Invariances in Gaussian processes

And how to learn them

ST John PROWLER.io

Outline

- 1. What are invariances?
- 2. Why do we want to make use of them?
- 3. How can we construct invariant GPs?
- 4. Where invariant GPs are actually crucial
- 5. How can we figure out what invariances to employ?

What are invariances?

Function $f(\cdot)$ does not change under some transformation $\mathbf{x} \to t(\mathbf{x})$ i.e. $f(\mathbf{x}) = f(t(\mathbf{x}))$ for $\forall \mathbf{x} \in \mathcal{X} \quad \forall t \in \mathcal{T}$

Can be discrete or continuous

- Translation
- Rotation
- Reflection
- Permutation

Invariance under discrete translation

$$\mathbf{x} \to \mathbf{x} + n\mathbf{d}$$
 $n \in \mathbb{Z}$ $f(\mathbf{x}) = f(\mathbf{x} + n\mathbf{d})$ $\forall \mathbf{x} \in \mathcal{X}$

Periodic functions



Invariance under discrete translation



5/53

Invariance under discrete rotation

Density of water molecules as a function of (x, y) point in plane

1/6th of the plane already predicts the function value everywhere



Invariance under reflection

Solar elevation measured as function of azimuth (for different days)



7/53



[100, 200, 1, 1, 1]



[1, 200, 1, 100, 1]



f(100, 200, 1, 1, 1) = f(1, 200, 1, 100, 1)

Different inputs but same function value





Invariance under continuous transformations

Translation

 $\mathbf{x} \to \mathbf{x} + n\mathbf{d}$ $n \in \mathbb{R}$









Class label as a function of image pixel matrix



17/53



Example: molecular energy



Approximately invariant...



Approximately invariant...





2. Why do we want to use invariances?

- Incorporate prior knowledge about the behaviour of a system
 - *Physical* symmetries, e.g. modelling total energy (and gradients, i.e. forces) of a set of atoms
- Helps generalisation
- Improved accuracy vs number of training points















Constructing invariant GPs

We want a prior over functions that obey the chosen symmetry.

Symmetrise the function: Can do this by

- a) appropriate **mapping** to invariant space
- b) **sum** over transformations

Permutation-invariant GPs: mapping construction

$$\mathbf{x} = (x_1, x_2) \quad \rightarrow \quad \varphi(\mathbf{x}) = \begin{pmatrix} (x_1 + x_2)/2 \\ |x_1 - x_2|/2 \end{pmatrix}$$

$$f(x_1, x_2) = g((x_1 + x_2)/2, |x_1 - x_2|/2)$$

$$f(\mathbf{x}) = g(\varphi(\mathbf{x})) \quad \rightarrow \quad k_f(\mathbf{x}, \mathbf{x}') = k_g(\varphi(\mathbf{x}), \varphi(\mathbf{x}'))$$

Permutation-invariant GPs: sum construction

$$f(x_1, x_2) = g(x_1, x_2) + g(x_2, x_1)$$

$$(x_1, x_2) \to (x_2, x_1)$$
:
 $f(x_2, x_1) = g(x_2, x_1) + g(x_1, x_2)$

Invariant sum kernel

$$\begin{aligned} k_f((x_1, x_2), (x_1', x_2')) &= \operatorname{Cov} \left(f(x_1, x_2), f(x_1', x_2') \right) = \mathbb{E} \left[f(x_1, x_2) f(x_1', x_2') \right] \\ &= \mathbb{E} \left[(g(x_1, x_2) + g(x_2, x_1))(g(x_1', x_2') + g(x_2', x_1')) \right] \\ &= \mathbb{E} \left[g(x_1, x_2) g(x_1', x_2') \right] + \mathbb{E} \left[g(x_1, x_2) g(x_2', x_1') \right] \\ &+ \mathbb{E} \left[g(x_2, x_1) g(x_1', x_2') \right] + \mathbb{E} \left[g(x_2, x_1) g(x_2', x_1') \right] \\ &= k_g \left((x_1, x_2), (x_1', x_2') \right) + k_g \left((x_1, x_2), (x_2', x_1') \right) \\ &+ k_g \left((x_2, x_1), (x_1', x_2') \right) + k_g \left((x_2, x_1), (x_2', x_1') \right) \end{aligned}$$

Samples from the prior



How can we generalise this?

Symmetry group

Transformations can be **composed**:

$$f(x) = f(t(x)) \quad \forall \mathbf{x} \in \mathcal{X} \qquad \forall t \in \mathcal{T}$$
$$f(t(x)) = f(t'(t(x)))$$
$$\implies f(x) = f(t'(t(x)) = f((t' \circ t)(x))$$

Set of all compositions of transformations is a group; corresponds to symmetries

Orbit of *x*: all points reachable by transformations

$$\mathcal{O}(\mathbf{x}) = \{t(\mathbf{x}) | t \in \mathcal{T}\}$$

Example: Permutation in 2D

$$\pi(\mathbf{x}) : (x_1, x_2) \mapsto (x_2, x_1)$$
$$\mathcal{T} = \{I(\mathbf{x}), \pi(\mathbf{x})\}$$

$$\mathcal{O}((x_1, x_2)) = \mathcal{O}((x_2, x_1)) = \{(x_1, x_2), (x_2, x_1)\}$$

Examples of orbits: permutation invariance

Orbit size = 2



Examples of orbits: six-fold rotation invariance

Orbit size = 6


Examples of orbits: permutation and six-fold rotation

Orbit size = 12



Examples of orbit: continuous rotation symmetry

Uncountably infinite



Orbit of a periodic function in 1D

Countably infinite



Constructing invariant GPs: sum revisited

$\pi(\mathbf{x}) : (x_1, x_2) \mapsto (x_2, x_1)$ $\mathcal{T} = \{I(\mathbf{x}), \pi(\mathbf{x})\}$ $\mathcal{O}(\mathbf{x}) = \{t(\mathbf{x}) | t \in \mathcal{T}\}$

$$f(x_1, x_2) = g(x_1, x_2) + g(x_2, x_1) \implies f(\mathbf{x}) = \sum_{\mathbf{x}_o \in \mathcal{O}(\mathbf{x})} g(\mathbf{x}_o)$$

$$k_f(\mathbf{x}, \mathbf{x}') = \mathbb{E}\left[\sum_{\mathbf{x}_o \in \mathcal{O}(\mathbf{x})} g(\mathbf{x}_o) \sum_{\mathbf{x}_o' \in \mathcal{O}(\mathbf{x}')} g(\mathbf{x}_{o'})\right] = \sum_{\mathbf{x}_o \in \mathcal{O}(\mathbf{x})} \sum_{\mathbf{x}_o' \in \mathcal{O}(\mathbf{x}')} k_g(\mathbf{x}_o, \mathbf{x}_{o'})$$

Applications

Molecular modelling

Time-evolution of the configuration (position of all atoms) of a system of atoms/molecules

Need Potential Energy Surface (PES)! Gradients = forces (easy with GPs)

$$E(\mathbf{r}_1,\ldots,\mathbf{r}_N)$$

Potential Energy Surface



Approximate as sum over k-mers (many-body expansion)

Invariance to rotation/translation of local environment/k-mer

Invariance under permutation of equivalent atoms

Many-body expansion, sum over k-mers:

$$E(\mathbf{r}_1,\ldots,\mathbf{r}_N) = \sum_i E^{(1)}(\mathbf{r}_i) + \sum_{i\neq j} E^{(2)}(\mathbf{r}_i,\mathbf{r}_j) + \sum_{i\neq j\neq k} E^{(3)}(\mathbf{r}_i,\mathbf{r}_j,\mathbf{r}_k) + \cdots$$

$$E(\mathbf{r}_1, \dots, \mathbf{r}_N) = \sum_i E^{(1)}(\{\mathbf{r}_m : m \in \mathcal{M}_i\}) + \sum_{i \neq j} E^{(2)}(\{\mathbf{r}_m, \mathbf{r}_n : m \in \mathcal{M}_i, n \in \mathcal{M}_j\})$$
$$+ \sum_{i \neq j \neq k} E^{(3)}(\{\mathbf{r}_m, \mathbf{r}_n, \mathbf{r}_p : m \in \mathcal{M}_i, n \in \mathcal{M}_j, p \in \mathcal{M}_k\}) + \cdots$$

Invariance to rotation/translation of local environment/k-mer:

Map to interatomic distances



$$E^{(1)}(\{\mathbf{r}_{m}: m \in \mathcal{M}\}) = E^{(1)}(\{|\mathbf{r}_{m} - \mathbf{r}_{n}|: m, n \in \mathcal{M}\})$$
$$E^{(2)}(\{\mathbf{r}_{m}, \mathbf{r}_{n}: m \in \mathcal{M}_{1}, n \in \mathcal{M}_{2}\}) = E^{(2)}(\{|\mathbf{r}_{m} - \mathbf{r}_{n}|: m, n \in (\mathcal{M}_{1} \cup \mathcal{M}_{2})\})$$
46/53

Invariance under permutation of equivalent atoms:

sum over them!



How can we find out if an invariance is helpful?

- As usual (like another kernel hyperparameter): *marginal* likelihood
- Unlike "regular" likelihood (equivalent to training-set RMSE):
 - Less overfitting
 - Related to generalisation

Marginal likelihood and generalisation

Measures how well part of the training set predicts the other training points:

$$p(\mathbf{y}|\theta) = p(\mathbf{y}_1|\theta)p(\mathbf{y}_2|\mathbf{y}_1,\theta)p(\mathbf{y}_3|\mathbf{y}_{1:2},\theta)\prod_{c=4}^C p(\mathbf{y}_c|\mathbf{y}_{1:c-1},\theta)$$

= how accurately the model generalises during inference, similar to cross-validation (but differentiable)

Marginal likelihood



Summary: we have seen...

How to constrain GPs to give **invariant functions**

When invariance **improves** a model's **generalisation**

When invariance increases the marginal likelihood

That invariances exist in **real-world problems**

Questions?

Next up: how to learn invariances...

Snowflake prior



Why not just data augmentation?

Used in deep learning...

Invariances are better:

- 1. Cubic scaling with number of data points vs **linear scaling** with invariances in prior
- 2. Data augmentation results in same predictive mean, but not variance
- 3. Invariances in the GP prior give us **invariant samples**

Learning Invariances with the Marginal Likelihood

Mark van der Wilk PROWLER.io

We discussed...



How to constrain GPs to give **invariant functions**



When invariance **improves** a model's **generalisation**



When invariance **increases the marginal likelihood**



From known invariances to learning them

We previously saw that known invariances were useful to modelling.

- How do we *exploit invariances* in a problem, if we *don't know them a-priori*?
- Can we *learn* a useful invariance *from the data*?

Model selection

• Invariances in a GP are *expressed in the kernel*

$$k_{\theta}(\mathbf{x}, \mathbf{x}') = \sum_{\mathbf{x}_a \in \mathcal{O}_{\theta}(\mathbf{x})} \sum_{\mathbf{x}'_a \in \mathcal{O}_{\theta}(\mathbf{x}')} k_g(\mathbf{x}_a, \mathbf{x}'_a)$$

- We use the marginal likelihood to select models $p(\mathbf{y} | X, \theta) = \mathcal{N}(\mathbf{y}; 0, \mathbf{K_{ff}} + \sigma^2 I), \quad [\mathbf{K_{ff}}]_{ij} = k_{\theta}(\mathbf{x}_i, \mathbf{x}_j)$ $\theta^* = \operatorname*{argmax}_{\theta} p(\mathbf{y} | X, \theta)$
- **Parameterising** the orbit is all that is left

Parameterising orbits is hard

Strict invariance requires:

 $f(\mathbf{x}) = f(t(\mathbf{x})) \qquad \forall \mathbf{x} \in \mathcal{X} \qquad \forall t \in \mathcal{T}$

which we can obtain using the construction

$$f(\mathbf{x}) = \sum_{\mathbf{x}_a \in \mathcal{O}(\mathbf{x})} g(\mathbf{x}_a)$$

I don't know how to parameterise orbits!

From orbits to distributions

• We sum over an arbitrary set of points $f(\mathbf{x}) = \sum_{\mathbf{x}_a \in \mathcal{A}(\mathbf{x})} g(\mathbf{x}_a)$ • Take the infinite limit $f(\mathbf{x}) = \int g(\mathbf{x}_a) p(\mathbf{x}_a \mid \mathbf{x}) d\mathbf{x}_a$ • Find kernel $k_{\theta}(\mathbf{x}, \mathbf{x}') = \iint k_g(\mathbf{x}_a, \mathbf{x}'_a) p_{\theta}(\mathbf{x}_a \mid \mathbf{x}) p_{\theta}(\mathbf{x}'_a \mid \mathbf{x}') d\mathbf{x}_a d\mathbf{x}'_a$

I do know how to parameterise distributions!

Insensitivity

• We lose exact invariance... but this may be a blessing!



Insensitivity

• We lose exact invariance... but this may be a blessing!



What we will do

• Parameterise a distribution that describes the insensitivity $p_{\theta}(\mathbf{x}_a \mid \mathbf{x})$

• Use this distribution to *define a kernel*

$$k_{\theta}(\mathbf{x}, \mathbf{x}') = \iint k_{g}(\mathbf{x}_{a}, \mathbf{x}'_{a}) p_{\theta}(\mathbf{x}_{a} \mid \mathbf{x}) p_{\theta}(\mathbf{x}'_{a} \mid \mathbf{x}') \mathrm{d}\mathbf{x}_{a} \mathrm{d}\mathbf{x}'_{a}$$

• Find invariance in the kernel by *optimising the hyperparameters*

$$p(\mathbf{y} | X, \theta) = \mathcal{N}(\mathbf{y}; 0, \mathbf{K}_{\mathbf{ff}} + \sigma^2 I), \qquad [\mathbf{K}_{\mathbf{ff}}]_{ij} = k_{\theta}(\mathbf{x}_i, \mathbf{x}_j)$$
$$\theta^* = \operatorname*{argmax}_{\theta} p(\mathbf{y} | X, \theta)$$

Obstacles to inference

$$k_{\theta}(\mathbf{x}, \mathbf{x}') = \iint k_{g}(\mathbf{x}_{a}, \mathbf{x}_{a}') p_{\theta}(\mathbf{x}_{a} | \mathbf{x}) p_{\theta}(\mathbf{x}_{a}' | \mathbf{x}') d\mathbf{x}_{a} d\mathbf{x}_{a}'$$

$$p(\mathbf{y} | X, \theta) = \int p(\mathbf{y} | f(X)) p(f(X)) df(X)$$

$$p(f(X)) = \mathcal{N}(0, \mathbf{K_{ff}}), \quad [\mathbf{K_{ff}}]_{ij} = k_{\theta}(\mathbf{x}_{i}, \mathbf{x}_{j})$$

$$\theta^{*} = \underset{\theta}{\operatorname{argmax}} p(\mathbf{y} | X, \theta)$$

- 1. For large datasets, the matrix operations of K_{ff} becomes infeasible (O(N³) time complexity)
- 2. We may have non-Gaussian likelihoods (classification!)
- 3. We can't even evaluate the kernel!

Variational inference

$$\mathcal{L} = \sum_{n=1}^{N} \mathbb{E}_{q(f(\mathbf{x}_n))}[\log p(y_n \mid f(\mathbf{x}_n))] - \mathrm{KL}[q(\mathbf{u})||p(\mathbf{u})]$$

$$q(f(\mathbf{x}_n)) = \mathcal{N}(f(\mathbf{x}_n), \mu_n, \sigma_n^2)$$

$$\mu_n = \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{uu}}^{-1} \mathbf{m}$$

$$\sigma_n^2 = k(\mathbf{x}_n, \mathbf{x}_n) - \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{uu}}^{-1} (\mathbf{K}_{\mathbf{uu}} - \mathbf{S}) \mathbf{K}_{\mathbf{uu}}^{-1} \mathbf{k}_{\mathbf{un}}$$

$$k_{\theta}(\mathbf{x}, \mathbf{x}') = \iint k_g(\mathbf{x}_a, \mathbf{x}'_a) p_{\theta}(\mathbf{x}_a \mid \mathbf{x}) p_{\theta}(\mathbf{x}'_a \mid \mathbf{x}') \mathrm{d}\mathbf{x}_a \mathrm{d}\mathbf{x}'_a$$

For large datasets, the matrix operations of K_{ff} becomes infeasible (O(N³) time complexity)
 We may have non Gaussian likelihoods (classification!)

3. We can't even evaluate the kernel! **Still needed for K**_{uu} and k_{un}

• Variational posterior is constructed by conditioning $q(f(\mathbf{x}_n)) = \int p(f(\mathbf{x}_n) \,|\, \mathbf{u}) q(\mathbf{u}) \mathrm{d}\mathbf{u}$ $\mathbf{u} = \{f(\mathbf{z}_m)\}_{m=1}^N$

• Gaussian conditioning requires covariances

$$Cov[f(\mathbf{x}_n), u_m] = \iint k_g(\mathbf{x}_a, \mathbf{z}_a) p_\theta(\mathbf{x}_a | \mathbf{x}_n) p_\theta(\mathbf{z}_a | \mathbf{z}_m) d\mathbf{x}_a d\mathbf{z}_a$$

$$Cov[u_m, u_{m'}] = \iint k_g(\mathbf{z}_a, \mathbf{z}'_a) p_\theta(\mathbf{z}_a | \mathbf{z}_m) p_\theta(\mathbf{z}'_a | \mathbf{z}_{m'}) d\mathbf{z}_a d\mathbf{z}'_a$$

$$\mu_n = \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{uu}}^{-1} \mathbf{m}$$

$$\sigma_n^2 = k(\mathbf{x}_n, \mathbf{x}_n) - \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{uu}}^{-1} (\mathbf{K}_{\mathbf{uu}} - \mathbf{S}) \mathbf{K}_{\mathbf{uu}}^{-1} \mathbf{k}_{\mathbf{un}}$$



- Variational posterior is constructed by **conditioning** $q(f(\mathbf{x}_n)) = \int p(f(\mathbf{x}_n) | \mathbf{u}) q(\mathbf{u}) d\mathbf{u}$ $\mathbf{u} = \{\underline{g(\mathbf{z}_m)}\}_{m=1}^N$
- Gaussian conditioning requires covariances

 $\operatorname{Cov}[f(\mathbf{x}_n), u_m] = \mathbb{E}_g[f(\mathbf{x}_n)g(\mathbf{z}_m)]$

 $\operatorname{Cov}[u_m, u_{m'}] = \mathbb{E}_g[g(\mathbf{z}_m)g(\mathbf{z}_{m'})]$

$$\mu_n = \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{m}$$

$$\sigma_n^2 = k(\mathbf{x}_n, \mathbf{x}_n) - \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} (\mathbf{K}_{\mathbf{u}\mathbf{u}} - \mathbf{S}) \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{k}_{\mathbf{u}\mathbf{u}}$$



• Variational posterior is constructed by **conditioning** $q(f(\mathbf{x}_n)) = \int p(f(\mathbf{x}_n) | \mathbf{u}) q(\mathbf{u}) d\mathbf{u}$ $\mathbf{u} = \{\underline{g(\mathbf{z}_m)}\}_{m=1}^N$

• Gaussian conditioning requires **covariances** $\operatorname{Cov}[f(\mathbf{x}_n), u_m] = \mathbb{E}_g \left[\int g(\mathbf{x}_a) p(\mathbf{x}_a \mid \mathbf{x}_n) \mathrm{d}\mathbf{x}_a g(\mathbf{z}_m) \right]$

$$\operatorname{Cov}[u_m, u_{m'}] = \mathbb{E}_g[g(\mathbf{z}_m)g(\mathbf{z}_{m'})]$$

$$\mu_n = \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{m}$$

$$\sigma_n^2 = k(\mathbf{x}_n, \mathbf{x}_n) - \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} (\mathbf{K}_{\mathbf{u}\mathbf{u}} - \mathbf{S}) \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{k}_{\mathbf{u}\mathbf{u}}$$



• Variational posterior is constructed by **conditioning** $q(f(\mathbf{x}_n)) = \int p(f(\mathbf{x}_n) | \mathbf{u}) q(\mathbf{u}) d\mathbf{u}$ $\mathbf{u} = \{\underline{g(\mathbf{z}_m)}\}_{m=1}^N$

• Gaussian conditioning requires covariances

 $\operatorname{Cov}[f(\mathbf{x}_n), u_m] = \int k_g(\mathbf{x}_a, \mathbf{z}_m) p_\theta(\mathbf{x}_a \mid \mathbf{x}_n) \mathrm{d}\mathbf{x}_a$

 $\operatorname{Cov}[u_m, u_{m'}] = k_g(\mathbf{z}_m, \mathbf{z}_{m'})$

$$\mu_n = \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{m}$$

$$\sigma_n^2 = k(\mathbf{x}_n, \mathbf{x}_n) - \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} (\mathbf{K}_{\mathbf{u}\mathbf{u}} - \mathbf{S}) \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{k}_{\mathbf{u}\mathbf{u}}$$



Unbiased estimation of the kernel

$$\mathcal{L} = \sum_{n=1}^{N} \mathbb{E}_{q(f(\mathbf{x}_n))} [\log p(y_n | f(\mathbf{x}_n))] - \mathrm{KL}[q(\mathbf{u})|| p(\mathbf{u})$$
$$= \sum_{n=1}^{N} \left[-\frac{1}{2} \log 2\pi\sigma^2 - \frac{(y_n - \mu_n)^2 + \sigma_n^2}{\sigma^2} \right] - \mathrm{KL}$$
$$\mu_n = \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{uu}}^{-1} \mathbf{m}$$
$$\sigma_n^2 = k(\mathbf{x}_n, \mathbf{x}_n) - \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{uu}}^{-1} (\mathbf{K}_{\mathbf{uu}} - \mathbf{S}) \mathbf{K}_{\mathbf{uu}}^{-1} \mathbf{k}_{\mathbf{un}}$$

Unbiased estimates of μ_n , μ_n^2 , σ_n^2 , give unbiased estimate of the ELBO!

Unbiased estimation of the kernel

$$\mu_{n} = \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{m}$$

$$\sigma_{n}^{2} = k(\mathbf{x}_{n}, \mathbf{x}_{n}) - \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} (\mathbf{K}_{\mathbf{u}\mathbf{u}} - \mathbf{S}) \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{k}_{\mathbf{u}n}$$

$$\operatorname{Cov}[f(\mathbf{x}_{n}), u_{m}] = \int k_{g}(\mathbf{x}_{a}, \mathbf{z}_{m}) p_{\theta}(\mathbf{x}_{a} \mid \mathbf{x}_{n}) \mathrm{d}\mathbf{x}_{a}$$

$$k_{\theta}(\mathbf{x}, \mathbf{x}') = \iint k_{g}(\mathbf{x}_{a}, \mathbf{x}'_{a}) p_{\theta}(\mathbf{x}_{a} \mid \mathbf{x}) p_{\theta}(\mathbf{x}'_{a} \mid \mathbf{x}') \mathrm{d}\mathbf{x}_{a} \mathrm{d}\mathbf{x}'_{a}$$

Unbiased estimation of the kernel

$$\mu_{n} = \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{m}$$

$$\sigma_{n}^{2} = k(\mathbf{x}_{n}, \mathbf{x}_{n}) - \mathbf{k}_{n\mathbf{u}} \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} (\mathbf{K}_{\mathbf{u}\mathbf{u}} - \mathbf{S}) \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{k}_{\mathbf{u}n}$$

$$\operatorname{Cov}[f(\mathbf{x}_{n}), u_{m}] = \int k_{g}(\mathbf{x}_{a}, \mathbf{z}_{m}) p_{\theta}(\mathbf{x}_{a} \mid \mathbf{x}_{n}) d\mathbf{x}_{a}$$

$$k_{\theta}(\mathbf{x}, \mathbf{x}') = \iint k_{g}(\mathbf{x}_{a}, \mathbf{x}'_{a}) p_{\theta}(\mathbf{x}_{a} \mid \mathbf{x}) p_{\theta}(\mathbf{x}'_{a} \mid \mathbf{x}') d\mathbf{x}_{a} d\mathbf{x}'_{a}$$
sample $p_{\theta}(\mathbf{x}_{a} \mid \mathbf{x})$

(We only need to sample one set from $p_{\theta}(\mathbf{x}_{a} \mid \mathbf{x})$, see paper for details)
What we did

• Parameterise a distribution that describes the insensitivity $p_{\theta}(\mathbf{x}_a \mid \mathbf{x})$

- Use this distribution to *define a kernel* $k(\mathbf{x}, \mathbf{x}') = \iint k_g(\mathbf{x}_a, \mathbf{x}'_a) p(\mathbf{x}_a \mid \mathbf{x}) p(\mathbf{x}'_a \mid \mathbf{x}') \mathrm{d}\mathbf{x}_a \mathrm{d}\mathbf{x}'_a$
- Approximate the marginal likelihood using *the variational evidence lower bound (ELBO)*
- Find an unbiased ELBO approximation, using unbiased estimates of the kernel
- Optimise the hyperparameters, using the gradients of the ELBO

Results

 Single model tunes itself automatically to multiple datasets



Method	Error in %
RBF	2.15 ± 0.03
rot. only	2.08 ± 0.06
all affine	1.35 ± 0.07
local def.	1.47 ± 0.05

• Fire off optimisation and watch it go



Conclusions & outlook

- We can *parameterise invariant kernels*
- We can *learn the parameters* with a marginal likelihood approximation
- Learned invariances *improve generalisation*