

Non-degenerate Priors for Arbitrarily Deep Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

A good latent representation captures all the relevant degrees of freedom of the manifold on which the data live. We show, for typical deep architectures, that as the number of layers increase, the representational capacity of the model tends to capture fewer degrees of freedom. In the limit, deep representations only retain a single degree of freedom locally. In addition, gradient-based learning becomes intractable. We propose several solutions to address these pathologies.

1 The Jacobian of Deep GPs

The derivatives of a function drawn from a GP prior with an isotropic SE kernel are i.i.d. Normal Because differentiation is a linear operator, the derivatives of a function drawn from a GP prior are also jointly Gaussian distributed, with covariance between derivatives w.r.t. different dimensions of \mathbf{x} given by:

$$\text{cov} \left(\frac{\partial f(\mathbf{x})}{\partial x_{d_1}}, \frac{\partial f(\mathbf{x})}{\partial x_{d_2}} \right) = \frac{\partial^2 k(\mathbf{x}, \mathbf{x}')}{\partial x_{d_1} \partial x'_{d_2}} \Big|_{\mathbf{x}=\mathbf{x}'} \quad (1)$$

[cite carl's paper?] If our kernel is a product over individual dimensions $k(\mathbf{x}, \mathbf{x}') = \prod_d k_d(x_d, x'_d)$, as in the case of the isotropic squared-exp kernel, then the diagonal covariances are given by $\frac{\sigma_d^2}{\ell_d^2}$, and the off-diagonal entries are zero. This means that elements are independent and identically distributed.

The elements of the Jacobian of a GP with an isotropic SE kernel are i.i.d. Gaussians The Jacobian of the ℓ^{th} function is:

$$J_{\mathbf{x} \rightarrow \mathbf{y}}^\ell(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1^\ell(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_1^\ell(\mathbf{x})}{\partial x_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_D^\ell(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_D^\ell(\mathbf{x})}{\partial x_D} \end{bmatrix} \quad (2)$$

Because we've assumed that the GP on each output dimension $f_d(\mathbf{x}) \sim \mathcal{GP}$ is independent, it follows that for a given \mathbf{x} , each row of $J_{\mathbf{x} \rightarrow \mathbf{y}}(\mathbf{x})$ is independent. Above, we showed that the elements of each row are independent. This means that each entry in the Jacobian of a GP-distributed transformation is i.i.d. Normal.

We also have that if $\mathbf{x} = f(\mathbf{y})$, then $J_{\mathbf{y} \rightarrow \mathbf{x}}^{-1}(\mathbf{y}) = J_{\mathbf{x} \rightarrow \mathbf{y}}(\mathbf{x})$.

The Jacobian of a deep GP is a product of random normal matrices By the multivariate chain rule, the derivative (Jacobian) of any compositions of functions is simply the product of the Jacobians

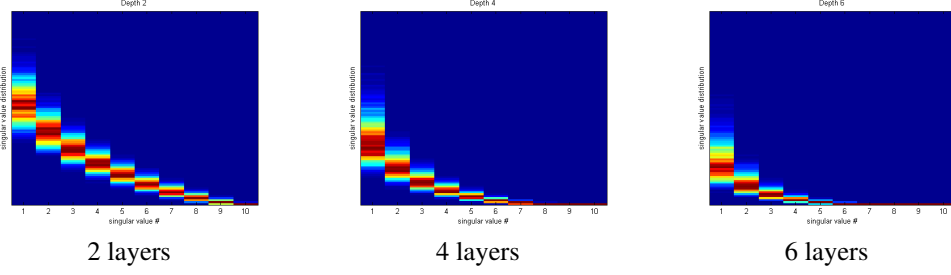


Figure 1: Eigenspectrum of the Jacobian of a deep GP. As the net gets deeper, the largest eigenvalue comes to dominate. This implies that there is only one degree of freedom in representation being computed.

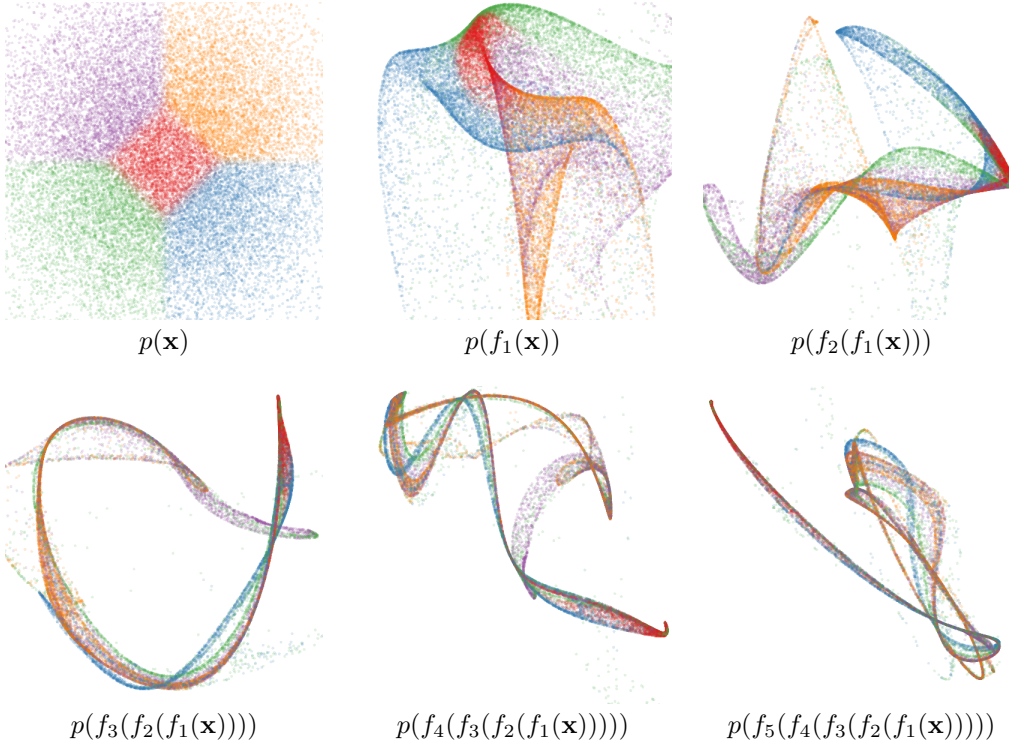


Figure 2: Draws from a deep GP. A distribution is warped by successive functions drawn from a GP prior. As the number of layers increases, the density exhibits a sort of filamentation.

of each function. and the Jacobian of the composed (deep) function is:

$$J^{1:L}(x) = \prod_{\ell=1}^L J^L(x) \quad (3)$$

1.1 Filamentation

Attempted definition of a filament A continuous, twice-differentiable region of a pdf is called a *filament* to the degree that, weighted by density, only one eigenvalue of the Hessian of the pdf are is small relative to the average eigenvalue.

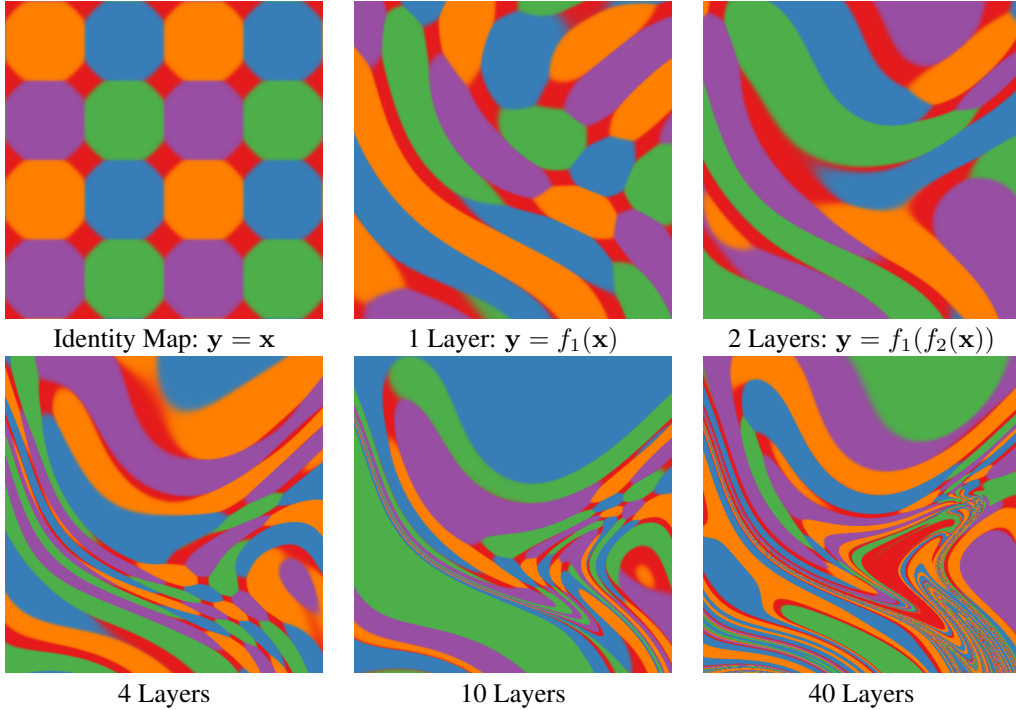


Figure 3: Feature Mapping of a deep GP. Shown here are the colors corresponding to the location $\mathbf{y} = f(\mathbf{x})$ that each point is mapped to after being warped by a deep GP. This figure can be seen as the inverse of figure 2. Just as the densities in 2 became locally one-dimensional, there is locally only one direction that one can move \mathbf{x} in to change \mathbf{y} .

2 Fixing the filamentation pathology

We can fix the pathologies exhibited in figures 2 and 3 by simply making each layer of computation depend not only on the output of the previous layer, but also on the original input \mathbf{x} . Draws from the resulting priors are shown in figures 4 and 5.

3 Conclusions

Deep neural networks and deep Gaussian processes are analyzable using random matrix theory. After proving that the Jacobian is an i.i.d. Gaussian matrix, many other forms of

If you want to use very deep nets, you won't be able to do so if you initialize/regularize all your weights independently We might want to think about different ways of

If you initialize independently, the density becomes fractal Points close in x -space can be very far in y -space, and vice versa.

A spikey eigenspectrum will lead to saturation Maybe we should initialize differently in order to avoid such saturation, like Martens' sparse initialization: http://www.cs.toronto.edu/~jmartens/docs/Deep_HessianFree.pdf

4 Arbitrarily Deep Kernels

[cite Youngmin Cho's thesis, and Radford Neal's]

One can derive a Gaussian process as a neural network: $f(x) = \alpha^T \Phi(x) = \sum_{i=1}^K \alpha_i \phi_i(x)$. We can consider applying the feature transform $\Phi(\cdot)$ to the features themselves: $\Phi_2 = \Phi(\Phi(\mathbf{x}))$,

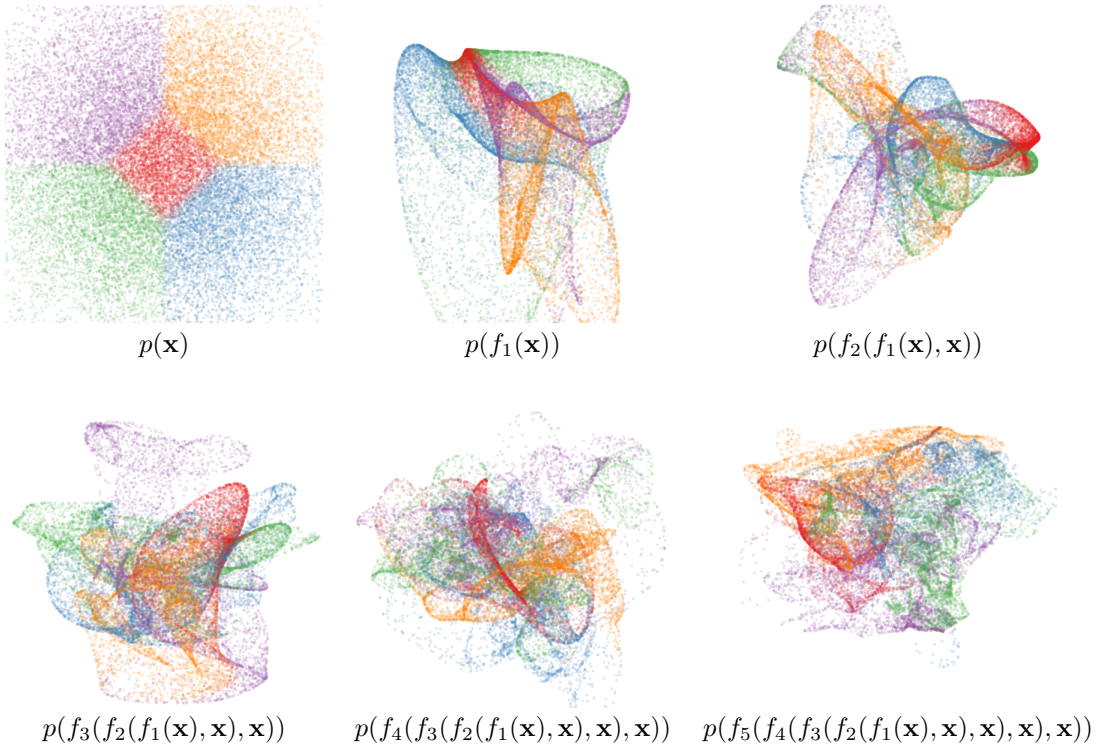


Figure 4: Draws from a deep GP, with each layer connected to the input \mathbf{x} . By always depending on the original input, the density becomes more complex without concentrating along filaments.

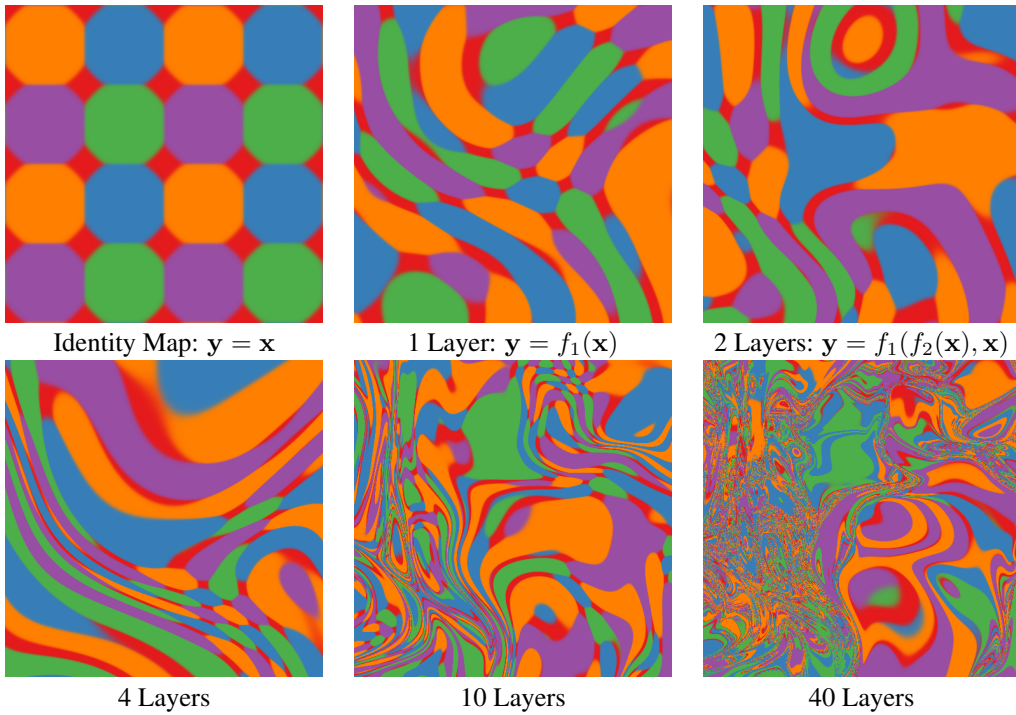


Figure 5: Feature Mapping of a deep GP with each layer connected to the input \mathbf{x} . Just as the densities in 4 became remain locally two-dimensional even after many transformations, in this mapping there are locally usually two directions that one can move \mathbf{x} in to change \mathbf{y} .

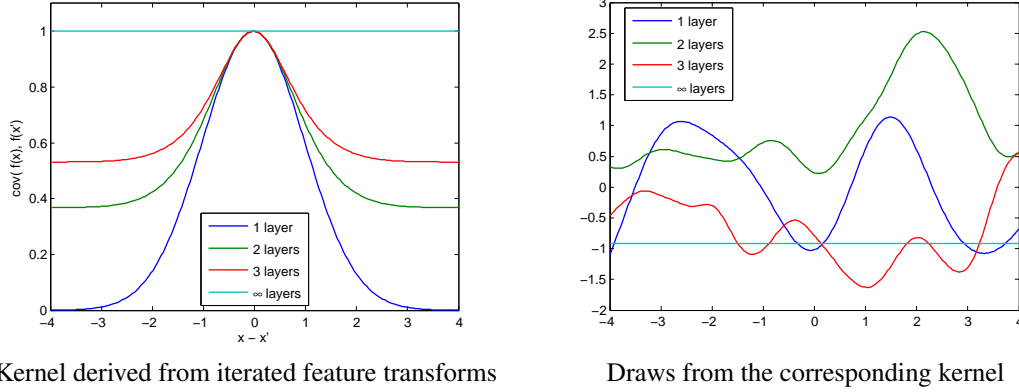


Figure 6: A degenerate kernel produced by repeatedly applying a feature transform.

In this section, we derive a kernel which corresponds to arbitrarily many compositions of the feature vectors corresponding to the squared-exp kernel:

$$k_2(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}\|\Phi(\mathbf{x}) - \Phi(\mathbf{x}')\|_2^2\right) \quad (4)$$

$$k_2(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}\sum_i [\phi_i(\mathbf{x}) - \phi_i(\mathbf{x}')]^2\right) \quad (5)$$

$$k_2(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}\sum_i [\phi_i(\mathbf{x})^2 - 2\phi_i(\mathbf{x})\phi_i(\mathbf{x}') + \phi_i(\mathbf{x}')^2]\right) \quad (6)$$

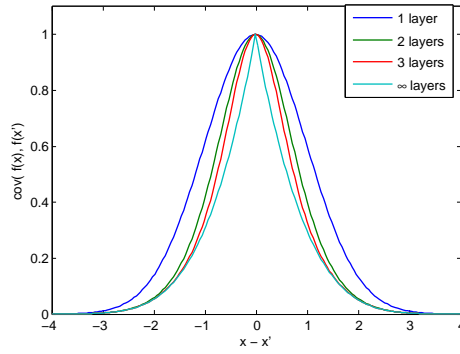
$$k_2(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}\left[\sum_i \phi_i(\mathbf{x})^2 - 2\sum_i \phi_i(\mathbf{x})\phi_i(\mathbf{x}') + \sum_i \phi_i(\mathbf{x}')^2\right]\right) \quad (7)$$

$$k_2(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}[k_1(\mathbf{x}, \mathbf{x}) - 2k_1(\mathbf{x}, \mathbf{x}') + k_1(\mathbf{x}', \mathbf{x}')] \right) \quad (8)$$

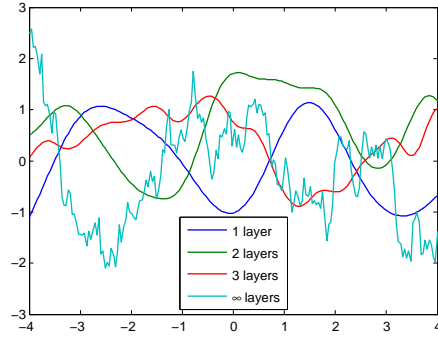
$$k_2(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}') - 1) \quad (9)$$

Thus if $k_1(x, y) = e^{-\|x-y\|^2/2}$, then the two-layer kernel is simply $k_2(x, y) = e^{k_1(x, y)-1}$. This formula is true for every layer: $k_{n+1}(x, y) = e^{k_n(x, y)-1}$. Note that nothing in this derivation depends on details of k_1 , except that $k_1(\mathbf{x}, \mathbf{x}) = 1$. Because this is true for k_2 as well, this recursion holds in general, and we have that $k_{n+1}(x, y) = e^{k_n(x, y)-1}$. In the infinite limit, this recursion converges to $k(x, y) = 1$ for all inputs.

Figure 6 shows this kernel at different depths, including the degenerate limit. One interpretation of why repeated feature transforms lead to this degenerate prior is that each layer can only lose information about the previous set of features. In the limit, the transformed features contain no information about the original input \mathbf{x} . Since the function doesn't depend on its input, it must be the same everywhere.



Kernel derived from iterated feature transforms
with all layers connected to the input



Draws from the corresponding kernel

Figure 7: A non-degenerate version of the infinitely deep feature transform kernel. By connecting the inputs \mathbf{x} to each layer, the function can still depend on its input even after arbitrarily many layers of computation.

4.1 Fixing the deep kernel

Follow a suggestion from Radford Neal’s thesis, we connect the inputs to each layer of features. We do this simply by augmenting the feature vector $\Phi_n(\mathbf{x})$ with the extra features \mathbf{x} at every layer:

$$k_{n+1}(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{1}{2} \left\| \begin{bmatrix} \Phi_n(\mathbf{x}) \\ \mathbf{x} \end{bmatrix} - \begin{bmatrix} \Phi_n(\mathbf{x}') \\ \mathbf{x}' \end{bmatrix} \right\|_2^2 \right) \quad (10)$$

$$k_{n+1}(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{1}{2} \sum_i [\phi_i(\mathbf{x}) - \phi_i(\mathbf{x}')]^2 - \frac{1}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2 \right) \quad (11)$$

$$k_{n+1}(\mathbf{x}, \mathbf{x}') = \exp \left(k_1(\mathbf{x}, \mathbf{x}') - 1 - \frac{1}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2 \right) \quad (12)$$

Thus, this kernel satisfies the recurrence $k - \log(k) = 1 + \frac{1}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2$.

Properties of this kernel The solution to this recurrence has no closed form, but it is continuous and differentiable everywhere except at $\mathbf{x} = \mathbf{x}'$.

Conjectures:

- Samples from a GP with this prior are not differentiable.
- This kernel has smaller covariance than the squared-exp everywhere except at $\mathbf{x} = \mathbf{x}'$.
- Samples from this kernel are fractal.

5 Related Work

[Ryan Adams, Wallach and Zoubin on nonparametric deep nets]

[Deep GP Kernels by Youngmin Cho] http://cseweb.ucsd.edu/~yoc002/paper/thesis_youngmincho.pdf

[GP Dynamical systems]

[Warping a 1d uniform distribution]

Layer-wise analysis of deep networks with Gaussian kernels: http://books.nips.cc/papers/files/nips23/NIPS2010_0206.pdf

Deep Gaussian Processes We introduce a generative non-parametric model to address this problem. Our approach is based on the GP-LVM [1, 2, 3], a flexible nonparametric density model.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

Acknowledgments

We would like to thank Andrew McHutchon, Neil Lawrence, Josh Tenenbaum, Andreas Diamanadou, James Lloyd, and Mark van der Wilk for helpful discussions.

References

[1] N.D. Lawrence. Gaussian process latent variable models for visualisation of high dimensional data. *Advances in Neural Information Processing Systems*, 16:329–336, 2004.

[2] M. Salzmann, R. Urtasun, and P. Fua. Local deformation models for monocular 3D shape recovery. In *IEEE Conference on Computer Vision and Pattern Recognition*, CVPR, pages 1–8, 2008.

[3] N.D. Lawrence and R. Urtasun. Non-linear matrix factorization with Gaussian processes. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 601–608. ACM, 2009.