

CONTINUOUS-TIME LEARNING OF PROBABILITY DISTRIBUTIONS VIA NEURAL ODES WITH APPLICATIONS IN CONTINUOUS GLUCOSE MONITORING DATA

Antonio Álvarez-López

Departamento de Matemáticas,
Universidad Autónoma de Madrid

December 16, 2025

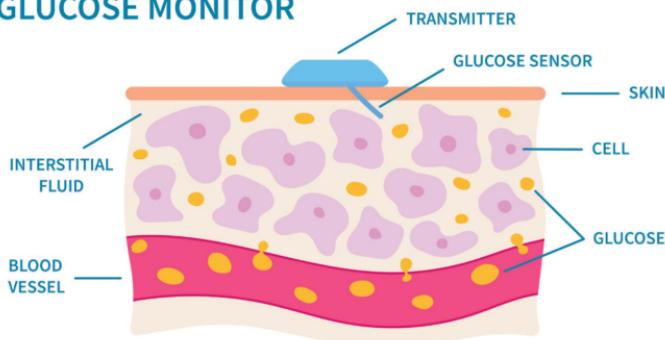
JOINT WORK WITH...



Marcos Matabuena (Harvard University)
Incoming professor in MBUZAI

MOTIVATION: DIGITAL HEALTH

CONTINUOUS GLUCOSE MONITOR



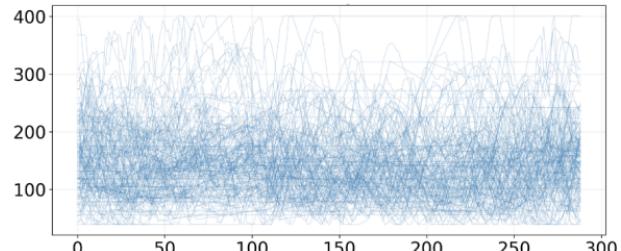
Source: (Messer et al. (2018)).

- ▶ Sends a measure every 5 minutes.
- ▶ **Advantages:**
 - Fewer finger-stick tests.
 - Alerts for hypo- and hyperglycaemia, or wide fluctuations in glucose levels
 - Treatment by closed-loop insulin pump.
- ▶ **Limitations:**
 - Periodic sensor replacement, higher cost.

DATA

Continuous Glucose Monitoring produces (irregular) data streams.

	0	5	10	15	20	25	30	35	40	45	50	55	60	65	I
1	NA	1													
2	216	222	226	246	250	254	256	250	248	246	242	238	234	230	2
3	114	114	112	112	110	108	108	108	106	104	104	102	100	100	3
4	178	180	182	184	186	188	188	190	190	190	192	192	192	186	4
5	178	174	170	168	168	168	168	168	166	166	164	162	164	166	5
6	196	196	190	186	190	194	196	196	192	190	186	186	188	190	6
7	70	74	76	76	76	74	74	74	74	76	78	78	78	78	7
8	186	186	186	188	188	190	190	192	194	196	198	200	202	204	8
9	222	220	216	214	212	210	210	210	206	202	202	204	208	210	9

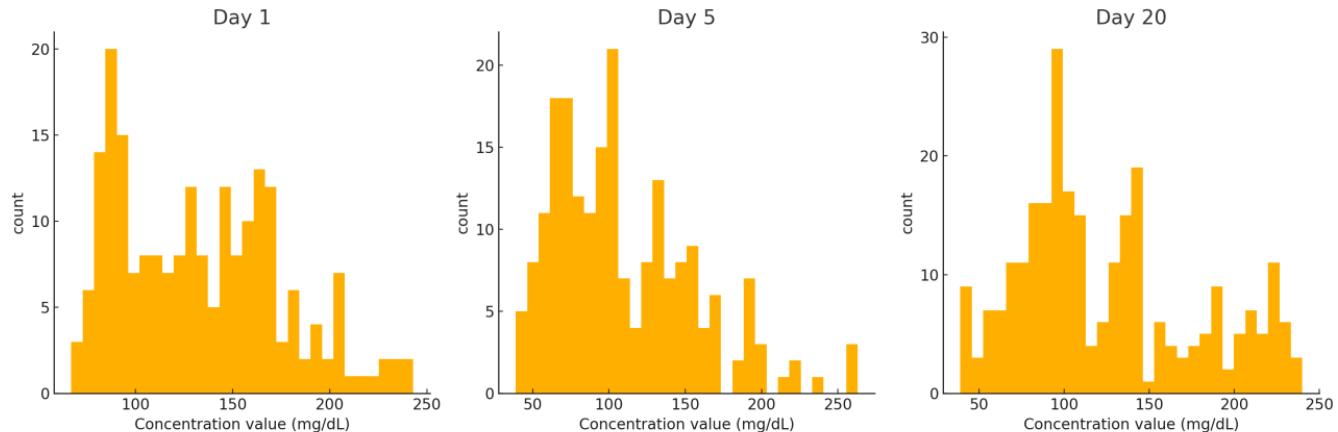


Glucose (mg/dL) vs. Measurement (total: 288/day)

Data source: PEDAP trial public dataset (Wadwa et al., NEJM 2023; Jaeb Center Public Study Websites).

DATA

We build **glucodensities** (Matabuena, Petersen, Vidal, Gude (2021)).



Idea: Treat CGM readings as samples and estimate density of % time spent at each glucose level.

GOALS

- ▶ **Mathematical:** Model an evolving probability law from samples observed at discrete times.
- ▶ **Clinical:** Track glucose distribution in time, reflecting disease progression or treatment efficacy.
- ▶ **Data difficulties:**
 - Empirical / discrete distributions from irregular streams with sensor noise.
 - Multimodality: Mixing times of day \Rightarrow different metabolic states.

GOALS

- ▶ **Mathematical:** Model an evolving probability law from samples observed at discrete times.
- ▶ **Clinical:** Track glucose distribution in time, reflecting disease progression or treatment efficacy.
- ▶ **Data difficulties:**
 - Empirical / discrete distributions from irregular streams with sensor noise.
 - Multimodality: Mixing times of day \Rightarrow different metabolic states.
- ▶ **Existing approaches:**
 - Classical density/CDF estimation (KDE, etc.), repeated over time (Chacón (2018))
 - Deep generative density models (e.g., normalizing flows) (Papamakarios et al. (2021))
 - Semiparametric time-varying models (GAMLSS) (Rigby & Stasinopoulos (2005))
 - Functional-quantile / multilevel distributional frameworks (Matabuena et al. (2025))
- ▶ **Main limitations:**
 - KDE: bandwidth sensitivity + curse of dimensionality.
 - Flow methods: often hard to interpret clinically.
 - GAMLSS: typically scalar responses and can impose rigid functional forms.
 - Functional-quantile dynamics: interpretable but often restricted to (near-)linear dynamics.

MODEL

Gaussian mixture with time-evolving weights

$$f_\theta(x, t) = \sum_{s=1}^K \alpha_s(t) \mathcal{N}(x | \mu_s, \Sigma_s), \quad x \in \mathbb{R}^d$$

where $\mu_s \in \mathbb{R}^d$, $\Sigma_s \in \mathbb{S}_+^d$ are fixed; and

$$(\alpha_1, \dots, \alpha_K) : [0, T] \longrightarrow \Delta_{K-1} := \left\{ \alpha \in \mathbb{R}^K \mid \alpha_i \geq 0, \sum_i \alpha_i = 1 \right\},$$
$$\dot{\alpha}(t) = \text{NODE}_\phi(\alpha(t), t).$$

¹Universal approximation of neural ODEs for dynamic behavior + Density of Gaussians (N. Wiener, Tauberian theorems, 1932).

MODEL

Gaussian mixture with time-evolving weights

$$f_\theta(x, t) = \sum_{s=1}^K \alpha_s(t) \mathcal{N}(x | \mu_s, \Sigma_s), \quad x \in \mathbb{R}^d$$

where $\mu_s \in \mathbb{R}^d$, $\Sigma_s \in \mathbb{S}_+^d$ are fixed; and

$$\begin{aligned} (\alpha_1, \dots, \alpha_K) : [0, T] &\longrightarrow \Delta_{K-1} := \left\{ \alpha \in \mathbb{R}^K \mid \alpha_i \geq 0, \sum_i \alpha_i = 1 \right\}, \\ \dot{\alpha}(t) &= \text{NODE}_\phi(\alpha(t), t). \end{aligned}$$

- ▶ **Interpretability:** Fixed Gaussian components $(m_s, \Sigma_s) \equiv$ Reference glycaemic states.
- ▶ By choosing K large enough, the parametric family of f_θ offers universal approximation¹ over

$$\left\{ f : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}_{\geq 0} \mid \int_{\mathbb{R}^d} f = 1 \right\}, \quad \|f\| := \sup_{t \in [0, T]} \|f(\cdot, t)\|_{L^1(\mathbb{R}^d)}.$$

¹Universal approximation of neural ODEs for dynamic behavior + Density of Gaussians (N. Wiener, Tauberian theorems, 1932).

TWO-STAGE ALGORITHM

- 1: **Input:** Time series $\{(t_i, X_i)\}_{i=1}^n$, number of Gaussians K
- 2: Construct aggregated data:

$$X = \bigcup_{i=1}^n X_i$$

- 3: **Initialization:** Run KMeans on X to obtain initial parameters $\{\alpha_s, \mu_s, \sigma_s^2\}_{s=1}^K$
- 4: **for** $\ell = 1$ to n_{iter} **do**
- 5: **Global GMM Fitting (Gradient Descent):** On n_{grad} iterations, find

$$\{\mu_s, \sigma_s^2\}_{s=1}^K = \arg \min_{\mu, \sigma^2} \text{MMD}^2 \left(P_X, \sum_{s=1}^K \alpha_s \mathcal{N}(\mu_s, \sigma_s^2) \right)$$

- 6: **Local Weight Estimation:** For each t_i , compute

$$(\alpha_1^{(i)}, \dots, \alpha_K^{(i)}) = \arg \min_{\alpha} \text{MMD}^2 \left(P_{X_i}, \sum_{s=1}^K \alpha_s \mathcal{N}(\mu_s^*, \sigma_s^{2*}) \right)$$

- 7: **end for**

- 8: **Neural ODE Modeling:** Define the weight dynamics

$$\frac{d\alpha(t)}{dt} = f(\alpha(t), \psi), \quad \alpha = (\alpha_1, \dots, \alpha_K).$$

- 9: **Parameter Estimation:** Find

$$\psi^* = \arg \min_{\psi} \sum_{i=1}^n \left\| \alpha^{(i)} - \alpha(t_i; \psi) \right\|^2,$$

where $\alpha(t_i; \psi)$ is the solution of the Neural ODE at time t_i .

Why two stages? Joint problem is strongly non-convex \Rightarrow
 Single pass converges to poor local minima.
 Alternating strategy yields stable updates

3–7. Discrete-time fit:

- Minimize a discrepancy (MMD²) for each t_i .
- Yields preliminary weights $\alpha^{(i)}$.

8–9. Continuous-time fit:

- Fit neural ODE interpolating the weights $\alpha^{(i)}$.
- Enforce temporal smoothness.

STAGE 1. DISCRETE-TIME FITTING: MAXIMUM MEAN DISCREPANCY

Kernel $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ symmetric, positive-definite. Common choice is the Gaussian kernel

$$k(x, x') = \exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right), \quad \sigma = \text{median}\{\|x_i - x_j\|\}_{1 \leq i < j \leq n} \quad (\text{heuristic pick})$$

Definition. Let μ, ν probability measures on \mathbb{R}^d . Define

$$\text{MMD}^2(\mu, \nu) = \iint k(x, x') d\mu(x) d\mu(x') + \iint k(y, y') d\nu(y) d\nu(y') - 2 \iint k(x, y) d\mu(x) d\nu(y).$$

MMD is a **distance** if and only if k is **characteristic**²; then $\text{MMD} = 0 \iff \mu = \nu$.

² $\mu \mapsto \mathbb{E}_\mu[k(x, \cdot)]$ injective. For instance, the Gaussian kernel.

STAGE 1. DISCRETE-TIME FITTING: MAXIMUM MEAN DISCREPANCY

Kernel $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ symmetric, positive-definite. Common choice is the Gaussian kernel

$$k(x, x') = \exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right), \quad \sigma = \text{median}\{\|x_i - x_j\|\}_{1 \leq i < j \leq n} \quad (\text{heuristic pick})$$

Definition. Let μ, ν probability measures on \mathbb{R}^d . Define

$$\text{MMD}^2(\mu, \nu) = \iint k(x, x') d\mu(x) d\mu(x') + \iint k(y, y') d\nu(y) d\nu(y') - 2 \iint k(x, y) d\mu(x) d\nu(y).$$

MMD is a **distance** if and only if k is **characteristic**²; then $\text{MMD} = 0 \iff \mu = \nu$.

Advantages:

- ▶ *Efficient*: For μ mixture and ν discrete, we use a closed-form expression for squared MMD.
- ▶ *Robust* to the presence of outliers. Also to non-overlapping supports of μ and ν
- ▶ *Differentiable*: gradients back-propagate through k .

² $\mu \mapsto \mathbb{E}_\mu[k(x, \cdot)]$ injective. For instance, the Gaussian kernel.

STAGE 2. WEIGHT EVOLUTION: NEURAL ODES

- **Continuous-time model.** Replace discrete network layers with an ODE:

$$\dot{\alpha}(t) = \text{NODE}_\phi(\alpha(t), t), \quad \alpha(0) = \alpha_0 \quad \rightarrow \quad \text{Output: } \alpha(t) = \text{ODESolve}(\alpha_0, t_0, t, \text{NODE}_\phi)$$

- **Projection to simplex:**

$$\alpha(t) \leftarrow \alpha(t)/\mathbf{1}^\top \alpha(t) \quad \forall t$$

- **Training loss.**

$$\mathcal{L}_{\text{NODE}}(\phi) = \sum_i \|\alpha(t_i; \phi) - \alpha^{(i)}\|^2 + \nu \|\phi\|^2, \quad \nu \geq 0 \text{ fixed}$$

$\nabla_\phi \mathcal{L}_{\text{NODE}}$ computed by adjoint method \Rightarrow constant-memory back-propagation.

- **Advantages**

1. Performance in high dimensions.
2. Parameter-efficient (one vector field = arbitrary depth).
3. Invertible flow (useful for generative models).
4. Smooth latent trajectories.

VALIDATION (SYNTHETIC DATA)

VALIDATION (SYNTHETIC DATA)

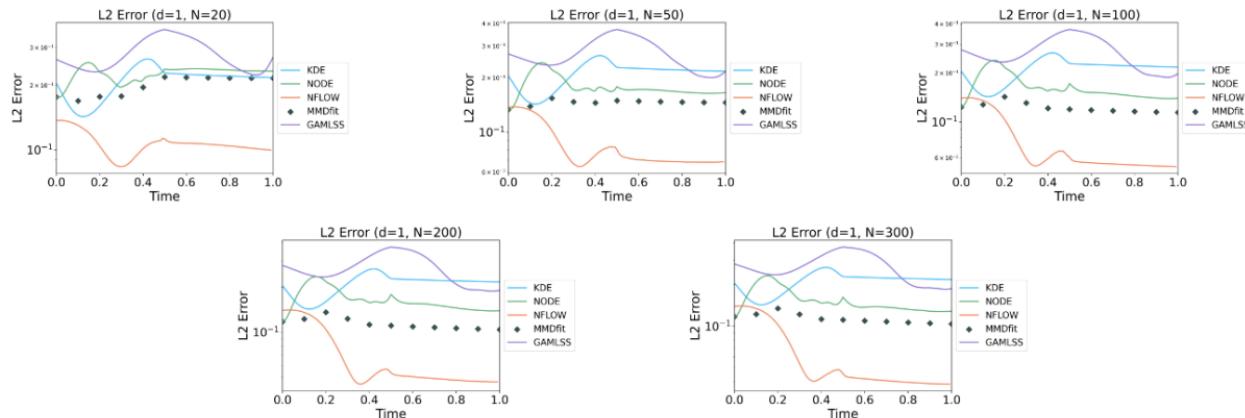
MMD Fitting Parameters

Parameter	Value	Parameter	Value
<i>Max iterations</i>	20	<i>Rel. tol. param.</i>	10^{-9}
<i>Rel. tol. err.</i>	10^{-7}	<i>Abs. tol.</i>	10^{-6}
<i>Components (K)</i>	10	<i>Learning rate</i>	0.01
<i>Grad. steps</i>	10	<i>Ridge reg. (λ)</i>	0.1

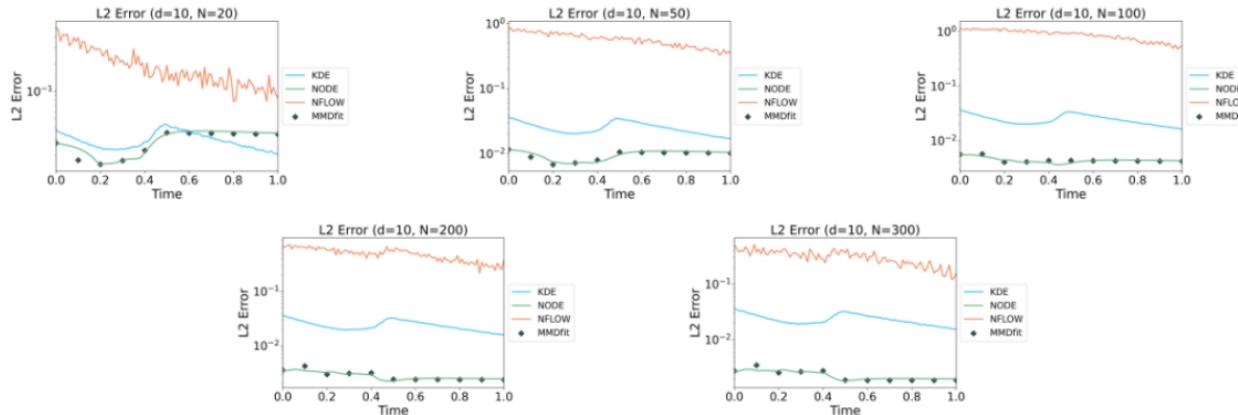
Neural ODE Training Parameters

Parameter	Value	Parameter	Value
<i>Seed</i>	42	<i>Hidden dim</i>	200
<i>Activation (σ)</i>	ReLU	<i>Layers</i>	2
<i>Optimizer</i>	Adam	<i>Integration time (T)</i>	1.0
<i>Step size (dt)</i>	0.01	<i>Integrator</i>	RK4
<i>Learning rate</i>	10^{-3}	<i>L2 reg. (λ)</i>	0.001
<i>Max epochs</i>	2000	<i>MC samples (n_{MC})</i>	10000
<i>Rel. tol.</i>	10^{-6}	<i>Abs. tol.</i>	10^{-6}

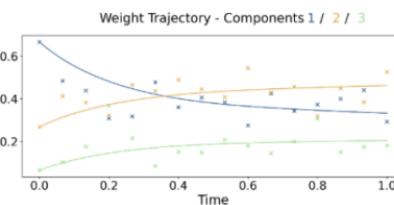
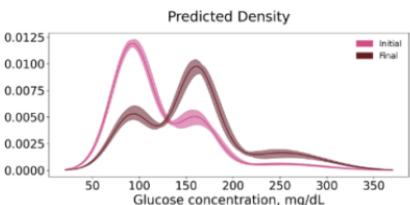
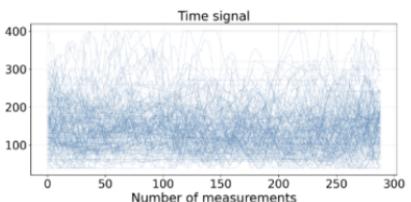
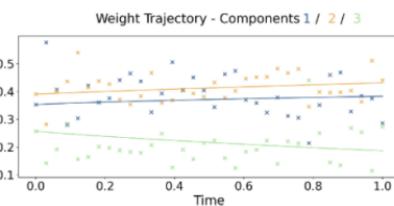
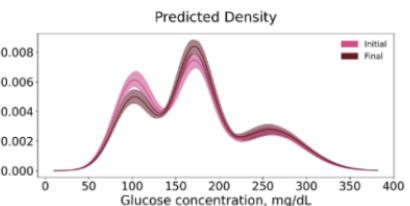
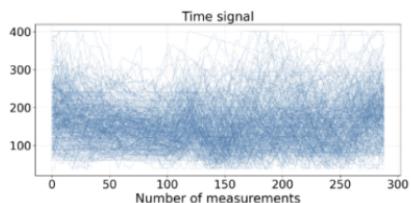
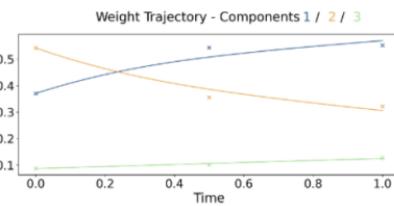
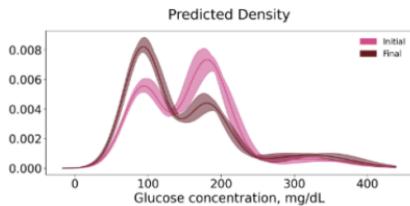
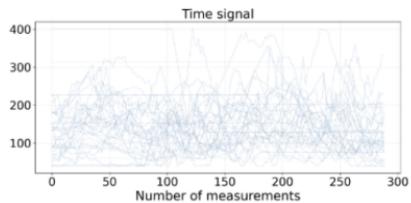
COMPARISON WITH OTHER MODELS: $d = 1$



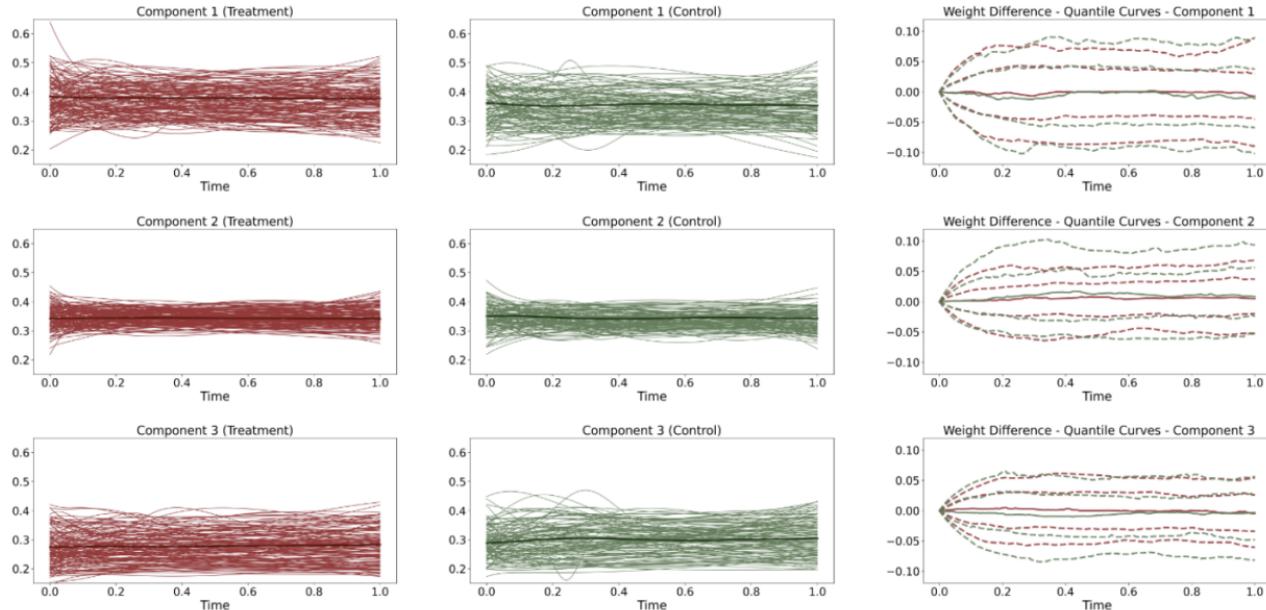
COMPARISON WITH OTHER MODELS: $d = 10$



CLINICAL CASE STUDY



CLINICAL CASE STUDY



Columns 1 and 2: Component weight trajectories in the treatment group (left) and control group (right).
Column 3: $(0.1, 0.25, 0.5, 0.75, 0.9)$ -quantiles of the process $Z_{is}(t) = \alpha_{is} - \alpha_{is}(0)$ ($s = 1, 2, 3$) for each group.

FUTURE WORK

- ▶ Use the model to predict clinical features/outcomes by regression.
- ▶ Generalize to purely **functional data** (each sample is an element of a Hilbert space) \implies probability over a functional space. This is motivated by biomechanics.
- ▶ Add **control** term to model action of insulin pumps and allow for design of controllers.

FUTURE WORK

- ▶ Use the model to predict clinical features/outcomes by regression.
- ▶ Generalize to purely **functional data** (each sample is an element of a Hilbert space) \implies probability over a functional space. This is motivated by biomechanics.
- ▶ Add **control** term to model action of insulin pumps and allow for design of controllers.

Thanks for the attention!



EXTRA I: RKHS, KDE AND MMD

1. Reproducing Kernel Hilbert Space (RKHS): A Hilbert space of functions where the evaluation functional is linear and bounded.

This defines a positive-definite kernel with the reproducing property:

$$f(x) = \langle f, k(\cdot, x) \rangle_{\mathcal{H}_k}, \quad \mathcal{H}_k := \overline{\text{span}}\{k(x, \cdot)\}_{x \in \mathcal{X}}.$$

Functions in \mathcal{H}_k can thus be evaluated via dot products \Rightarrow kernel trick.

2. Kernel Mean Embedding (KME): Image of a distribution P in the RKHS (it codifies all the means of functions in \mathcal{H}_k)

$$\mu_P = \mathbb{E}_{x \sim P}[k(x, \cdot)] \in \mathcal{H}_k, \quad \hat{\mu}_P = \frac{1}{m} \sum_{i=1}^m k(x_i, \cdot).$$

If the kernel is characteristic, the map $P \mapsto \mu_P$ is injective.

3. Maximum Mean Discrepancy (MMD): Distance between two KMEs

$$\text{MMD}_k(P, Q) = \|\mu_P - \mu_Q\|_{\mathcal{H}_k}, \quad \text{MMD}_k(P, Q) = 0 \iff P = Q \quad (\text{for characteristic } k).$$

Kernel \Rightarrow RKHS \Rightarrow KME \Rightarrow MMD

EXTRA II: COMPUTE MMD

Closed-form expression for the MMD between a Gaussian mixture and an empirical distribution:

$$f_i(x) = \sum_{s=1}^K w_s \mathcal{N}(x | m_s, \Sigma_s).$$

At each $t_i \in \tau$, we use a Gaussian kernel k_i such as (2), with $\sigma_i^2 \approx (\text{median}_{j \neq k} \|X_{t_i,j} - X_{t_i,k}\|)^2$. Thus,

$$\text{MMD}^2(P_{t_i}, Q) = \sum_{s=1}^K \sum_{r=1}^K w_s w_r I_{i,s,r} - \frac{2}{n_i} \sum_{s=1}^K \sum_{j=1}^{n_i} w_s J_{i,s,j} + \frac{1}{n_i^2} \sum_{j=1}^{n_i} \sum_{\ell=1}^{n_i} k_i(X_{t_i,j}, X_{t_i,\ell}).$$

The two first terms admit closed-form expressions:

$$I_{i,s,r} = \frac{(\sigma_i^2)^{d/2}}{\sqrt{\det(\Sigma_s + \Sigma_r + \sigma_i^2 \mathbf{Id})}} \exp\left(-\frac{1}{2} (m_s - m_r)^\top (\Sigma_s + \Sigma_r + \sigma_i^2 \mathbf{Id})^{-1} (m_s - m_r)\right),$$

$$J_{i,s,j} = \frac{(\sigma_i^2)^{d/2}}{\sqrt{\det(\Sigma_s + \sigma_i^2 \mathbf{Id})}} \exp\left(-\frac{1}{2} (X_{t_i,j} - m_s)^\top (\Sigma_s + \sigma_i^2 \mathbf{Id})^{-1} (X_{t_i,j} - m_s)\right),$$