# Compositional Model Design:
# High-Dimensional Multi-Output Regression

Wessel Bruinsma

University of Cambridge, CBL

13 May 2020

Model $p(x)$ for $x: \mathbb{R} \to \mathbb{R}^m$:

✗ computational complexity
✗ model complexity

✓ Existing learning and inference routines.

✗ Feasible for moderate number of outputs ($\sim$10–100).

Problem: Data may be very high-dimensional ($\sim$1000–10000).

• Example: daily temperature measurements around Europe.

• Often fewer ($\ll 1000$) underlying "mechanisms".

Goal: Design a procedure that "wraps" $p(x)$ to scale to many outputs.

Approach

Compose $p(x)$ with a likelihood $p(y \mid x)$ to scale to many outputs:

high-dim. model ⟵ $p(y) = \int p(y \mid x)p(x)\,\mathrm{d}x.$

maps into high-dim. space ⟵     ⟶ available model

Principled probabilistic approach: hope for well-calibrated uncertainty!

Desiderata:
- ✓ Learning and inference for $p(y)$ use existing routines for $p(x)$.
- ✓ Favourable scaling in number of outputs.
- ✓ Ability to deal with missing data.

```python
from dream import DeepGP, MultiOutput

dgp   = DeepGP(num_outputs=10)              # p(x)
model = MultiOutput(dgp, num_outputs=10000) # p(y)

model.fit(x, y)  # Uses `dgp.fit`.

pred = model.predict(x)  # Uses `dgp.predict`.
```

# Likelihood Model: $p(y \mid x)$

Data: $y(t) \in \mathbb{R}^p$.         Model: $x(t) \in \mathbb{R}^m$.

- Data is high-dimensional: $p \gg m$.
- Likelihood model $p(y \mid x)$ needs to transform $x(t)$ to $y(t)$.

Linear model:

$$
\overset{\mathbb{R}^p}{y(t)} = \overset{\mathbb{R}^p}{h_1} \overset{\mathbb{R}}{x_1(t)} + \cdots + \overset{\mathbb{R}^p}{h_m} \overset{\mathbb{R}}{x_m(t)} + \overset{\mathbb{R}^p}{\varepsilon(t)}
$$

$$
= \underset{\mathbb{R}^{p \times m}}{H} \underset{\mathbb{R}^m}{x(t)} + \varepsilon(t).
$$

- Data lives around $m$-dimensional linear subspace $\mathrm{col}(H) \subseteq \mathbb{R}^p$.

Full generative model:

$$x \sim p(x), \qquad\qquad\qquad \text{(latent model)}$$
$$f(t) \,|\, H, x(t) = Hx(t), \qquad\qquad \text{(mixing mechanism)}$$
$$y(t) \,|\, f(t) \sim \mathcal{N}(f(t), \Sigma), \qquad\qquad \text{(observation model)}$$

$$x \text{: "latent processes"},$$
$$H \text{: "basis" or "mixing matrix"}.$$

✓ Scales $p(x)$ to many outputs.

? Learning and inference for $p(y)$ use existing routines for $p(x)$.

? Favourable scaling in number of outputs.

? Deal with missing data.

We consider $x \sim \mathrm{GPAR}$:

$$
\begin{aligned}
x_1(t) &= f_1(t), & f_1 &\sim \mathcal{GP}(0, k_1), \\
x_2(t) &= f_2(t, x_1(t)), & f_2 &\sim \mathcal{GP}(0, k_2), \\
x_3(t) &= f_3(t, x_1(t), x_2(t)), & f_3 &\sim \mathcal{GP}(0, k_3).
\end{aligned}
$$

✓ Captures nonlinear dependencies between outputs.

✓ Learning and inference exact and closed form.

✓ Feasible for moderate number of outputs.

✓ Depends on ordering of outputs. (Implicit in $H$!)

✗ Cannot be further composed with Gaussian likelihood.

Inference in $p(y)$ in Terms of Inference in $p(x)$

Projection of the data:

$$T = \overset{\mathbb{R}^p}{\overbrace{y}} \mapsto \underset{\text{``observation for } p(x)\text{''}}{\underbrace{\overset{\mathbb{R}^m}{\overbrace{(H^{\mathsf{T}}\Sigma^{-1}H)^{-1}H^{\mathsf{T}}\Sigma^{-1}y}}}}.$$

Then use inference for $p(x)$!

**Proposition:**

1. $p(f \mid \overset{\mathbb{R}^{p \times n}}{\overbrace{Y}}) = p(f \mid \overset{\mathbb{R}^{m \times n}}{\overbrace{TY}})$ with $Ty_i \mid x_i \sim \mathcal{N}(x_i, \underset{\text{``projected noise''}}{\Sigma_T})$.

2. $p(Y) = \left[\prod_{i=1}^{n} \dfrac{\mathcal{N}(y_i \mid 0, \Sigma)}{\mathcal{N}(Ty_i \mid 0, \Sigma_T)}\right] \displaystyle\int p(x) \prod_{i=1}^{n} \mathcal{N}(Ty_i \mid x_i, \Sigma_T)\, \mathrm{d}x.$
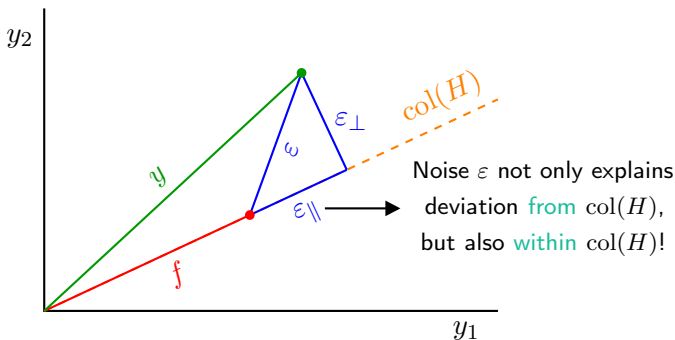
$$\log p(Y) \simeq \overbrace{\log \int p(x) \prod_{i=1}^{n} \mathcal{N}(Ty_i \,|\, x_i, \Sigma_T) \, \mathrm{d}x}^{\text{likelihood under } p(x) \text{ with additional noise}}$$
$$- \underbrace{\frac{1}{2} \sum_{i=1}^{n} \|y_i - HTy_i\|_{\Sigma}^2}_{\substack{\text{data "lost" by projection} \\ \text{(reconstruction error)}}} - \underbrace{\frac{1}{2} n \log \frac{|\Sigma|}{|\Sigma_T|}}_{\substack{\text{noise "lost" by} \\ \text{projection}}}.$$

- Learn $H \implies$ learn $T \implies$ learn a transform of the data!
  - "Regularisation terms" prevent underfitting.

✗ Requires additional projected noise $\Sigma_T$.

- Can we eliminate $\Sigma_T$?

- Consider case $y(t) \in \mathbb{R}^2$ and $x(t) \in \mathbb{R}$:



Noise $\varepsilon$ not only explains deviation from $\mathrm{col}(H)$, but also within $\mathrm{col}(H)$!

- $\varepsilon_{\parallel}$ is responsible for $\Sigma_T$!
- Idea: Set $\varepsilon_{\parallel} = 0$ to obtain $\Sigma_T = 0$.

- Consider $\Sigma = \sigma^2 I$. Then $T = H^\dagger$.

1. Decompose

$$\sigma^2 I = \sigma_{\parallel}^2 \underbrace{U U^{\mathsf{T}}}_{\text{orth proj. onto } \mathrm{col}(H)} + \sigma_{\perp}^2 \overbrace{N N^{\mathsf{T}}}^{\text{orth proj. onto } \mathrm{col}(H)^{\perp}}.$$

2. Take $\sigma_{\parallel}^2 \to 0$. Then $\Sigma_T \to 0$.

3. Optimise over $\sigma_{\perp}^2$:

measure of goodness of fit for $H$

$$\sigma_{\perp}^2 = \frac{1}{n(p-m)} \|(I - H H^\dagger) Y\|_F^2.$$

$$\log p_y(Y) \simeq \log p_x(H^\dagger Y)$$
$$- \frac{1}{2} n \log |H^\mathsf{T} H| - \frac{1}{2} n(p - m) \log(\|(I - HH^\dagger)Y\|_F^2)$$

✓ Learning and inference for $p(y)$ use existing routines for $p(x)$.

✓ Linear scaling in number of outputs.

• Possible to generalise to general $\Sigma$.

Missing data is tricky:

❶ Cannot compute $H^\dagger Y$.

❷ No mechanism to "tell $p(x)$".

**Proposition:** Missing data is equivalent to adding noise to $p(x)$:

$$\Sigma_{\mathsf{miss}} = \sigma_\perp^2 \Big[ ((L^\mathsf{T} H)^\mathsf{T} (L^\mathsf{T} H))^{-1} - (H^\mathsf{T} H)^{-1} \Big].$$
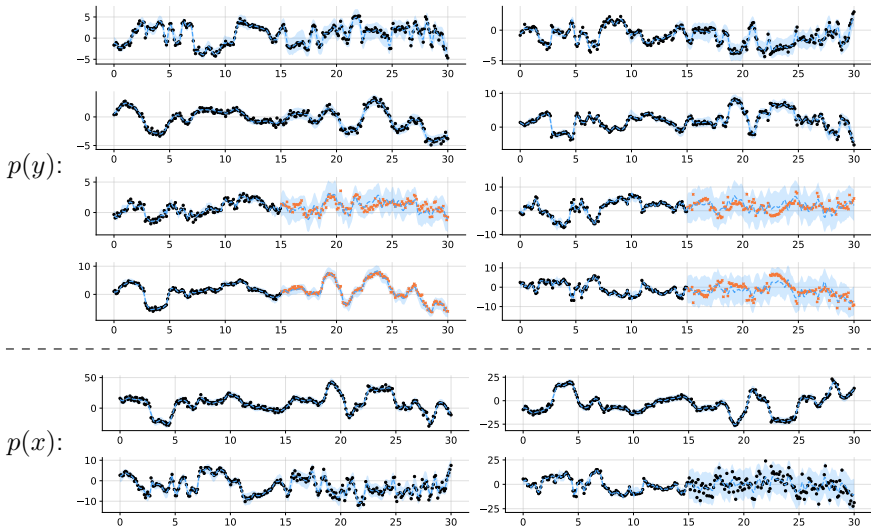
$x \sim \text{GPAR}$:

✗ Learning requires subtle approximation.

✓ Can produce exact posterior samples.

- $y(t) \in \mathbb{R}^8$ and $x(t) \in \mathbb{R}^4$.
- Markov GPAR: $x_i(t)$ depends nonlinearly only on $(t, x_{i-1}(t))$.
- Data generated by sample from model with random $H$.

- Task: Impute outputs 5–8 from outputs 1–4.
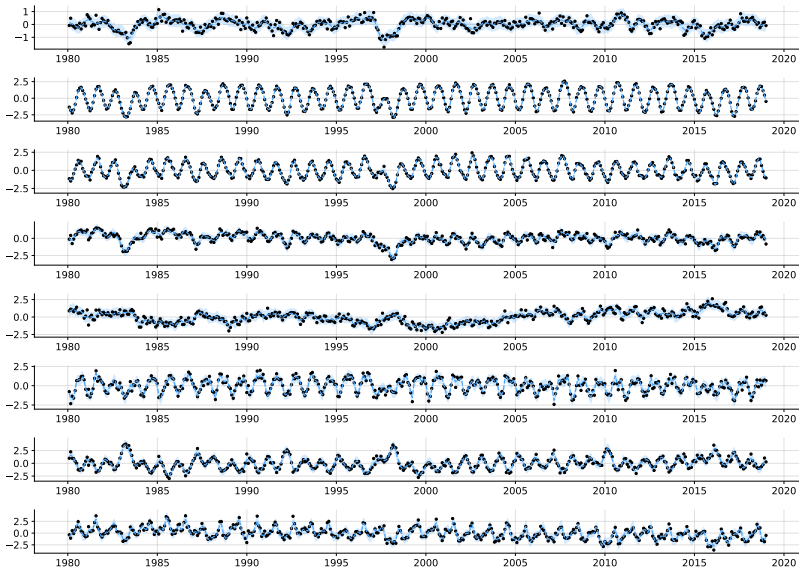- Parameters randomly initialised and learned.

$p(y)$:

$p(x)$:

- Average monthly temp. in Peru from 1980 to 2018 ($n = 468$).
- Measured at $p = 2808$ locations.
- GPAR with $m = 8$ outputs and nonlinear dependencies.
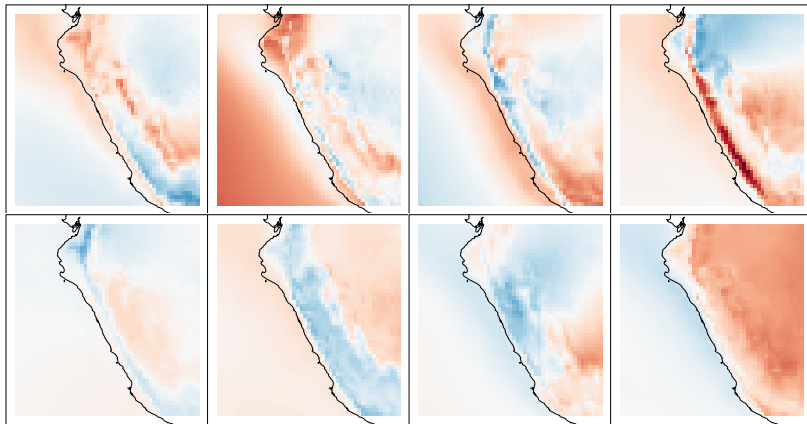- Basis $h_1, \ldots, h_8$ constrained to be orthogonal.

- Compose $p(x)$ with a likelihood $p(y \mid x)$ to scale to many outputs:

high-dim. model $\longleftarrow$ $p(y) = \int p(y \mid x)p(x)\,\mathrm{d}x.$

      maps into high-dim. space $\longleftarrow$     $\longrightarrow$ available model

- Likelihood is linear and uses orthogonal noise:

$$y(t) = h_1 x_1(t) + \cdots + h_m x_m(t) + \varepsilon_\perp(t).$$

$\checkmark$ Learning and inference for $p(y)$ use existing routines for $p(x)$.

These slides: https://wessel.page.link/compositional.