# Random Feature Stein Discrepancies

## Jonathan Huggins
### Department of Biostatistics, Harvard University

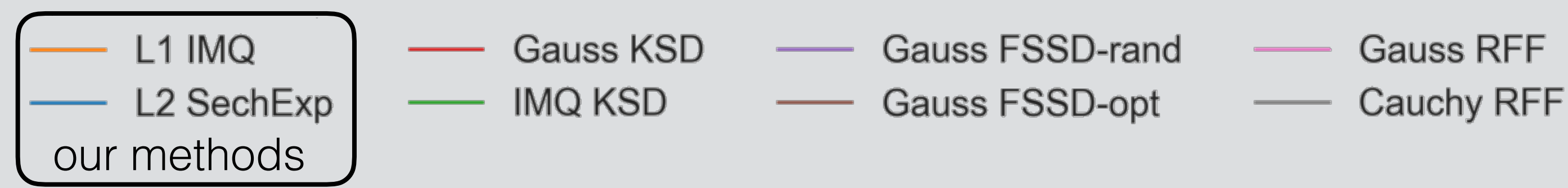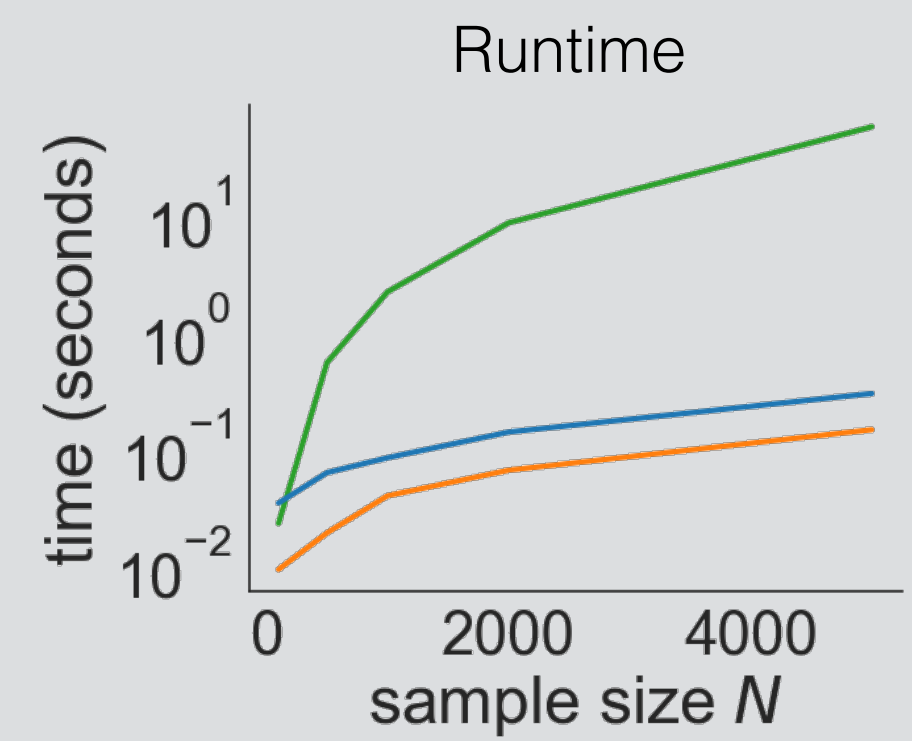## Lester Mackey
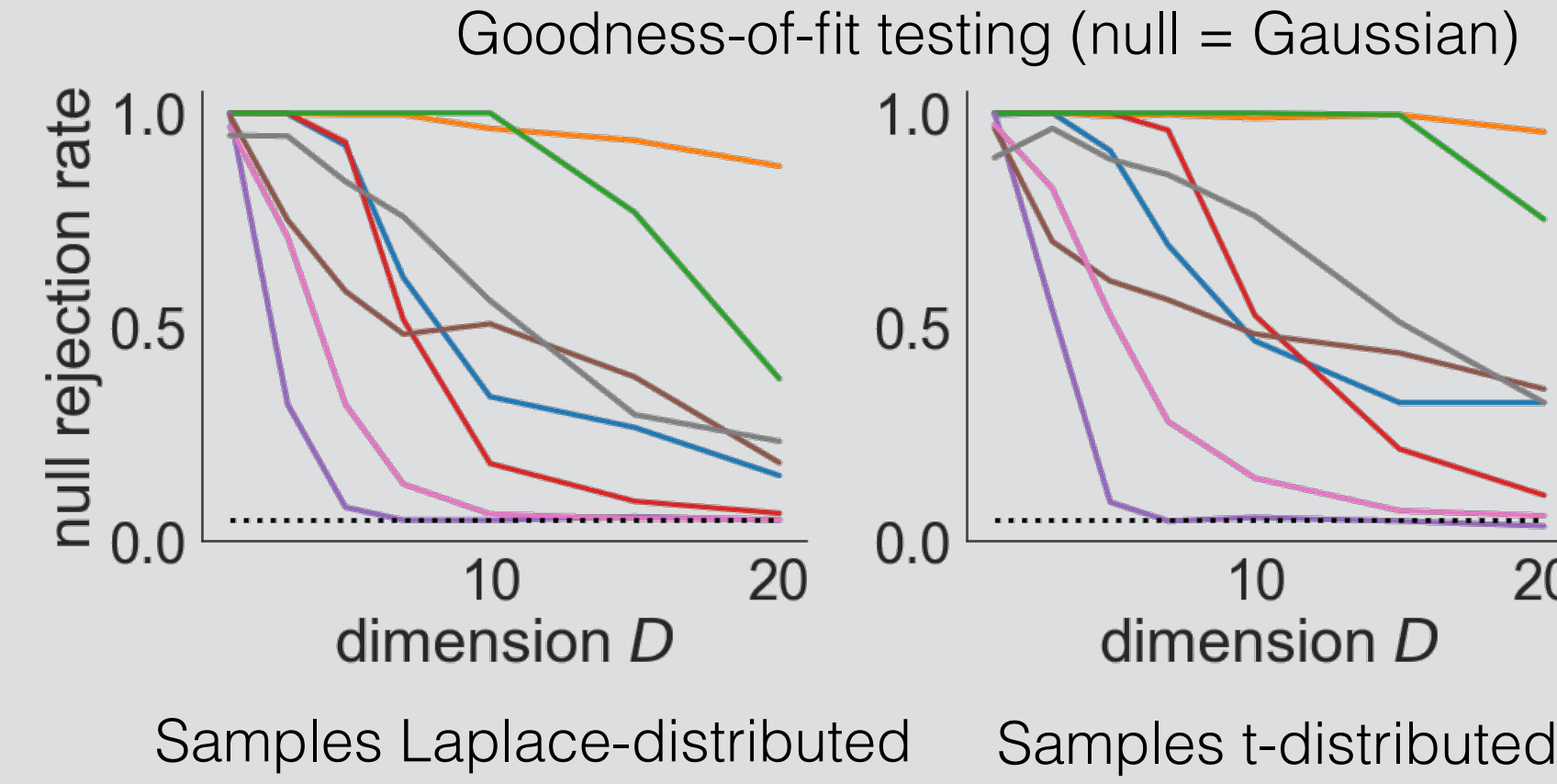### Microsoft Research New England

## Summary

- **Computable Stein discrepancies** used for…
  - sampler selection
  - posterior inference
  - goodness-of-fit testing
- **But** computation scales **quadratically** with sample size
- We introduce **random feature Stein discrepancies**, which…
  - retain excellent theoretical properties of existing Stein discrepancies
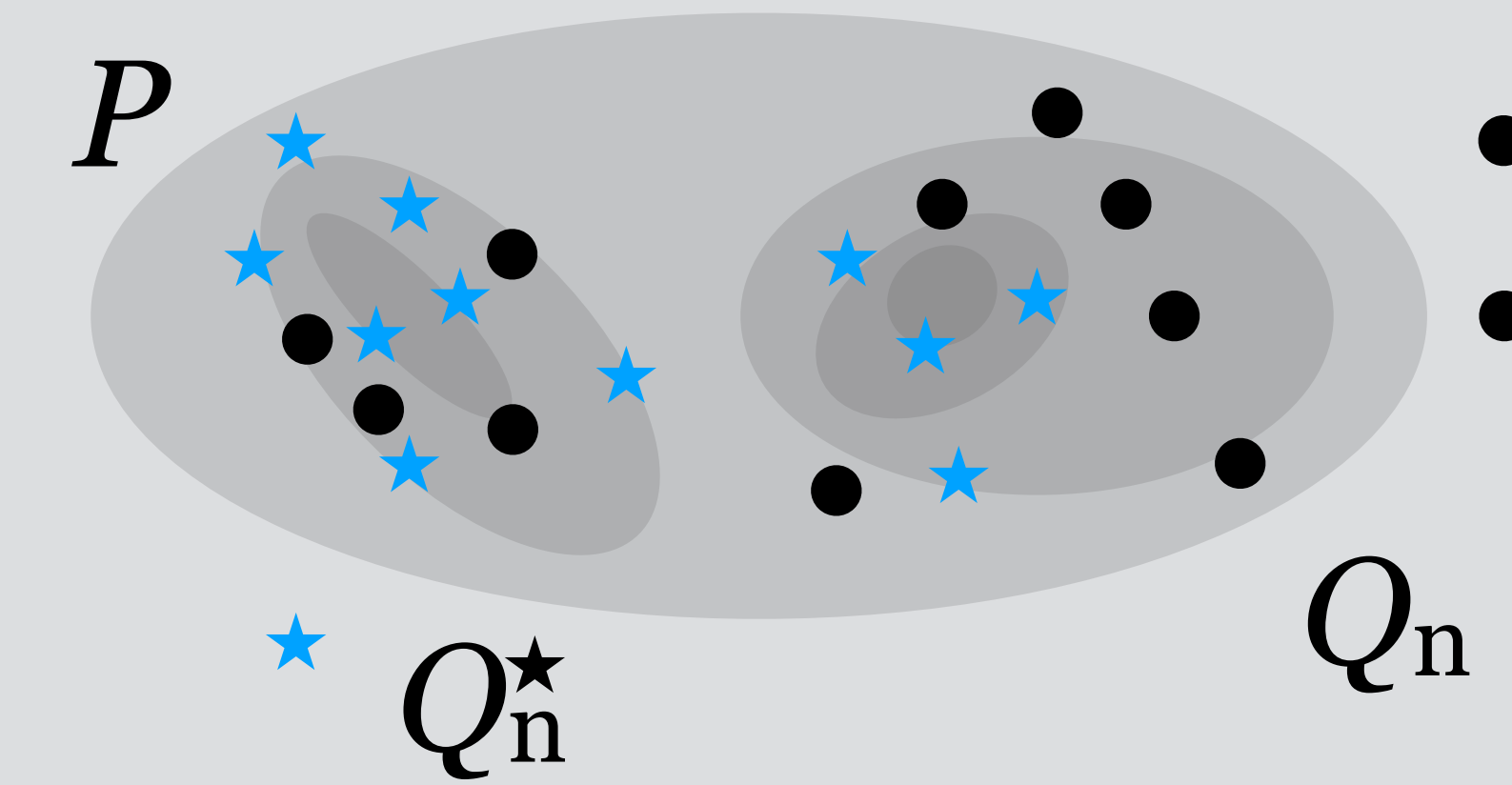  - are computable in **linear time**

| | our methods | | |
|---|---|---|---|
| L1 IMQ | Gauss KSD | Gauss FSSD-rand | Gauss RFF |
| L2 SechExp | IMQ KSD | Gauss FSSD-opt | Cauchy RFF |

faster…

more powerful…

Runtime



Goodness-of-fit testing (null = Gaussian)

Samples Laplace-distributed    Samples t-distributed

## The Big Picture

### I. Motivation



$P$

$Q_n^\star$    $Q_n$

sampler selection       goodness-of-fit
★ **versus** ●        $P = Q_n$ **versus** $P \neq Q_n$

### II. Stein discrepancies

$$d_{\mathcal{T}\mathcal{G}}(Q_n, P) = \sup_{g \in \mathcal{G}} |Q_n(\mathcal{T}g) - \cancel{P(\mathcal{T}g)}|$$
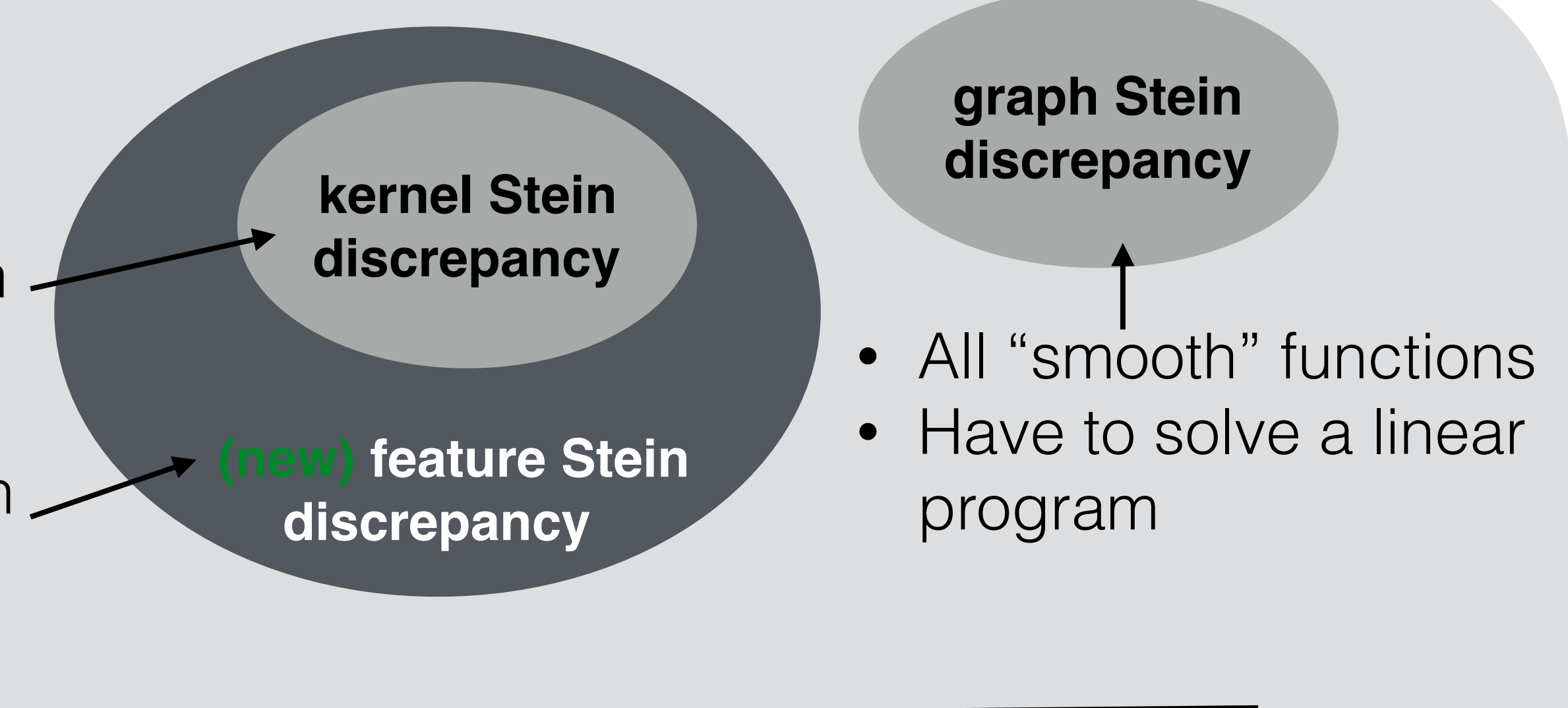
Stein operator: $P(\mathcal{T}g) = 0$

### III. Desiderata

1. Detect convergence of $Q_n$ to $P$
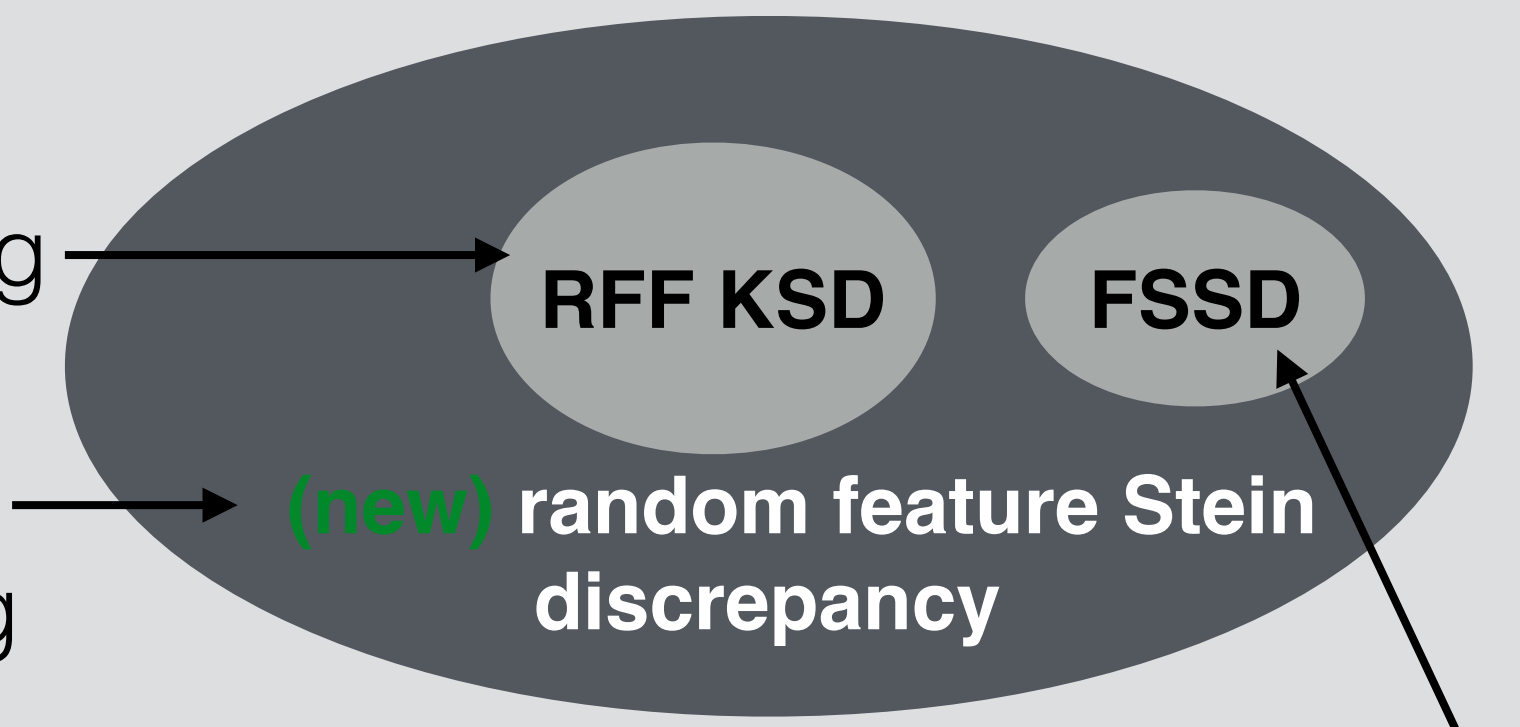2. Detect non-convergence
3. Computationally efficient

### IV. Choices for $\mathcal{G}$

- Uses a kernel function
- Closed-form

  **kernel Stein discrepancy**

  **[new] feature Stein discrepancy**

- Uses a feature function
- Approximate with importance sampling

**graph Stein discrepancy**

- All "smooth" functions
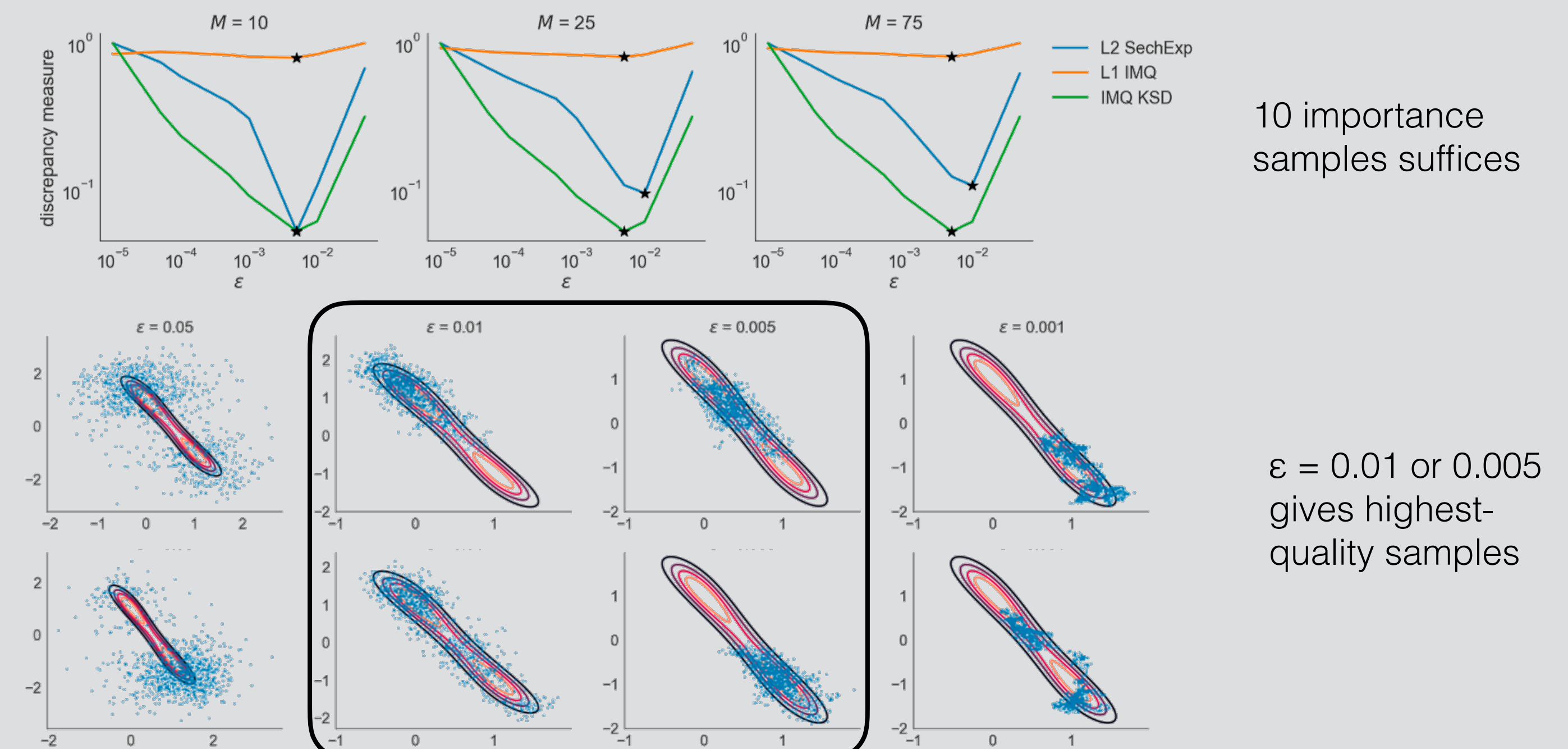- Have to solve a linear program

### V. Linear-time approximations

- Not convergence-determining

  **RFF KSD**    **FSSD**

- Can be convergence-determining
- Provably close to ΦSD when using near-linear number of importance samples [$\Omega(n^\kappa)$ for any $\kappa > 0$]

  **[new] random feature Stein discrepancy**

- Underpowered

### VI. Experiment: Selecting SGLD step size



$M = 10$    $M = 25$    $M = 75$

10 importance samples suffices

$\varepsilon = 0.05$    $\varepsilon = 0.01$    $\varepsilon = 0.005$    $\varepsilon = 0.001$

$\varepsilon = 0.01$ or $0.005$ gives highest-quality samples

## All the Details

### I. Defining the (random) feature Stein discrepancy

$\mathcal{G}_{\Phi,r} \triangleq \left\{ g : \mathbb{R}^D \to \mathbb{R} \mid g_d(x) = \int \Phi(x,z)\overline{f_d(z)}\,\mathrm{d}z \text{ with } \sum_{d=1}^D \|f_d\|_{L^s}^2 \leq 1 \text{ for } s = \frac{r}{r-1} \right\}.$

**Feature Stein discrepancy:** $\Phi\mathrm{SD}_{\Phi,r}^2(\mu, P) \triangleq \sup_{g \in \mathcal{G}_{\Phi,r}} |\mu(\mathcal{T}g)|^2 = \sum_{d=1}^D \|\mu(\mathcal{T}_d\Phi)\|_{L^r}^2$

**Random feature Stein discrepancy:** for $Z_1, \ldots, Z_M \overset{\text{i.i.d.}}{\sim} \nu$,

$\mathrm{R}\Phi\mathrm{SD}_{\Phi,r,\nu,M}^2(\mu, P) \triangleq \sum_{d=1}^D \left( M^{-1} \sum_{m=1}^M \nu(Z_m)^{-1} |\mu(\mathcal{T}_d\Phi)(Z_m)|^r \right)^{2/r}$

### II. Assumed form for feature function

**Assumption:** The base kernel has the form $k(x,y) = A_n(x)\Psi(x-y)A_n(y)$ for $A_n(x) \triangleq A(x - m_n)$ and $m_n \triangleq \mathbb{E}_{X \sim Q_n}[X]$

**Assumption:** $\Phi(x,z) = A_n(x)F(x-z)$

### III. Detecting convergence and non-convergence

**Proposition.** If the tilted Wasserstein distance

$\mathcal{W}_{A_n}(Q_n, P) \triangleq \sup_{h \in \mathcal{H}} |Q_n(A_n h) - P(A_n h)| \quad (\mathcal{H} \triangleq \{h : \|\nabla h(x)\|_2 \leq 1, \forall x \in \mathbb{R}^D\})$

converges to zero, then $\Phi\mathrm{SD}_{\Phi,r}(Q_n, P) \to 0$ and $\mathrm{R}\Phi\mathrm{SD}_{\Phi,r,\nu_n,M_n}(Q_n, P) \overset{P}{\to} 0$ for any choices of $r \in [1,2]$, $\nu_n$, and $M_n \geq 1$.

**Proposition** (KSD-ΦSD inequality)**.** If $k(x,y) = \int \mathcal{F}(\Phi(x,\cdot))(\omega)\overline{\mathcal{F}(\Phi(y,\cdot))(\omega)}\rho(\omega)\,\mathrm{d}\omega$, $r \in [1,2]$, and $\rho \in L^t$ for $t = r/(2-r)$, then

$\mathrm{KSD}_k^2(Q_n, P) \leq \|\rho\|_{L^t} \Phi\mathrm{SD}_{\Phi,r}^2(Q_n, P).$

### IV. Constructing convergence-determining RΦSDs

$(C,\gamma)$ **second moments:** We say $(\Phi, r, \nu)$ yields $(C,\gamma)$ second moments for $P$ and $Q_n$ if $\mathbb{E}[Y_{n,d}^2] \leq C\mathbb{E}[Y_{n,d}]^{2-\gamma}$ for $Y_{n,d} \triangleq |(Q_n\mathcal{T}_d\Phi)(Z)|^r/\nu(Z)$ and $Z \sim \nu$.

**Proposition.** Suppose $(\Phi, r, \nu)$ yields $(C,\gamma)$ second moments for $P$ and $Q_n$. If the reference $\mathrm{KSD}_k(Q_n, P) = \Omega(n^{-1/2})$ then a sample size $M = \Omega(n^{\gamma r/2})$ suffices to have, with high probability,

$$2\|\rho\|_{L^t}^{1/2} \mathrm{R}\Phi\mathrm{SD}_{\Phi,r,\nu,M}(Q_n, P) \geq \mathrm{KSD}_k(Q_n, P).$$

**Theorem.** Under regularity conditions, there exists a smoothness parameter $\overline{\lambda} \in (1/2, 1]$ and a constant $b \in [0,1)$ such that the following holds. For any $\xi \in (0, 1-b)$, $c > 0$, and $\alpha > 2(1-\overline{\lambda})$, if $\nu(z) \geq c\,\Psi(z - m_n)^{\xi r}$, then there exists a constant $C_\alpha > 0$ such that $(\Phi, r, \nu)$ yields $(C_\alpha, \gamma_\alpha)$ second moments for $P$ and $Q_n$, where $\gamma_\alpha \triangleq \alpha + (2-\alpha)\xi/(2 - b - \xi)$.

**Tilted hyperbolic secant kernel:**

$\Psi(x) = \Psi_a^{\mathrm{sech}}(x) \triangleq \prod_{d=1}^D \mathrm{sech}\left(\sqrt{\frac{\pi}{2}}ax_d\right) \quad \text{and} \quad A(x) = \prod_{d=1}^D e^{c\sqrt{1+x_d^2}}$

$L^2$ **tilted hyperbolic secant RΦSD:**

$F = \Psi_{2a}^{\mathrm{sech}}, \quad r = 2, \quad \text{and} \quad \nu(z) \propto \Psi_{4a\xi}^{\mathrm{sech}}(z - m_n)$

**Inverse multiquadric kernel:** $\Psi_{c,\beta}^{\mathrm{IMQ}}(x) \triangleq (c^2 + \|x\|_2^2)^\beta$ for some $\beta < 0$

$L^r$ **IMQ RΦSD:**

$F = \Psi_{c',\beta'}^{\mathrm{IMQ}}, \quad r = -D/(2\beta'\underline{\xi}), \quad \text{and} \quad \nu(z) \propto \Psi_{c',\beta'}^{\mathrm{IMQ}}(z - m_N)^{\xi r},$

where $c' = \overline{\lambda}c/2$, $\beta' \in [-D/(2\underline{\xi}), -\beta/(2\underline{\xi}) - D/(2\underline{\xi})]$, and $\xi \in (\underline{\xi}, 1)$

Simplest choice: $\beta' = -D/(2\underline{\xi})$ yields $r = 1$.

## More Experiments

### Efficiency



$\gamma = 0.500$
$\gamma = 0.333$
$\gamma = 0.250$
$\gamma = 0.200$

L1 IMQ    L2 SechExp    M necessary for stdev(RFSD) < FSD/2

### Goodness-of-fit testing



Size is controlled    Perturbed RBM

| | | |
|---|---|---|
| L2 SechExp | | Gauss KSD |
| L1 IMQ | | IMQ KSD |
| Gauss RFF | | |
| Cauchy RFF | | |
| Gauss FSSD-rand | | |
| Gauss FSSD-opt | | |