqrt{d}\ell}{\sqrt{D}} \left(L + \mathbb{E} \max_{i=1, \dots, D} \|\omega_i\| \right)$$

using Dudley's entropy integral.

Concentration

$$\begin{aligned} \Pr(\|f_\phi\|_\infty \geq \mathbb{E} \|f_\phi\|_\infty + \varepsilon) \\ \leq 2 \exp\left(-\frac{D\varepsilon^2}{D\mathbb{E} \|f_\phi\|_\infty + \frac{1}{2}\sigma_w^2 + \frac{1}{6}D\varepsilon}\right), \end{aligned}$$

$$\begin{aligned} \Pr(\|f_\psi\|_\infty \geq \mathbb{E} \|f_\psi\|_\infty + \varepsilon) \\ \leq 2 \exp\left(-\frac{D\varepsilon^2}{\frac{4}{9}D\mathbb{E} \|f_\psi\|_\infty + \frac{1}{81}(1 + \sigma_w^2) + \frac{2}{27}D\varepsilon}\right) \end{aligned}$$

via Bousquet's inequality. f_ψ concentrates more tightly, but its mean is higher, both in the bound and empirically.

Numerical results with $d = 1$

Gaussian kernel, $\sigma = 1$. ϕ has solid lines, ψ dashed.

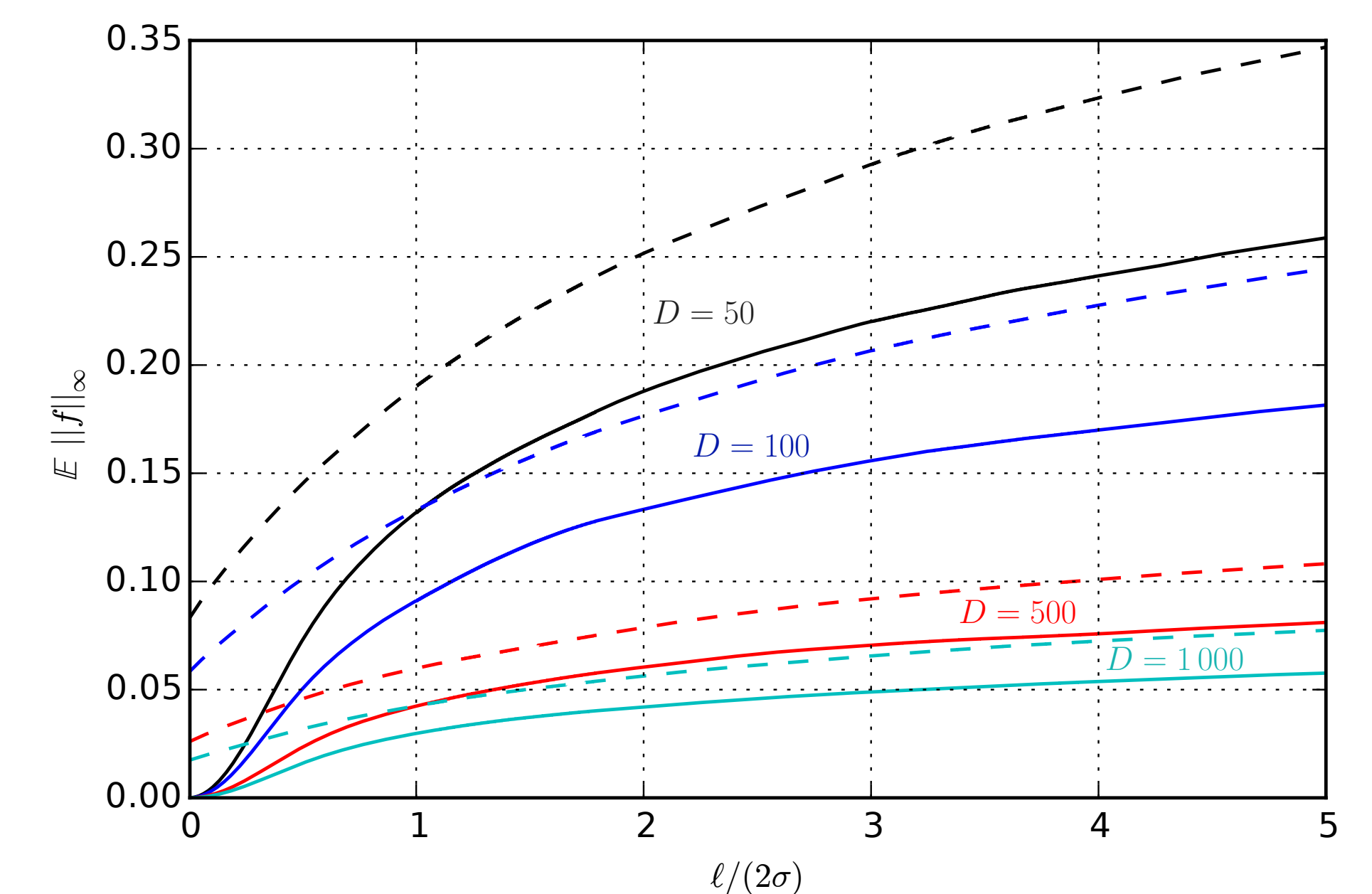


Figure 3: Average max error within a given radius.

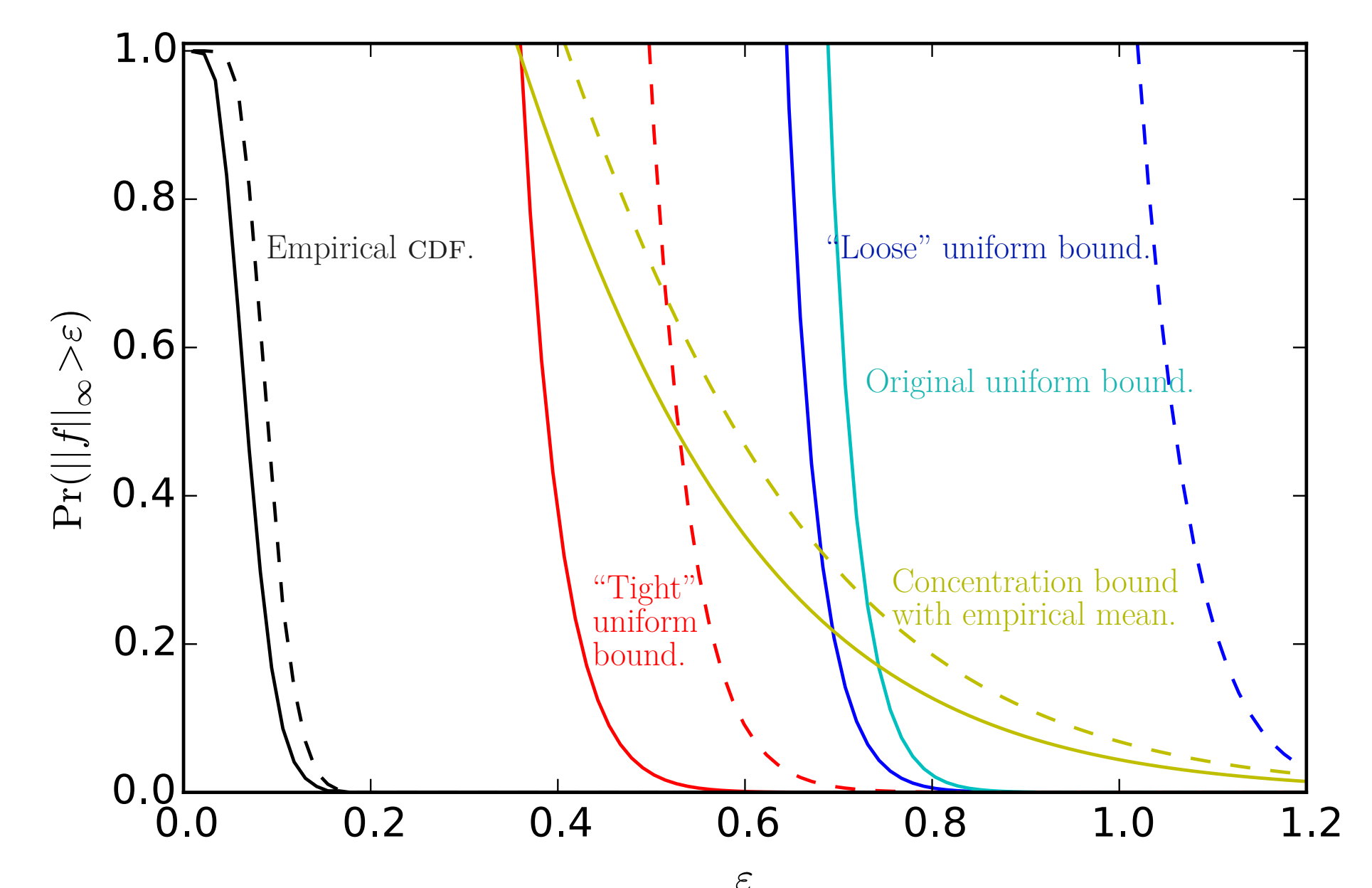


Figure 4: $\Pr(\|f\|_\infty > \varepsilon)$ on $[-3, 3]$ with $D = 500$.

Further results

The paper also has:

- Exact expectations and concentration bounds of squared L_2 error, for any measure.
- Bounds on changes in the outputs of ridge regression, SVM, and maximum mean discrepancy tests due to the features.
- More experiments.