Introduction to Gaussian Processes

Neil D. Lawrence

MLSS, Cadiz 12th May 2016



Outline

Gaussian Processes

GP Non-Gaussian

GP Limitations

Kalman Filter

Dimensionality Reduction

Outline

Gaussian Processes

GP Non-Gaussian

GP Limitations

Kalman Filter

Dimensionality Reduction



Rasmussen and Williams (2006)

y = mx + c















y = mx + c

point 1:
$$x = 1, y = 3$$

 $3 = m + c$
point 2: $x = 3, y = 1$
 $1 = 3m + c$
point 3: $x = 2, y = 2.5$
 $2.5 = 2m + c$



riens. L'opinion contraire est une illusion de l'esprit qui, perdant de vue les raisons fugitives du choix de la volonté dans les choses indifférentes, se persuade qu'elle s'est déterminée d'ellemême et sans motifs.

Nous devons donc envisager l'état présent de l'univers, comme l'effet de son état antérieur, et comme la cause de celui qui va suivre. Une intelligence qui, pour un instant donné, connaîtrait toutes les forces dont la nature est animée, et la situation respective des êtres qui la composent, si d'ailleurs elle était assez vaste pour soumettre ces données à l'analyse, embrasserait dans la même formule les mouvemens des plus grands corps de l'univers et ceux du plus léger atome : rien ne serait incertain pour elle, et l'avenir comme le passé, serait présent à ses yeux. L'esprit humain offre, dans la perfection qu'il a su donner à l'Astronomie, une faible esquisse de cette intelligence. Ses découvertes en Mécanique et en Géométrie, jointes à celle de la pesanteur universelle, l'ont mis à portée de comprendre dans les mêmes expressions analytiques , les états passés et futurs du système du monde. En appliquant la même méthode à quelques autres objets de ses connaissances , il est parvetiu à ramener à des lois générales, les phénomènes observés, et à prévoir ceux que des circonstances données doivent faire éclore. Tous ces efforts dans la recherche de la vérité, tendent à le rapprocher sans cesse de l'intelligence que nous venons de concevoir, mais dont il restera toujours infiniment éloigné. Cette tendance propre à l'espèce humaine, est ce qui la rend supérieure aux animaux; et ses progrès en ce genre, distinguent les nations et les siècles, et font leur véritable gloire.

Rappelons-nous qu'autrefois, et à une époque qui

4 A PHILOSOPHICAL ESSAY ON PROBABILITIES.

other, we say that its choice is an effect without a cause. It is then, says Leibnitz, the blind chance of the Epicureans. The contrary opinion is an illusion of the mind, which, losing sight of the evasive reasons of the choice of the will in indifferent things, believes that choice is determined of itself and without motives.

We ought then to regard the present state of the universe as the effect of its anterior state and as the cause of the one which is to follow. Given for one instant an intelligence which could comprehend all the forces by which nature is animated and the respective situation of the beings who compose it-an intelligence sufficiently vast to submit these data to analysis-it would embrace in the same formula the movements of the greatest bodies of the universe and those of the lightest atom; for it, nothing would be uncertain and the future, as the past, would be present to its eyes. The human mind offers, in the perfection which it has been able to give to astronomy, a feeble idea of this intelligence. Its discoveries in mechanics and geometry, added to that of universal gravity, have enabled it to comprehend in the same analytical expressions the past and future states of the system of the world. Applying the same method to some other objects of its knowledge, it has succeeded in referring to general laws observed phenomena and in foreseeing those which given circumstances ought to produce. All these efforts in the search for truth tend to lead it back continually to the vast intelligence which we have just mentioned. but from which it will always remain infinitely removed. This tendency, peculiar to the human race, is that which renders it superior to animals; and their progress

6 A PHILOSOPHICAL ESSAY ON PROBABILITIES.

height: "The day will come when, by study pursued through several ages, the things now concealed will appear with evidence; and posterity will be astonished that truths so clear had escaped us." Clairaut then undertook to submit to analysis the perturbations which the comet had experienced by the action of the two great planets, Jupiter and Saturn; after immense calculations he fixed its next passage at the perihelion toward the beginning of April, 1759, which was actually verified by observation. The regularity which astronomy shows us in the movements of the comets doubtless exists also in all phenomena.

The curve described by a simple molecule of air or vapor is regulated in a manner just as certain as the planetary orbits; the only difference between them is that which comes from our ignorance.

Probability is relative, in part to this ignorance, in part to our knowledge. We know that of three or a greater number of events a single one ought to occur; but nothing induces us to believe that one of them will occur rather than the others. In this state of indecision it is impossible for us to announce their occurrence with certainty. It is, however, probable that one of these events, chosen at will, will not occur because we see several cases equally possible which exclude its occurrence, while only a single one favors it.

The theory of chance consists in reducing all the events of the same kind to a certain number of cases equally possible, that is to say, to such as we may be equally undecided about in regard to their existence, and in determining the number of cases favorable to the event whose probability is sought. The ratio of $y = mx + c + \epsilon$

point 1:
$$x = 1, y = 3$$

 $3 = m + c + \epsilon_1$
point 2: $x = 3, y = 1$
 $1 = 3m + c + \epsilon_2$
point 3: $x = 2, y = 2.5$
 $2.5 = 2m + c + \epsilon_3$

What about two unknowns and *one* observation?

$$y_1 = mx_1 + c$$





Can compute *m* given *c*.

 $c = 1.75 \Longrightarrow m = 1.25$



Can compute *m* given *c*.

$$c = -0.777 \Longrightarrow m = 3.78$$



Can compute *m* given *c*.

 $c = -4.01 \Longrightarrow m = 7.01$



Can compute *m* given *c*.

 $c = -0.718 \Longrightarrow m = 3.72$



Can compute *m* given *c*.

 $c = 2.45 \Longrightarrow m = 0.545$



Can compute *m* given *c*.

 $c = -0.657 \Longrightarrow m = 3.66$



Can compute *m* given *c*.

 $c = -3.13 \Longrightarrow m = 6.13$



Can compute *m* given *c*.

$$c = -1.47 \Longrightarrow m = 4.47$$



Can compute *m* given *c*. Assume

$$c \sim \mathcal{N}(0,4)$$
,

we find a distribution of solutions.



Gaussian Process

$$y_i(\mathbf{x}_i) = f(\mathbf{x}_i) + \epsilon_i$$

- Place a prior over the process as well as the noise.
- Leads to models that are not i.i.d.
- Contrast with classical model's objective function:

$$\sum_{i=1}^{n} (1 - y_i (\mathbf{w}^{\top} \mathbf{x}_i - b))_+ + \lambda \mathbf{w}^{\top} \mathbf{w}$$

- I'm keen on the idea of a conceptual separation model and algorithm.
- Model is how you encode the regularities of the universe.
- Algorithm is how you combine that model with data.

```
data + model \rightarrow prediction
```

 Of course often we are restricted in modeling choice due to lack of algorithms.









Multi-variate Gaussians

- We will consider a Gaussian with a particular structure of covariance matrix.
- Generate a single sample from this 25 dimensional Gaussian distribution, $\mathbf{f} = [f_1, f_2 \dots f_{25}]$.
- We will plot these points against their index.


(a) A 25 dimensional correlated random variable (values ploted against index)

(b) colormap *i*showing correlations between dimensions.



(a) A 25 dimensional correlated random variable (values ploted against index)

(b) colormap *i*showing correlations between dimensions.



(a) A 25 dimensional correlated random variable (values ploted against index)



0.8

(b) colormap showing correlations between dimensions.





(a) A 25 dimensional correlated random variable (values ploted against index)

(b) colormap showing correlations between dimensions.



(a) A 25 dimensional correlated random variable (values ploted against index)

(b) colormap showing correlations between dimensions.

0.8

0.6

0.4

0.2

0



0.8 0.6 0.4 0.2 0.2

(a) A 25 dimensional correlated random variable (values ploted against index)

(b) colormap showing correlations between dimensions.



(a) A 25 dimensional correlated random variable (values ploted against index)



(b) colormap showing correlations between dimensions.





The single contour of the Gaussian density represents the joint distribution, p(f₁, f₂).



- ► The single contour of the Gaussian density represents the joint distribution, p(f₁, f₂).
- We observe that $f_1 = -0.313$.



- ► The single contour of the Gaussian density represents the joint distribution, p(f₁, f₂).
- We observe that $f_1 = -0.313$.
- Conditional density: $p(f_2|f_1 = -0.313)$.



- The single contour of the Gaussian density represents the joint distribution, p(f₁, f₂).
- We observe that $f_1 = -0.313$.
- Conditional density: $p(f_2|f_1 = -0.313)$.

Prediction with Correlated Gaussians

- ▶ Prediction of *f*₂ from *f*₁ requires *conditional density*.
- Conditional density is *also* Gaussian.

$$p(f_2|f_1) = \mathcal{N}\left(f_2|\frac{k_{1,2}}{k_{1,1}}f_1, k_{2,2} - \frac{k_{1,2}^2}{k_{1,1}}\right)$$

where covariance of joint density is given by

$$\mathbf{K} = \begin{bmatrix} k_{1,1} & k_{1,2} \\ k_{2,1} & k_{2,2} \end{bmatrix}$$



The single contour of the Gaussian density represents the joint distribution, p(f₁, f₅).



- The single contour of the Gaussian density represents the joint distribution, p(f₁, f₅).
- We observe that $f_1 = -0.313$.



- The single contour of the Gaussian density represents the joint distribution, p(f₁, f₅).
- We observe that $f_1 = -0.313$.
- Conditional density: $p(f_5|f_1 = -0.313)$.



- The single contour of the Gaussian density represents the joint distribution, p(f₁, f₅).
- We observe that $f_1 = -0.313$.
- Conditional density: $p(f_5|f_1 = -0.313)$.

Prediction with Correlated Gaussians

- Prediction of f* from f requires multivariate *conditional density*.
- Multivariate conditional density is *also* Gaussian.

$$p(\mathbf{f}_{*}|\mathbf{f}) = \mathcal{N}\left(\mathbf{f}_{*}|\mathbf{K}_{*,f}\mathbf{K}_{f,f}^{-1}\mathbf{f}, \mathbf{K}_{*,*} - \mathbf{K}_{*,f}\mathbf{K}_{f,f}^{-1}\mathbf{K}_{f,*}\right)$$

Here covariance of joint density is given by

$$\mathbf{K} = \begin{bmatrix} \mathbf{K}_{f,f} & \mathbf{K}_{*,f} \\ \mathbf{K}_{f,*} & \mathbf{K}_{*,*} \end{bmatrix}$$

Prediction with Correlated Gaussians

- Prediction of f* from f requires multivariate *conditional density*.
- Multivariate conditional density is *also* Gaussian.

$$p(\mathbf{f}_*|\mathbf{f}) = \mathcal{N}(\mathbf{f}_*|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
$$\boldsymbol{\mu} = \mathbf{K}_{*,f}\mathbf{K}_{\mathbf{f},\mathbf{f}}^{-1}\mathbf{f}$$
$$\boldsymbol{\Sigma} = \mathbf{K}_{*,*} - \mathbf{K}_{*,f}\mathbf{K}_{\mathbf{f},\mathbf{f}}^{-1}\mathbf{K}_{\mathbf{f},*}$$
$$\blacktriangleright \text{ Here covariance of joint density is given by}$$

$$\mathbf{K} = \begin{bmatrix} \mathbf{K}_{\mathbf{f},\mathbf{f}} & \mathbf{K}_{*,\mathbf{f}} \\ \mathbf{K}_{\mathbf{f},*} & \mathbf{K}_{*,*} \end{bmatrix}$$

Where did this covariance matrix come from?

Exponentiated Quadratic Kernel Function (RBF, Squared Exponential, Gaussian)

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right)$$

- Covariance matrix is built using the *inputs* to the function x.
- For the example above it was based on Euclidean distance.
- The covariance function is also know as a kernel.



Where did this covariance matrix come from?

Exponentiated Quadratic Kernel Function (RBF, Squared Exponential, Gaussian)

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right)$$

- Covariance matrix is built using the *inputs* to the function x.
- For the example above it was based on Euclidean distance.
- The covariance function is also know as a kernel.

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = -3.0, x_1 = -3.0$$

$$k_{1,1} = 1.00 \times \exp\left(-\frac{(-3.0 - -3.0)^2}{2 \times 2.00^2}\right)$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = -3.0, x_1 = -3.0$$

$$k_{1,1} = 1.00 \times \exp\left(-\frac{(-3.0 - -3.0)^2}{2 \times 2.00^2}\right)$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_1 = -3.0$$

$$k_{2,1} = 1.00 \times \exp\left(-\frac{(1.20 - 3.0)^2}{2 \times 2.00^2}\right)$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{\|x_{i} - x_{j}\|^{2}}{2\ell^{2}}\right)$$

$$x_{2} = 1.20, x_{1} = -3.0$$

$$k_{2,1} = 1.00 \times \exp\left(-\frac{(1.20 - 3.0)^{2}}{2 \times 2.00^{2}}\right)$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_1 = -3.0$$

$$k_{2,1} = 1.00 \times \exp\left(-\frac{(1.20 - 3.0)^2}{2 \times 2.00^2}\right)$$

$$1.00 \quad 0.110$$

$$0.110$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_2 = 1.20$$

$$k_{2,2} = 1.00 \times \exp\left(-\frac{(1.20 - 1.20)^2}{2 \times 2.00^2}\right)$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_2 = 1.20$$

$$k_{2,2} = 1.00 \times \exp\left(-\frac{(1.20 - 1.20)^2}{2 \times 2.00^2}\right)$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_1 = -3.0$$

$$k_{3,1} = 1.00 \times \exp\left(-\frac{(1.40 - -3.0)^2}{2 \times 2.00^2}\right)$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{\|x_{i} - x_{j}\|^{2}}{2\ell^{2}}\right)$$

$$x_{3} = 1.40, x_{1} = -3.0$$

$$k_{3,1} = 1.00 \times \exp\left(-\frac{(1.40 - 3.0)^{2}}{2 \times 2.00^{2}}\right)$$

$$0.0889$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_1 = -3.0$$

$$k_{3,1} = 1.00 \times \exp\left(-\frac{(1.40 - 3.0)^2}{2 \times 2.00^2}\right)$$

$$0.0889$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_2 = 1.20$$

$$k_{3,2} = 1.00 \times \exp\left(-\frac{(1.40 - 1.20)^2}{2 \times 2.00^2}\right)$$

$$0.0889$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_2 = 1.20$$

$$k_{3,2} = 1.00 \times \exp\left(-\frac{(1.40 - 1.20)^2}{2 \times 2.00^2}\right)$$

$$0.0889$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{||x_i - x_j||^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_2 = 1.20$$

$$k_{3,2} = 1.00 \times \exp\left(-\frac{(1.40 - 1.20)^2}{2 \times 2.00^2}\right)$$

$$0.0889 \quad 0.995$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{\|x_{i} - x_{j}\|^{2}}{2\ell^{2}}\right)$$

$$x_{3} = 1.40, x_{3} = 1.40$$

$$k_{3,3} = 1.00 \times \exp\left(-\frac{(1.40 - 1.40)^{2}}{2 \times 2.00^{2}}\right)$$

$$0.0889 \quad 0.995$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_3 = 1.40$$

$$k_{3,3} = 1.00 \times \exp\left(-\frac{(1.40 - 1.40)^2}{2 \times 2.00^2}\right)$$

$$1.00$$

$$1.00 \quad 0.110 \quad 0.0889$$

$$0.995$$

$$1.00$$
Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$



Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = -3, x_1 = -3$$

$$k_{1,1} = 1.0 \times \exp\left(-\frac{(-3--3)^2}{2 \times 2.0^2}\right)$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{||x_{i} - x_{j}||^{2}}{2\ell^{2}}\right)$$

$$x_{1} = -3, x_{1} = -3$$

$$k_{1,1} = 1.0 \times \exp\left(-\frac{(-3 - 3)^{2}}{2 \times 2.0^{2}}\right)$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{||x_{i} - x_{j}||^{2}}{2\ell^{2}}\right)$$

$$x_{2} = 1.2, x_{1} = -3$$

$$k_{2,1} = 1.0 \times \exp\left(-\frac{(1.2 - 3)^{2}}{2 \times 2.0^{2}}\right)$$

Where did this covariance matrix come from?

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{||x_{i} - x_{j}||^{2}}{2\ell^{2}}\right)$$

$$x_{2} = 1.2, x_{1} = -3$$

$$k_{2,1} = 1.0 \times \exp\left(-\frac{(1.2 - 3)^{2}}{2 \times 2.0^{2}}\right)$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{||x_{i}-x_{j}||^{2}}{2\ell^{2}}\right)$$

$$x_{2} = 1.2, x_{2} = 1.2$$

$$k_{2,2} = 1.0 \times \exp\left(-\frac{(1.2-1.2)^{2}}{2\times 2.0^{2}}\right)$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{||x_{i} - x_{j}||^{2}}{2\ell^{2}}\right)$$

$$x_{2} = 1.2, x_{2} = 1.2$$

$$k_{2,2} = 1.0 \times \exp\left(-\frac{(1.2 - 1.2)^{2}}{2 \times 2.0^{2}}\right)$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{||x_{i} - x_{j}||^{2}}{2\ell^{2}}\right)$$

$$x_{3} = 1.4, x_{1} = -3$$

$$k_{3,1} = 1.0 \times \exp\left(-\frac{(1.4 - -3)^{2}}{2 \times 2.0^{2}}\right)$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.4, x_1 = -3$$

$$k_{3,1} = 1.0 \times \exp\left(-\frac{(1.4 - 3)^2}{2 \times 2.0^2}\right)$$

$$1.0 \quad 0.11$$

$$0.11 \quad 1.0$$

$$0.089$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{\|x_{i} - x_{j}\|^{2}}{2\ell^{2}}\right)$$

$$x_{3} = 1.4, x_{1} = -3$$

$$k_{3,1} = 1.0 \times \exp\left(-\frac{(1.4 - 3)^{2}}{2 \times 2.0^{2}}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 \\ 0.089 \end{bmatrix}$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{\|x_{i} - x_{j}\|^{2}}{2\ell^{2}}\right)$$

$$x_{3} = 1.4, x_{2} = 1.2$$

$$k_{3,2} = 1.0 \times \exp\left(-\frac{(1.4 - 1.2)^{2}}{2 \times 2.0^{2}}\right)$$

$$1.0 = 0.11 = 0.089$$

$$0.11 = 1.0$$

$$0.089$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{||x_i - x_j||^2}{2\ell^2}\right)$$

$$x_3 = 1.4, x_2 = 1.2$$

$$k_{3,2} = 1.0 \times \exp\left(-\frac{(1.4 - 1.2)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 \\ 0.089 & 1.0 \end{bmatrix}$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.4, x_2 = 1.2$$

$$k_{3,2} = 1.0 \times \exp\left(-\frac{(1.4 - 1.2)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 & 1.0 \\ 0.089 & 1.0 \end{bmatrix}$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.4, x_3 = 1.4$$

$$k_{3,3} = 1.0 \times \exp\left(-\frac{(1.4 - 1.4)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 & 1.0 \\ 0.089 & 1.0 \end{bmatrix}$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{\|x_{i} - x_{j}\|^{2}}{2\ell^{2}}\right)$$

$$x_{3} = 1.4, x_{3} = 1.4$$

$$k_{3,3} = 1.0 \times \exp\left(-\frac{(1.4 - 1.4)^{2}}{2 \times 2.0^{2}}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 & 1.0 \\ 0.089 & 1.0 & 1.0 \end{bmatrix}$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_1 = -3$$

$$k_{4,1} = 1.0 \times \exp\left(-\frac{(2.0 - 3)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 \\ 0.11 & 1.0 & 1.0 \\ 0.089 & 1.0 & 1.0 \end{bmatrix}$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_1 = -3$$

$$k_{4,1} = 1.0 \times \exp\left(-\frac{(2.0 - 3)^2}{2 \times 2.0^2}\right)$$

$$1.0 \quad 0.011 \quad 0.089$$

$$1.0 \quad 1.0$$

$$0.044$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_1 = -3$$

$$k_{4,1} = 1.0 \times \exp\left(-\frac{(2.0 - 3)^2}{2 \times 2.0^2}\right)$$

$$1.0 \quad 0.11 \quad 0.089 \quad 0.044$$

$$0.11 \quad 1.0 \quad 1.0$$

$$0.044$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_2 = 1.2$$

$$k_{4,2} = 1.0 \times \exp\left(-\frac{(2.0 - 1.2)^2}{2 \times 2.0^2}\right)$$

$$1.0 \quad 0.11 \quad 0.089 \quad 0.044$$

$$0.11 \quad 1.0 \quad 1.0$$

$$0.044$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_2 = 1.2$$

$$k_{4,2} = 1.0 \times \exp\left(-\frac{(2.0 - 1.2)^2}{2 \times 2.0^2}\right)$$

$$1.0 \quad 0.044$$

$$0.11 \quad 1.0 \quad 1.0$$

$$0.044 \quad 0.92$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_2 = 1.2$$

$$k_{4,2} = 1.0 \times \exp\left(-\frac{(2.0 - 1.2)^2}{2 \times 2.0^2}\right)$$

$$1.0 \quad 0.11 \quad 0.089 \quad 0.044$$

$$0.11 \quad 1.0 \quad 0.92$$

$$0.089 \quad 1.0 \quad 1.0$$

$$0.044 \quad 0.92$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_3 = 1.4$$

$$k_{4,3} = 1.0 \times \exp\left(-\frac{(2.0 - 1.4)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 & 0.044 \\ 0.11 & 1.0 & 1.0 & 0.92 \\ 0.089 & 1.0 & 1.0 \\ 0.044 & 0.92 \end{bmatrix}$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_3 = 1.4$$

$$k_{4,3} = 1.0 \times \exp\left(-\frac{(2.0 - 1.4)^2}{2 \times 2.0^2}\right)$$

$$1.0 \quad 0.11 \quad 0.089 \quad 0.044$$

$$0.11 \quad 1.0 \quad 0.92$$

$$0.089 \quad 1.0 \quad 1.0$$

$$0.044 \quad 0.92 \quad 0.96$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_3 = 1.4$$

$$k_{4,3} = 1.0 \times \exp\left(-\frac{(2.0 - 1.4)^2}{2 \times 2.0^2}\right)$$

$$\begin{bmatrix} 1.0 & 0.11 & 0.089 & 0.044 \\ 0.11 & 1.0 & 1.0 & 0.92 \\ 0.089 & 1.0 & 1.0 & 0.96 \\ 0.044 & 0.92 & 0.96 \end{bmatrix}$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_4 = 2.0$$

$$k_{4,4} = 1.0 \times \exp\left(-\frac{(2.0 - 2.0)^2}{2 \times 2.0^2}\right)$$

$$1.0 \quad 0.11 \quad 0.089 \quad 0.044$$

$$0.11 \quad 1.0 \quad 0.92$$

$$0.089 \quad 1.0 \quad 1.0 \quad 0.96$$

$$0.044 \quad 0.92 \quad 0.96$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_4 = 2.0, x_4 = 2.0$$

$$k_{4,4} = 1.0 \times \exp\left(-\frac{(2.0 - 2.0)^2}{2 \times 2.0^2}\right)$$

$$1.0 \quad 0.11 \quad 0.089 \quad 0.044$$

$$0.11 \quad 1.0 \quad 0.92$$

$$0.089 \quad 1.0 \quad 1.0 \quad 0.96$$

$$0.044 \quad 0.92 \quad 0.96 \quad 1.0$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$



Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_1 = -3.0, x_1 = -3.0$$

$$k_{1,1} = 4.00 \times \exp\left(-\frac{(-3.0 - -3.0)^2}{2 \times 5.00^2}\right)$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{\|x_{i} - x_{j}\|^{2}}{2\ell^{2}}\right)$$

$$x_{1} = -3.0, x_{1} = -3.0$$

$$k_{1,1} = 4.00 \times \exp\left(-\frac{(-3.0 - 3.0)^{2}}{2\times 5.00^{2}}\right)$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_1 = -3.0$$

$$k_{2,1} = 4.00 \times \exp\left(-\frac{(1.20 - 3.0)^2}{2 \times 5.00^2}\right)$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{\|x_{i} - x_{j}\|^{2}}{2\ell^{2}}\right)$$

$$x_{2} = 1.20, x_{1} = -3.0$$

$$k_{2,1} = 4.00 \times \exp\left(-\frac{(1.20 - 3.0)^{2}}{2\times 5.00^{2}}\right)$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{\|x_{i} - x_{j}\|^{2}}{2\ell^{2}}\right)$$

$$x_{2} = 1.20, x_{1} = -3.0$$

$$k_{2,1} = 4.00 \times \exp\left(-\frac{(1.20 - 3.0)^{2}}{2 \times 5.00^{2}}\right)$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_2 = 1.20$$

$$k_{2,2} = 4.00 \times \exp\left(-\frac{(1.20 - 1.20)^2}{2 \times 5.00^2}\right)$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_2 = 1.20, x_2 = 1.20$$

$$k_{2,2} = 4.00 \times \exp\left(-\frac{(1.20 - 1.20)^2}{2 \times 5.00^2}\right)$$

Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{\|x_{i} - x_{j}\|^{2}}{2\ell^{2}}\right)$$

$$x_{3} = 1.40, x_{1} = -3.0$$

$$k_{3,1} = 4.00 \times \exp\left(-\frac{(1.40 - -3.0)^{2}}{2\times 5.00^{2}}\right)$$
Where did this covariance matrix come from?

$$k(x_{i}, x_{j}) = \alpha \exp\left(-\frac{\|x_{i} - x_{j}\|^{2}}{2\ell^{2}}\right)$$

$$x_{3} = 1.40, x_{1} = -3.0$$

$$k_{3,1} = 4.00 \times \exp\left(-\frac{(1.40 - 3.0)^{2}}{2\times 5.00^{2}}\right)$$

$$2.72$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_1 = -3.0$$

$$k_{3,1} = 4.00 \times \exp\left(-\frac{(1.40 - 3.0)^2}{2 \times 5.00^2}\right)$$

$$2.72$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_2 = 1.20$$

$$k_{3,2} = 4.00 \times \exp\left(-\frac{(1.40 - 1.20)^2}{2 \times 5.00^2}\right)$$

$$2.72$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_2 = 1.20$$

$$k_{3,2} = 4.00 \times \exp\left(-\frac{(1.40 - 1.20)^2}{2 \times 5.00^2}\right)$$

$$2.72$$

$$4.00$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_2 = 1.20$$

$$k_{3,2} = 4.00 \times \exp\left(-\frac{(1.40 - 1.20)^2}{2 \times 5.00^2}\right)$$

$$2.72 \quad 4.00$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_3 = 1.40$$

$$k_{3,3} = 4.00 \times \exp\left(-\frac{(1.40 - 1.40)^2}{2 \times 5.00^2}\right)$$

$$2.72 \quad 4.00$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$

$$x_3 = 1.40, x_3 = 1.40$$

$$k_{3,3} = 4.00 \times \exp\left(-\frac{(1.40 - 1.40)^2}{2 \times 5.00^2}\right)$$

$$k_{3,3} = 4.00 \times \exp\left(-\frac{(1.40 - 1.40)^2}{2 \times 5.00^2}\right)$$

Where did this covariance matrix come from?

$$k(x_i, x_j) = \alpha \exp\left(-\frac{\|x_i - x_j\|^2}{2\ell^2}\right)$$



















Gaussian noise model,

$$p(y_i|f_i) = \mathcal{N}(y_i|f_i,\sigma^2)$$

where σ^2 is the variance of the noise.

• Equivalent to a covariance function of the form

$$k(\mathbf{x}_i, \mathbf{x}_j) = \delta_{i,j} \sigma^2$$

where $\delta_{i,j}$ is the Kronecker delta function.

 Additive nature of Gaussians means we can simply add this term to existing covariance matrices.



















Gaussian Process Fit to Olympic Marathon Data



Gaussian Processes

GP Non-Gaussian

GP Limitations

Kalman Filter

Dimensionality Reduction

General Noise Models

Graph of a GP

- Relates input variables,
 X, to vector, y, through f given kernel parameters
 θ.
- Plate notation indicates independence of y_i|f_i.
- ► In general p (y_i|f_i) is non-Gaussian.
- We approximate with Gaussian $p(y_i|f_i) \approx \mathcal{N}(m_i|f_i, \beta_i^{-1}).$



Figure: The Gaussian process depicted graphically.



Figure: Inclusion of a data point with Gaussian noise.



Figure: Inclusion of a data point with Gaussian noise.



Figure: Inclusion of a data point with Gaussian noise.

Local Moment Matching

- Easiest to consider a single previously unseen data point, y_{*}, x_{*}.
- ► Before seeing data point, prediction of *f*^{*} is a GP, *q*(*f*^{*}|**y**, **X**).
- Update prediction using Bayes' Rule,

$$p(f_*|\mathbf{y}, y_*, \mathbf{X}, \mathbf{x}_*) = \frac{p(y_*|f_*) p(f_*|\mathbf{y}, \mathbf{X}, \mathbf{x}_*)}{p(\mathbf{y}, y_*|\mathbf{X}, \mathbf{x}_*)}.$$

This posterior is not a Gaussian process if $p(y_*|f_*)$ is non-Gaussian.

Classification Noise Model

Probit Noise Model



Figure: The probit model (classification). The plot shows $p(y_i|f_i)$ for different values of y_i . For $y_i = 1$ we have $p(y_i|f_i) = \phi(f_i) = \int_{-\infty}^{f_i} \mathcal{N}(z|0, 1) dz$.

Match Moments

- Idea behind EP approximate with a Gaussian process at this stage by matching moments.
- This is equivalent to minimizing the following KL divergence where q (f_{*}|y, y_{*}, X, x_{*}) is constrained to be a GP.

 $q\left(f_{*}|\mathbf{y},y_{*}\mathbf{X},\mathbf{x}_{*}\right) = \operatorname{argmin}_{q\left(f_{*}|\mathbf{y},y_{*}\mathbf{X},\mathbf{x}_{*}\right)} \operatorname{KL}\left(p\left(f_{*}|\mathbf{y},y_{*}\mathbf{X},\mathbf{x}_{*}\right) \| q\left(f_{*}|\mathbf{y},y_{*},\mathbf{X},\mathbf{x}_{*}\right)\right)$

This is equivalent to setting

$$\langle f_* \rangle_{q(f_*|\mathbf{y}, y_*, \mathbf{X}, \mathbf{x}_*)} = \langle f_* \rangle_{p(f_*|\mathbf{y}, y_*, \mathbf{X}, \mathbf{x}_*)}$$
$$\langle f_*^2 \rangle_{q(f_*|\mathbf{y}, y_*, \mathbf{X}, \mathbf{x}_*)} = \langle f_*^2 \rangle_{p(f_*|\mathbf{y}, y_*, \mathbf{X}, \mathbf{x}_*)}$$

Equivalent Gaussian

► This is achieved by replacing p (y_{*}|f_{*}) with a Gaussian distribution

$$p(f_*|\mathbf{y}, y_*, \mathbf{X}, \mathbf{x}_*) = \frac{p(y_*|f_*)p(f_*|\mathbf{y}, \mathbf{X}, \mathbf{x}_*)}{p(\mathbf{y}, y_*|\mathbf{X}, \mathbf{x}_*)}$$

becomes

$$q\left(f_{*}|\mathbf{y}, y_{*}, \mathbf{X}, \mathbf{x}_{*}\right) = \frac{\mathcal{N}\left(m_{*}|f_{*}, \beta_{m}^{-1}\right) p\left(f_{*}|\mathbf{y}, \mathbf{X}, \mathbf{x}_{*}\right)}{p\left(\mathbf{y}, y_{*}|\mathbf{X}, \mathbf{x}_{*}\right)}.$$


Figure: An EP style update with a classification noise model.



Figure: An EP style update with a classification noise model.



Figure: An EP style update with a classification noise model.



Figure: An EP style update with a classification noise model.

Ordinal Noise Model

Ordered Categories



Figure: The ordered categorical noise model (ordinal regression). The plot shows $p(y_i|f_i)$ for different values of y_i . Here we have assumed three categories.

- Equivalent Gaussian is found by making a local 2nd order Taylor approximation at the mode.
- Laplace was the first to suggest this¹, so it's known as the Laplace approximation.

Can we determine covariance parameters from the data?

$$\mathcal{N}(\mathbf{y}|\mathbf{0},\mathbf{K}) = \frac{1}{(2\pi)^{\frac{n}{2}}|\mathbf{K}|^{\frac{1}{2}}} \exp\left(-\frac{\mathbf{y}^{\mathsf{T}}\mathbf{K}^{-1}\mathbf{y}}{2}\right)$$

The parameters are *inside* the covariance function (matrix).

Can we determine covariance parameters from the data?

$$\mathcal{N}(\mathbf{y}|\mathbf{0},\mathbf{K}) = \frac{1}{(2\pi)^{\frac{n}{2}}|\mathbf{K}|^{\frac{1}{2}}} \exp\left(-\frac{\mathbf{y}^{\mathsf{T}}\mathbf{K}^{-1}\mathbf{y}}{2}\right)$$

The parameters are *inside* the covariance function (matrix).

Can we determine covariance parameters from the data?

$$\log \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}) = -\frac{1}{2} \log |\mathbf{K}| - \frac{\mathbf{y}^{\mathsf{T}} \mathbf{K}^{-1} \mathbf{y}}{2} - \frac{n}{2} \log 2\pi$$

The parameters are *inside* the covariance function (matrix).

Can we determine covariance parameters from the data?

$$E(\boldsymbol{\theta}) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^{\top} \mathbf{K}^{-1} \mathbf{y}}{2}$$

The parameters are *inside* the covariance function (matrix).

Eigendecomposition of Covariance

A useful decomposition for understanding the objective function.

 $\mathbf{K} = \mathbf{R} \boldsymbol{\Lambda}^2 \mathbf{R}^\top$



Diagonal of Λ represents distance along axes. **R** gives a rotation of these axes.

1 4 4 14 1 4 1 5 1 5 5 5



















$$|\mathbf{\Lambda}| = \lambda_1 \lambda_2 \lambda_3$$





 $|\mathbf{R}\mathbf{\Lambda}| = \lambda_1 \lambda_2$





 y_1





 y_1





 y_1



$$E(\boldsymbol{\theta}) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^{\mathsf{T}} \mathbf{K}^{-1} \mathbf{y}}{2}$$



$$E(\boldsymbol{\theta}) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^{\mathsf{T}} \mathbf{K}^{-1} \mathbf{y}}{2}$$



$$E(\boldsymbol{\theta}) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^{\mathsf{T}} \mathbf{K}^{-1} \mathbf{y}}{2}$$



$$E(\boldsymbol{\theta}) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^{\mathsf{T}} \mathbf{K}^{-1} \mathbf{y}}{2}$$



$$E(\boldsymbol{\theta}) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^{\mathsf{T}} \mathbf{K}^{-1} \mathbf{y}}{2}$$



$$E(\boldsymbol{\theta}) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^{\mathsf{T}} \mathbf{K}^{-1} \mathbf{y}}{2}$$



$$E(\boldsymbol{\theta}) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^{\mathsf{T}} \mathbf{K}^{-1} \mathbf{y}}{2}$$



$$E(\boldsymbol{\theta}) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^{\mathsf{T}} \mathbf{K}^{-1} \mathbf{y}}{2}$$



$$E(\boldsymbol{\theta}) = \frac{1}{2} \log |\mathbf{K}| + \frac{\mathbf{y}^{\mathsf{T}} \mathbf{K}^{-1} \mathbf{y}}{2}$$

- Given given expression levels in the form of a time series from Della Gatta et al. (2008).
- Want to detect if a gene is expressed or not, fit a GP to each gene (Kalaitzis and Lawrence, 2011).



RESEARCH ARTICLE

Open Access

A Simple Approach to Ranking Differentially Expressed Gene Expression Time Courses through Gaussian Process Regression

Alfredo A Kalaitzis" and Neil D Lawrence"

Abstract

Background: The analysis of gene expression from time series underpins many biological studies. Two basic forms of analysis recur for data of this type: removing inactive (quiet) genes from the study and determining which genes are differentially expressed. Often these analysis stages are applied disregarding the fact that the data is drawn from a time series. In this paper we propose a simple model for accounting for the underlying temporal nature of the data based on a Gaussian process.

Results: We review Gaussian process (GP) regression for estimating the continuous trajectories underlying in gene expression time-series. We present a simple approach which can be used to filter quiet genes, or for the case of time series in the form of expression ratios, quantify differential expression. We assess via ROC curves the rankings produced by our regression framework and compare them to a recently proposed hierarchical Bayesian model for the analysis of gene expression time-series (BATS). We compare on both simulated and experimental data showing that the proposed approach considerably outperforms the current state of the art.



Contour plot of Gaussian process likelihood.


Optima: length scale of 1.2221 and \log_{10} SNR of 1.9654 log likelihood is -0.22317.



Optima: length scale of 1.5162 and \log_{10} SNR of 0.21306 log likelihood is -0.23604.



Optima: length scale of 2.9886 and \log_{10} SNR of -4.506 log likelihood is -2.1056.

Basis Function Form

Radial basis functions commonly have the form

$$\phi_k(\mathbf{x}_i) = \exp\left(-\frac{|\mathbf{x}_i - \boldsymbol{\mu}_k|^2}{2\ell^2}\right)$$



Figure: A set of radial basis functions with width $\ell = 2$ and location parameters $\mu = [-4 \ 0 \ 4]^{\top}$.

Represent a function by a linear sum over a basis,

$$f(\mathbf{x}_{i,:};\mathbf{w}) = \sum_{k=1}^{m} w_k \phi_k(\mathbf{x}_{i,:}), \qquad (1)$$

• Here: *m* basis functions and $\phi_k(\cdot)$ is *k*th basis function and

$$\mathbf{w} = [w_1, \ldots, w_m]^\top$$

• For standard linear model: $\phi_k(\mathbf{x}_{i,:}) = x_{i,k}$.

Random Functions

Functions derived using:

$$f(x) = \sum_{k=1}^m w_k \phi_k(x),$$



$$w_k \sim \mathcal{N}(0, \alpha)$$
.



х

Figure: Functions sampled using the basis set from figure 9. Each line is a separate sample, generated by a weighted sum of the basis set. The weights, **w** are sampled from a Gaussian density with variance $\alpha = 1$.

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

$$y \sim \mathcal{N}(\mu, \sigma^2)$$

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

$$y \sim \mathcal{N}\left(\mu,\sigma^2\right)$$

$$wy \sim \mathcal{N}\left(w\mu, w^2\sigma^2\right)$$









► If

 $\mathbf{y} = \mathbf{W}\mathbf{x}$



Covariance Functions

RBF Basis Functions

$$k(\mathbf{x}, \mathbf{x}') = \alpha \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

$$\phi_k(x) = \exp\left(-\frac{\left\|x - \mu_k\right\|_2^2}{\ell^2}\right)$$
$$\mu = \begin{bmatrix}-1\\0\\1\end{bmatrix}$$



Covariance Functions

RBF Basis Functions

$$k(\mathbf{x}, \mathbf{x}') = \alpha \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$





1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

$$y \sim \mathcal{N}(\mu, \sigma^2)$$

1. Sum of Gaussian variables is also Gaussian.

$$y_i \sim \mathcal{N}\left(\mu_i, \sigma_i^2\right)$$

$$\sum_{i=1}^{n} y_i \sim \mathcal{N}\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right)$$

$$y \sim \mathcal{N}\left(\mu,\sigma^2\right)$$

$$wy \sim \mathcal{N}\left(w\mu, w^2\sigma^2\right)$$









► If

 $\mathbf{y} = \mathbf{W}\mathbf{x}$



Need to choose

- 1. location of centers
- 2. number of basis functions

Restrict analysis to 1-D input, *x*.

Consider uniform spacing over a region:

$$k(x_i, x_j) = \alpha \phi_k(x_i)^\top \phi_k(x_j)$$

Need to choose

- 1. location of centers
- 2. number of basis functions

Restrict analysis to 1-D input, *x*.

• Consider uniform spacing over a region:

$$k(x_i, x_j) = \alpha \sum_{k=1}^m \phi_k(x_i)\phi_k(x_j)$$

Need to choose

- 1. location of centers
- 2. number of basis functions

Restrict analysis to 1-D input, *x*.

Consider uniform spacing over a region:

$$k(x_i, x_j) = \alpha \sum_{k=1}^{m} \exp\left(-\frac{(x_i - \mu_k)^2}{2\ell^2}\right) \exp\left(-\frac{(x_j - \mu_k)^2}{2\ell^2}\right)$$

Need to choose

- 1. location of centers
- 2. number of basis functions

Restrict analysis to 1-D input, *x*.

• Consider uniform spacing over a region:

$$k(x_i, x_j) = \alpha \sum_{k=1}^{m} \exp\left(-\frac{(x_i - \mu_k)^2}{2\ell^2} - \frac{(x_j - \mu_k)^2}{2\ell^2}\right)$$

Need to choose

- 1. location of centers
- 2. number of basis functions

Restrict analysis to 1-D input, *x*.

Consider uniform spacing over a region:

$$k(x_i, x_j) = \alpha \sum_{k=1}^{m} \exp\left(-\frac{x_i^2 + x_j^2 - 2\mu_k(x_i + x_j) + 2\mu_k^2}{2\ell^2}\right),$$

Uniform Basis Functions

Set each center location to

$$\mu_k = a + \Delta \mu \cdot (k - 1).$$

Uniform Basis Functions

Set each center location to

$$\mu_k = a + \Delta \mu \cdot (k-1).$$

Specify the basis functions in terms of their indices,

$$k(x_{i}, x_{j}) = \alpha' \Delta \mu \sum_{k=1}^{m} \exp\left(-\frac{x_{i}^{2} + x_{j}^{2}}{2\ell^{2}} - \frac{2(a + \Delta \mu \cdot (k - 1))(x_{i} + x_{j}) + 2(a + \Delta \mu \cdot (k - 1))^{2}}{2\ell^{2}}\right)$$

Uniform Basis Functions

Set each center location to

$$\mu_k = a + \Delta \mu \cdot (k-1).$$

Specify the basis functions in terms of their indices,

$$k(x_{i}, x_{j}) = \alpha' \Delta \mu \sum_{k=1}^{m} \exp\left(-\frac{x_{i}^{2} + x_{j}^{2}}{2\ell^{2}} - \frac{2(a + \Delta \mu \cdot (k - 1))(x_{i} + x_{j}) + 2(a + \Delta \mu \cdot (k - 1))^{2}}{2\ell^{2}}\right).$$

• Here we've scaled variance of process by $\Delta \mu$.

Take

$$\mu_1 = a$$
 and $\mu_m = b$ so $b = a + \Delta \mu \cdot (m - 1)$

Take

$$\mu_1 = a$$
 and $\mu_m = b$ so $b = a + \Delta \mu \cdot (m - 1)$

This implies

$$b - a = \Delta \mu (m - 1)$$

Take

$$\mu_1 = a$$
 and $\mu_m = b$ so $b = a + \Delta \mu \cdot (m - 1)$

This implies

$$b-a=\Delta\mu(m-1)$$

and therefore

$$m = \frac{b-a}{\Delta\mu} + 1$$

Take

$$\mu_1 = a$$
 and $\mu_m = b$ so $b = a + \Delta \mu \cdot (m - 1)$

This implies

$$b-a=\Delta\mu(m-1)$$

and therefore

$$m = \frac{b-a}{\Delta \mu} + 1$$

• Take limit as $\Delta \mu \rightarrow 0$ so $m \rightarrow \infty$
Infinite Basis Functions

Take

$$\mu_1 = a$$
 and $\mu_m = b$ so $b = a + \Delta \mu \cdot (m - 1)$

This implies

$$b-a=\Delta\mu(m-1)$$

and therefore

$$m = \frac{b-a}{\Delta\mu} + 1$$

• Take limit as $\Delta \mu \rightarrow 0$ so $m \rightarrow \infty$

$$k(x_i, x_j) = \alpha' \int_a^b \exp\left(-\frac{x_i^2 + x_j^2}{2\ell^2} + \frac{2\left(\mu - \frac{1}{2}\left(x_i + x_j\right)\right)^2 - \frac{1}{2}\left(x_i + x_j\right)^2}{2\ell^2}\right) d\mu,$$

where we have used $a + k \cdot \Delta \mu \rightarrow \mu$.

Result

Performing the integration leads to

$$k(x_i, x_j) = \alpha' \sqrt{\pi \ell^2} \exp\left(-\frac{\left(x_i - x_j\right)^2}{4\ell^2}\right)$$
$$\times \frac{1}{2} \left[\operatorname{erf}\left(\frac{\left(b - \frac{1}{2}\left(x_i + x_j\right)\right)}{\ell}\right) - \operatorname{erf}\left(\frac{\left(a - \frac{1}{2}\left(x_i + x_j\right)\right)}{\ell}\right) \right],$$

Result

Performing the integration leads to

$$k(x_i, x_j) = \alpha' \sqrt{\pi \ell^2} \exp\left(-\frac{\left(x_i - x_j\right)^2}{4\ell^2}\right)$$
$$\times \frac{1}{2} \left[\operatorname{erf}\left(\frac{\left(b - \frac{1}{2}\left(x_i + x_j\right)\right)}{\ell}\right) - \operatorname{erf}\left(\frac{\left(a - \frac{1}{2}\left(x_i + x_j\right)\right)}{\ell}\right) \right],$$

• Now take limit as $a \to -\infty$ and $b \to \infty$

Result

Performing the integration leads to

$$k(x_i, x_j) = \alpha' \sqrt{\pi \ell^2} \exp\left(-\frac{\left(x_i - x_j\right)^2}{4\ell^2}\right)$$
$$\times \frac{1}{2} \left[\operatorname{erf}\left(\frac{\left(b - \frac{1}{2}\left(x_i + x_j\right)\right)}{\ell}\right) - \operatorname{erf}\left(\frac{\left(a - \frac{1}{2}\left(x_i + x_j\right)\right)}{\ell}\right) \right],$$

• Now take limit as $a \to -\infty$ and $b \to \infty$

$$k(x_i, x_j) = \alpha \exp\left(-\frac{(x_i - x_j)^2}{4\ell^2}\right).$$

where $\alpha = \alpha' \sqrt{\pi \ell^2}$.

 An RBF model with infinite basis functions is a Gaussian process.

- An RBF model with infinite basis functions is a Gaussian process.
- The covariance function is given by the exponentiated quadratic covariance function.

$$k(x_i, x_j) = \alpha \exp\left(-\frac{(x_i - x_j)^2}{4\ell^2}\right).$$

Infinite Feature Space

- An RBF model with infinite basis functions is a Gaussian process.
- The covariance function is the exponentiated quadratic.
- Note: The functional form for the covariance function and basis functions are similar.
 - this is a special case,
 - in general they are very different

Similar results can obtained for multi-dimensional input models Williams (1998); Neal (1996).

RBF Basis Functions

$$k(\mathbf{x}, \mathbf{x}') = \alpha \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

$$\phi_k(x) = \exp\left(-\frac{\left\|x - \mu_k\right\|_2^2}{\ell^2}\right)$$
$$\mu = \begin{bmatrix}-1\\0\\1\end{bmatrix}$$



RBF Basis Functions

$$k(\mathbf{x}, \mathbf{x}') = \alpha \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$





Where did this covariance matrix come from?

Exponentiated Quadratic Kernel Function (RBF, Squared Exponential, Gaussian)

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right)$$

- Covariance matrix is built using the *inputs* to the function x.
- For the example above it was based on Euclidean distance.
- The covariance function is also know as a kernel.



Where did this covariance matrix come from?

Exponentiated Quadratic Kernel Function (RBF, Squared Exponential, Gaussian)

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right)$$

- Covariance matrix is built using the *inputs* to the function x.
- For the example above it was based on Euclidean distance.
- The covariance function is also know as a kernel.

MLP Covariance Function

$$k(\mathbf{x}, \mathbf{x}') = \alpha \operatorname{asin}\left(\frac{w\mathbf{x}^{\top}\mathbf{x}' + b}{\sqrt{w\mathbf{x}^{\top}\mathbf{x} + b + 1}\sqrt{w\mathbf{x}'^{\top}\mathbf{x}' + b + 1}}\right)$$

 Based on infinite neural network model.

$$w = 40$$
$$b = 4$$



MLP Covariance Function

$$k(\mathbf{x}, \mathbf{x}') = \alpha \operatorname{asin}\left(\frac{w\mathbf{x}^{\top}\mathbf{x}' + b}{\sqrt{w\mathbf{x}^{\top}\mathbf{x} + b + 1}\sqrt{w\mathbf{x}'^{\top}\mathbf{x}' + b + 1}}\right)$$

 Based on infinite neural network model.

$$w = 40$$
$$b = 4$$

Constructing Covariance Functions

Sum of two covariances is also a covariance function.

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$$

Constructing Covariance Functions

Product of two covariances is also a covariance function.

 $k(\mathbf{x},\mathbf{x}')=k_1(\mathbf{x},\mathbf{x}')k_2(\mathbf{x},\mathbf{x}')$

Multiply by Deterministic Function

- If $f(\mathbf{x})$ is a Gaussian process.
- $g(\mathbf{x})$ is a deterministic function.
- $h(\mathbf{x}) = f(\mathbf{x})g(\mathbf{x})$
- Then

$$k_h(\mathbf{x}, \mathbf{x}') = g(\mathbf{x})k_f(\mathbf{x}, \mathbf{x}')g(\mathbf{x}')$$

where k_h is covariance for $h(\cdot)$ and k_f is covariance for $f(\cdot)$.

MLP Covariance Function

$$k(\mathbf{x}, \mathbf{x}') = \alpha \operatorname{asin}\left(\frac{w\mathbf{x}^{\top}\mathbf{x}' + b}{\sqrt{w\mathbf{x}^{\top}\mathbf{x} + b + 1}\sqrt{w\mathbf{x}'^{\top}\mathbf{x}' + b + 1}}\right)$$

 Based on infinite neural network model.

$$w = 40$$
$$b = 4$$



MLP Covariance Function

$$k(\mathbf{x}, \mathbf{x}') = \alpha \operatorname{asin}\left(\frac{w\mathbf{x}^{\top}\mathbf{x}' + b}{\sqrt{w\mathbf{x}^{\top}\mathbf{x} + b + 1}\sqrt{w\mathbf{x}'^{\top}\mathbf{x}' + b + 1}}\right)$$

 Based on infinite neural network model.

$$w = 40$$
$$b = 4$$

Linear Covariance Function

$$k(\mathbf{x}, \mathbf{x}') = \alpha \mathbf{x}^\top \mathbf{x}'$$



$$\alpha = 1$$



Linear Covariance Function

$$k(\mathbf{x}, \mathbf{x}') = \alpha \mathbf{x}^\top \mathbf{x}'$$

Bayesian linear regression.

$$\alpha = 1$$

Outline

Gaussian Processes

GP Non-Gaussian

GP Limitations

Kalman Filter

Dimensionality Reduction

- ► Inference is O(n³) due to matrix inverse (in practice use Cholesky).
- Gaussian processes don't deal well with discontinuities (financial crises, phosphorylation, collisions, edges in images).
- Widely used exponentiated quadratic covariance (RBF) can be too smooth in practice (but there are many alternatives!!).

Outline

Gaussian Processes

GP Non-Gaussian

GP Limitations

Kalman Filter

Dimensionality Reduction

Simple Markov Chain

- Assume 1-d latent state, a vector over time, $\mathbf{x} = [x_1 \dots x_T]$.
- Markov property,

$$\begin{aligned} x_i &= x_{i-1} + \epsilon_i, \\ \epsilon_i &\sim \mathcal{N}(0, \alpha) \\ \implies x_i &\sim \mathcal{N}(x_{i-1}, \alpha) \end{aligned}$$

Initial state,

 $x_0 \sim \mathcal{N}(0, \alpha_0)$

- If $x_0 \sim \mathcal{N}(0, \alpha)$ we have a Markov chain for the latent states.
- Markov chain it is specified by an initial distribution (Gaussian) and a transition distribution (Gaussian).



















Multivariate Gaussian Properties: Reminder

If $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ and $\mathbf{x} = \mathbf{W}\mathbf{z} + \mathbf{b}$ then $\mathbf{x} \sim \mathcal{N}(\mathbf{W}\boldsymbol{\mu} + \mathbf{b}, \mathbf{W}\mathbf{C}\mathbf{W}^{\mathsf{T}})$

Multivariate Gaussian Properties: Reminder



Matrix Representation of Latent Variables



 $x_1 = \epsilon_1$


 $x_2 = \epsilon_1 + \epsilon_2$



 $x_3 = \epsilon_1 + \epsilon_2 + \epsilon_3$



 $x_4 = \epsilon_1 + \epsilon_2 + \epsilon_3 + \epsilon_4$



 $x_5 = \epsilon_1 + \epsilon_2 + \epsilon_3 + \epsilon_4 + \epsilon_5$

 $\mathbf{x} = \mathbf{L}_1 \times \boldsymbol{\epsilon}$

- Since x is linearly related to ε we know x is a Gaussian process.
- Trick: we only need to compute the mean and covariance of x to determine that Gaussian.

$\mathbf{x} = \mathbf{L}_1 \boldsymbol{\epsilon}$

$\langle x angle = \langle L_1 \epsilon angle$

$\langle x \rangle = L_1 \langle \epsilon \rangle$

$\langle x \rangle = L_1 \langle \epsilon \rangle$

$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$

$\langle x\rangle = L_1 0$

$\langle x \rangle = 0$

$\mathbf{x}\mathbf{x}^{\top} = \mathbf{L}_{1}\boldsymbol{\epsilon}\boldsymbol{\epsilon}^{\top}\mathbf{L}_{1}^{\top}$ $\mathbf{x}^{\top} = \boldsymbol{\epsilon}^{\top}\mathbf{L}^{\top}$

 $\left\langle \mathbf{x}\mathbf{x}^{\top}\right\rangle =\left\langle \mathbf{L}_{1}\boldsymbol{\epsilon}\boldsymbol{\epsilon}^{\top}\mathbf{L}_{1}^{\top}\right\rangle$

$\langle \mathbf{x}\mathbf{x}^{\top} \rangle = \mathbf{L}_1 \langle \epsilon \epsilon^{\top} \rangle \mathbf{L}_1^{\top}$

$\langle \mathbf{x}\mathbf{x}^{\top} \rangle = \mathbf{L}_1 \langle \epsilon \epsilon^{\top} \rangle \mathbf{L}_1^{\top}$

 $\boldsymbol{\epsilon} \sim \mathcal{N}\left(\mathbf{0}, \alpha \mathbf{I}\right)$

$\langle \mathbf{x}\mathbf{x}^{\top} \rangle = \alpha \mathbf{L}_{1}\mathbf{L}_{1}^{\top}$

$\mathbf{x} = \mathbf{L}_1 \boldsymbol{\epsilon}$

$\mathbf{x} = \mathbf{L}_{1} \boldsymbol{\epsilon}$ $\boldsymbol{\epsilon} \sim \mathcal{N} \left(\mathbf{0}, \alpha \mathbf{I} \right)$

$\mathbf{x} = \mathbf{L}_{1} \boldsymbol{\epsilon}$ $\boldsymbol{\epsilon} \sim \mathcal{N} \left(\mathbf{0}, \alpha \mathbf{I} \right)$

 $\mathbf{x} = \mathbf{L}_{1} \boldsymbol{\epsilon}$ $\boldsymbol{\epsilon} \sim \mathcal{N} \left(\mathbf{0}, \alpha \mathbf{I} \right)$ \Longrightarrow $\mathbf{x} \sim \mathcal{N} \left(\mathbf{0}, \alpha \mathbf{L}_{1} \mathbf{L}_{1}^{\top} \right)$

- Make the variance dependent on time interval.
- Assume variance grows *linearly* with time.
- Justification: sum of two Gaussian distributed random variables is distributed as Gaussian with sum of variances.
- If variable's movement is additive over time (as described) variance scales linearly with time.

• Given $\epsilon \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I}) \Longrightarrow \epsilon \sim \mathcal{N}\left(\mathbf{0}, \alpha \mathbf{L}_{1}\mathbf{L}_{1}^{\top}\right).$ Then $\epsilon \sim \mathcal{N}\left(\mathbf{0}, \Delta t \alpha \mathbf{I}\right) \Longrightarrow \epsilon \sim \mathcal{N}\left(\mathbf{0}, \Delta t \alpha \mathbf{L}_{1}\mathbf{L}_{1}^{\top}\right).$

where Δt is the time interval between observations.

$$\boldsymbol{\epsilon} \sim \mathcal{N}\left(0, \alpha \Delta t \mathbf{I}\right), \quad \mathbf{x} \sim \mathcal{N}\left(0, \alpha \Delta t \mathbf{L}_{1} \mathbf{L}_{1}^{\top}\right)$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}\left(0, \alpha \Delta t \mathbf{I}\right), \quad \mathbf{x} \sim \mathcal{N}\left(0, \alpha \Delta t \mathbf{L}_{1} \mathbf{L}_{1}^{\top}\right)$$

 $\mathbf{K} = \alpha \Delta t \mathbf{L}_{\mathbf{1}} \mathbf{L}_{\mathbf{1}}^{\top}$

$$\boldsymbol{\epsilon} \sim \mathcal{N}\left(0, \alpha \Delta t \mathbf{I}\right), \quad \mathbf{x} \sim \mathcal{N}\left(0, \alpha \Delta t \mathbf{L}_{1} \mathbf{L}_{1}^{\top}\right)$$

$$\mathbf{K} = \alpha \Delta t \mathbf{L}_{\mathbf{1}} \mathbf{L}_{\mathbf{1}}^{\mathsf{T}}$$

$$k_{i,j} = \alpha \Delta t \mathbf{l}_{:,i}^{\top} \mathbf{l}_{:,j}$$

where $\mathbf{l}_{:,k}$ is a vector from the *k*th row of \mathbf{L}_1 : the first *k* elements are one, the next T - k are zero.

$$\boldsymbol{\epsilon} \sim \mathcal{N}\left(0, \alpha \Delta t \mathbf{I}\right), \quad \mathbf{x} \sim \mathcal{N}\left(0, \alpha \Delta t \mathbf{L}_{1} \mathbf{L}_{1}^{\top}\right)$$

$$\mathbf{K} = \alpha \Delta t \mathbf{L}_{\mathbf{1}} \mathbf{L}_{\mathbf{1}}^{\mathsf{T}}$$

$$k_{i,j} = \alpha \Delta t \mathbf{l}_{:,i}^{\top} \mathbf{l}_{:,j}$$

where $\mathbf{l}_{:,k}$ is a vector from the *k*th row of \mathbf{L}_1 : the first *k* elements are one, the next T - k are zero.

 $k_{i,j} = \alpha \Delta t \min(i, j)$ define $\Delta ti = t_i$ so $k_{i,j} = \alpha \min(t_i, t_j) = k(t_i, t_j)$

Where did this covariance matrix come from?

Markov Process

$$k(t,t') = \alpha \min(t,t')$$

 Covariance matrix is built using the *inputs* to the function *t*.



Where did this covariance matrix come from?

Markov Process

$$k(t,t') = \alpha \min(t,t')$$

 Covariance matrix is built using the *inputs* to the function *t*.



Where did this covariance matrix come from?

Markov Process

Visualization of inverse covariance (precision).

- Precision matrix is sparse: only neighbours in matrix are non-zero.
- This reflects *conditional* independencies in data.
- In this case *Markov* structure.



Where did this covariance matrix come from?

Exponentiated Quadratic Kernel Function (RBF, Squared Exponential, Gaussian)

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right)$$

- Covariance matrix is built using the *inputs* to the function x.
- For the example above it was based on Euclidean distance.
- The covariance function is also know as a kernel.



Where did this covariance matrix come from?

Exponentiated Quadratic Kernel Function (RBF, Squared Exponential, Gaussian)

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right)$$

- Covariance matrix is built using the *inputs* to the function x.
- For the example above it was based on Euclidean distance.
- The covariance function is also know as a kernel.

Where did this covariance matrix come from?

Exponentiated Quadratic

Visualization of inverse covariance (precision).

- Precision matrix is not sparse.
- Each point is dependent on all the others.
- In this case non-Markovian.



Where did this covariance matrix come from?

Markov Process

Visualization of inverse covariance (precision).

- Precision matrix is sparse: only neighbours in matrix are non-zero.
- This reflects *conditional* independencies in data.
- In this case *Markov* structure.



Simple Kalman Filter I

• We have state vector $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_q] \in \mathbb{R}^{T \times q}$ and if each state evolves independently we have

$$p(\mathbf{X}) = \prod_{i=1}^{q} p(\mathbf{x}_{:,i})$$
$$p(\mathbf{x}_{:,i}) = \mathcal{N}(\mathbf{x}_{:,i}|\mathbf{0}, \mathbf{K}).$$

• We want to obtain outputs through:

$$\mathbf{y}_{i,:} = \mathbf{W}\mathbf{x}_{i,:}$$

Stacking and Kronecker Products I

Represent with a 'stacked' system:

$$p(\mathbf{x}) = \mathcal{N}\left(\mathbf{x}|\mathbf{0}, \mathbf{I} \otimes \mathbf{K}\right)$$

where the stacking is placing each column of **X** one on top of another as

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_{:,1} \\ \mathbf{x}_{:,2} \\ \vdots \\ \mathbf{x}_{:,q} \end{bmatrix}$$
Kronecker Product



Kronecker Product



Stacking and Kronecker Products I

Represent with a 'stacked' system:

$$p(\mathbf{x}) = \mathcal{N}\left(\mathbf{x}|\mathbf{0}, \mathbf{I} \otimes \mathbf{K}\right)$$

where the stacking is placing each column of **X** one on top of another as

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_{:,1} \\ \mathbf{x}_{:,2} \\ \vdots \\ \mathbf{x}_{:,q} \end{bmatrix}$$

Column Stacking















Can also stack each row of **X** to form column vector:

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_{1,:} \\ \mathbf{x}_{2,:} \\ \vdots \\ \mathbf{x}_{T,:} \end{bmatrix}$$

 $p(\mathbf{x}) = \mathcal{N}\left(\mathbf{x}|\mathbf{0}, \mathbf{K} \otimes \mathbf{I}\right)$

Row Stacking













The observations are related to the latent points by a linear mapping matrix,

$$\mathbf{y}_{i,:} = \mathbf{W}\mathbf{x}_{i,:} + \boldsymbol{\epsilon}_{i,:}$$
$$\boldsymbol{\epsilon} \sim \mathcal{N}\left(0, \sigma^2 \mathbf{I}\right)$$

Mapping from Latent Process to Observed



This leads to a covariance of the form

 $(\mathbf{I} \otimes \mathbf{W})(\mathbf{K} \otimes \mathbf{I})(\mathbf{I} \otimes \mathbf{W}^{\top}) + \mathbf{I}\sigma^{2}$ Using $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{A}\mathbf{C} \otimes \mathbf{B}\mathbf{D}$ This leads to $\mathbf{K} \otimes \mathbf{W}\mathbf{W}^{\top} + \mathbf{I}\sigma^{2}$

or

$$\mathbf{y} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{W}\mathbf{W}^\top \otimes \mathbf{K} + \mathbf{I}\sigma^2\right)$$

Kernels for Vector Valued Outputs: A Review

Foundations and Trends[®] in Machine Learning Vol. 4, No. 3 (2011) 195–266 © 2012 M. A. Álvarez, L. Rosasco and N. D. Lawrence DOI: 10.1561/2200000036



Kernels for Vector-Valued Functions: A Review

By Mauricio A. Álvarez, Lorenzo Rosasco and Neil D. Lawrence This Kronecker structure leads to several published models.

$$(\mathbf{K}(\mathbf{x},\mathbf{x}'))_{j,j'}=k(\mathbf{x},\mathbf{x}')k_T(j,j'),$$

where *k* has **x** and k_T has *i* as inputs.

- Can think of multiple output covariance functions as covariances with augmented input.
- Alongside x we also input the *j* associated with the *output* of interest.

► Taking B = WW^T we have a matrix expression across outputs.

$$\mathbf{K}(\mathbf{x},\mathbf{x}')=k(\mathbf{x},\mathbf{x}')\mathbf{B},$$

where **B** is a $p \times p$ symmetric and positive semi-definite matrix.

- **B** is called the *coregionalization* matrix.
- We call this class of covariance functions *separable* due to their product structure.

Sum of Separable Covariance Functions

 In the same spirit a more general class of kernels is given by

$$\mathbf{K}(\mathbf{x},\mathbf{x}')=\sum_{j=1}^{q}k_{j}(\mathbf{x},\mathbf{x}')\mathbf{B}_{j}.$$

This can also be written as

$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \sum_{j=1}^{q} \mathbf{B}_{j} \otimes k_{j}(\mathbf{X},\mathbf{X}),$$

- This is like several Kalman filter-type models added together, but each one with a different set of latent functions.
- We call this class of kernels sum of separable kernels (SoS kernels).

- Use of GPs in Geostatistics is called kriging.
- These multi-output GPs pioneered in geostatistics: prediction over vector-valued output data is known as *cokriging*.
- The model in geostatistics is known as the *linear model of coregionalization* (LMC, Journel and Huijbregts (1978); Goovaerts (1997)).
- Most machine learning multitask models can be placed in the context of the LMC model.

Weighted sum of Latent Functions

- In the linear model of coregionalization (LMC) outputs are expressed as linear combinations of independent random functions.
- In the LMC, each component f_i is expressed as a linear sum

$$f_j(\mathbf{x}) = \sum_{j=1}^q w_{j,j} u_j(\mathbf{x}).$$

where the latent functions are independent and have covariance functions $k_i(\mathbf{x}, \mathbf{x}')$.

► The processes $\{f_j(\mathbf{x})\}_{j=1}^q$ are independent for $q \neq j'$.

Kalman Filter Special Case

- The Kalman filter is an example of the LMC where $u_i(\mathbf{x}) \rightarrow x_i(t)$.
- I.e. we've moved form time input to a more general input space.
- In matrix notation:
 - 1. Kalman filter

 $\mathbf{F} = \mathbf{W}\mathbf{X}$

2. LMC

 $\mathbf{F} = \mathbf{W}\mathbf{U}$

where the rows of these matrices **F**, **X**, **U** each contain *q* samples from their corresponding functions at a different time (Kalman filter) or spatial location (LMC).

- If one covariance used for latent functions (like in Kalman filter).
- This is called the intrinsic coregionalization model (ICM, Goovaerts (1997)).
- The kernel matrix corresponding to a dataset **X** takes the form

- ► If outputs are noise-free, maximum likelihood is equivalent to independent fits of **B** and *k*(**x**, **x**') (Helterbrand and Cressie, 1994).
- In geostatistics this is known as autokrigeability (Wackernagel, 2003).
- In multitask learning its the cancellation of intertask transfer (Bonilla et al., 2008).

$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{w}\mathbf{w}^{\top} \otimes k(\mathbf{X},\mathbf{X}).$$

$$\mathbf{w} = \begin{bmatrix} 1\\5 \end{bmatrix}$$
$$\mathbf{B} = \begin{bmatrix} 1 & 5\\5 & 25 \end{bmatrix}$$



$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{w}\mathbf{w}^{\top} \otimes k(\mathbf{X},\mathbf{X}).$$





$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{w}\mathbf{w}^{\top} \otimes k(\mathbf{X},\mathbf{X}).$$





$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{w}\mathbf{w}^{\top} \otimes k(\mathbf{X},\mathbf{X}).$$





$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{w}\mathbf{w}^{\top} \otimes k(\mathbf{X},\mathbf{X}).$$





$$\mathbf{B} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1.5 \end{bmatrix}$$



$$\mathbf{B} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1.5 \end{bmatrix}$$



$$\mathbf{B} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1.5 \end{bmatrix}$$


Intrinsic Coregionalization Model

 $\mathbf{K}(\mathbf{X},\mathbf{X})=\mathbf{B}\otimes k(\mathbf{X},\mathbf{X}).$

$$\mathbf{B} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1.5 \end{bmatrix}$$



Intrinsic Coregionalization Model

 $\mathbf{K}(\mathbf{X},\mathbf{X})=\mathbf{B}\otimes k(\mathbf{X},\mathbf{X}).$

$$\mathbf{B} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1.5 \end{bmatrix}$$



$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{B}_1 \otimes k_1(\mathbf{X},\mathbf{X}) + \mathbf{B}_2 \otimes k_2(\mathbf{X},\mathbf{X})$$

$$\mathbf{B}_{1} = \begin{bmatrix} 1.4 & 0.5 \\ 0.5 & 1.2 \end{bmatrix}$$
$$\ell_{1} = 1$$
$$\mathbf{B}_{2} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1.3 \end{bmatrix}$$
$$\ell_{2} = 0.2$$



$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{B}_1 \otimes k_1(\mathbf{X},\mathbf{X}) + \mathbf{B}_2 \otimes k_2(\mathbf{X},\mathbf{X})$$

$$\mathbf{B}_{1} = \begin{bmatrix} 1.4 & 0.5 \\ 0.5 & 1.2 \end{bmatrix}$$
$$\ell_{1} = 1$$
$$\mathbf{B}_{2} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1.3 \end{bmatrix}$$
$$\ell_{2} = 0.2$$



$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{B}_1 \otimes k_1(\mathbf{X},\mathbf{X}) + \mathbf{B}_2 \otimes k_2(\mathbf{X},\mathbf{X})$$

$$\mathbf{B}_{1} = \begin{bmatrix} 1.4 & 0.5 \\ 0.5 & 1.2 \end{bmatrix}$$
$$\ell_{1} = 1$$
$$\mathbf{B}_{2} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1.3 \end{bmatrix}$$
$$\ell_{2} = 0.2$$



$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{B}_1 \otimes k_1(\mathbf{X},\mathbf{X}) + \mathbf{B}_2 \otimes k_2(\mathbf{X},\mathbf{X})$$

$$\mathbf{B}_{1} = \begin{bmatrix} 1.4 & 0.5 \\ 0.5 & 1.2 \end{bmatrix}$$
$$\ell_{1} = 1$$
$$\mathbf{B}_{2} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1.3 \end{bmatrix}$$
$$\ell_{2} = 0.2$$



$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{B}_1 \otimes k_1(\mathbf{X},\mathbf{X}) + \mathbf{B}_2 \otimes k_2(\mathbf{X},\mathbf{X})$$

$$\mathbf{B}_{1} = \begin{bmatrix} 1.4 & 0.5 \\ 0.5 & 1.2 \end{bmatrix}$$
$$\ell_{1} = 1$$
$$\mathbf{B}_{2} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1.3 \end{bmatrix}$$
$$\ell_{2} = 0.2$$



LMC in Machine Learning and Statistics

- Used in machine learning for GPs for multivariate regression and in statistics for computer emulation of expensive multivariate computer codes.
- Imposes the correlation of the outputs explicitly through the set of coregionalization matrices.
- Setting B = I_p assumes outputs are conditionally independent given the parameters θ. (Minka and Picard, 1997; Lawrence and Platt, 2004; Yu et al., 2005).
- More recent approaches for multiple output modeling are different versions of the linear model of coregionalization.

Semiparametric Latent Factor Model

 Coregionalization matrices are rank 1 Teh et al. (2005). rewrite equation (??) as

$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \sum_{j=1}^{q} \mathbf{w}_{:,j} \mathbf{w}_{:,j}^{\top} \otimes k_{j}(\mathbf{X},\mathbf{X}).$$

- Like the Kalman filter, but each latent function has a *different* covariance.
- Authors suggest using an exponentiated quadratic characteristic length-scale for each input dimension.

$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{w}_{:,1}\mathbf{w}_{:,1}^{\top} \otimes k_1(\mathbf{X},\mathbf{X}) + \mathbf{w}_{:,2}\mathbf{w}_{:,2}^{\top} \otimes k_2(\mathbf{X},\mathbf{X})$$

$$\mathbf{w}_1 = \begin{bmatrix} 0.5\\1 \end{bmatrix}$$
$$\mathbf{w}_2 = \begin{bmatrix} 1\\0.5 \end{bmatrix}$$



$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{w}_{:,1}\mathbf{w}_{:,1}^{\top} \otimes k_1(\mathbf{X},\mathbf{X}) + \mathbf{w}_{:,2}\mathbf{w}_{:,2}^{\top} \otimes k_2(\mathbf{X},\mathbf{X})$$





$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{w}_{:,1}\mathbf{w}_{:,1}^{\top} \otimes k_1(\mathbf{X},\mathbf{X}) + \mathbf{w}_{:,2}\mathbf{w}_{:,2}^{\top} \otimes k_2(\mathbf{X},\mathbf{X})$$





$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{w}_{:,1}\mathbf{w}_{:,1}^{\top} \otimes k_1(\mathbf{X},\mathbf{X}) + \mathbf{w}_{:,2}\mathbf{w}_{:,2}^{\top} \otimes k_2(\mathbf{X},\mathbf{X})$$





$$\mathbf{K}(\mathbf{X},\mathbf{X}) = \mathbf{w}_{:,1}\mathbf{w}_{:,1}^{\top} \otimes k_1(\mathbf{X},\mathbf{X}) + \mathbf{w}_{:,2}\mathbf{w}_{:,2}^{\top} \otimes k_2(\mathbf{X},\mathbf{X})$$





Gaussian processes for Multi-task, Multi-output and Multi-class

- ► Bonilla et al. (2008) suggest ICM for multitask learning.
- ► Use a PPCA form for **B**: similar to our Kalman filter example.
- Refer to the autokrigeability effect as the cancellation of inter-task transfer.
- Also discuss the similarities between the multi-task GP and the ICM, and its relationship to the SLFM and the LMC.

Multitask Classification

- Mostly restricted to the case where the outputs are conditionally independent given the hyperparameters φ (Minka and Picard, 1997; Williams and Barber, 1998; Lawrence and Platt, 2004; Seeger and Jordan, 2004; Yu et al., 2005; Rasmussen and Williams, 2006).
- Intrinsic coregionalization model has been used in the multiclass scenario. Skolidis and Sanguinetti (2011) use the intrinsic coregionalization model for classification, by introducing a probit noise model as the likelihood.
- Posterior distribution is no longer analytically tractable: approximate inference is required.

- A statistical model used as a surrogate for a computationally expensive computer model.
- Higdon et al. (2008) use the linear model of coregionalization to model images representing the evolution of the implosion of steel cylinders.
- In Conti and O'Hagan (2009) use the ICM to model a vegetation model: called the Sheffield Dynamic Global Vegetation Model (Woodward et al., 1998).

References I

- W. V. Baxter and K.-I. Anjyo. Latent doodle space. In *EUROGRAPHICS*, volume 25, pages 477–485, Vienna, Austria, September 4-8 2006.
- C. M. Bishop, M. Svensén, and C. K. I. Williams. GTM: the Generative Topographic Mapping. *Neural Computation*, 10(1):215–234, 1998. [DOI].
- E. V. Bonilla, K. M. Chai, and C. K. I. Williams. Multi-task Gaussian process prediction. In J. C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20, Cambridge, MA, 2008. MIT Press.
- S. Conti and A. O'Hagan. Bayesian emulation of complex multi-output and dynamic computer models. *Journal of Statistical Planning and Inference*, 140 (3):640–651, 2009. [DOI].
- A. Damianou, C. H. Ek, M. K. Titsias, and N. D. Lawrence. Manifold relevance determination. In J. Langford and J. Pineau, editors, *Proceedings of the International Conference in Machine Learning*, volume 29, San Francisco, CA, 2012. Morgan Kauffman. [PDF].
- A. Damianou, M. K. Titsias, and N. D. Lawrence. Variational Gaussian process dynamical systems. In P. Bartlett, F. Peirrera, C. Williams, and J. Lafferty, editors, *Advances in Neural Information Processing Systems*, volume 24, Cambridge, MA, 2011. MIT Press. [PDF].

References II

- G. Della Gatta, M. Bansal, A. Ambesi-Impiombato, D. Antonini, C. Missero, and D. di Bernardo. Direct targets of the trp63 transcription factor revealed by a combination of gene expression profiling and reverse engineering. *Genome Research*, 18(6):939–948, Jun 2008. [URL]. [DOI].
- C. H. Ek, J. Rihan, P. Torr, G. Rogez, and N. D. Lawrence. Ambiguity modeling in latent spaces. In A. Popescu-Belis and R. Stiefelhagen, editors, *Machine Learning for Multimodal Interaction (MLMI 2008)*, LNCS, pages 62–73. Springer-Verlag, 28–30 June 2008a. [PDF].
- C. H. Ek, P. H. Torr, and N. D. Lawrence. Gaussian process latent variable models for human pose estimation. In A. Popescu-Belis, S. Renals, and H. Bourlard, editors, *Machine Learning for Multimodal Interaction (MLMI* 2007), volume 4892 of *LNCS*, pages 132–143, Brno, Czech Republic, 2008b. Springer-Verlag. [PDF].
- Z. Ghahramani, editor. *Proceedings of the International Conference in Machine Learning*, volume 24, 2007. Omnipress. [Google Books].
- P. Goovaerts. *Geostatistics For Natural Resources Evaluation*. Oxford University Press, 1997. [Google Books].
- K. Grochow, S. L. Martin, A. Hertzmann, and Z. Popovic. Style-based inverse kinematics. In ACM Transactions on Graphics (SIGGRAPH 2004), pages 522–531, 2004.

References III

- J. D. Helterbrand and N. A. C. Cressie. Universal cokriging under intrinsic coregionalization. *Mathematical Geology*, 26(2):205–226, 1994.
- D. M. Higdon, J. Gattiker, B. Williams, and M. Rightley. Computer model calibration using high dimensional output. *Journal of the American Statistical Association*, 103(482):570–583, 2008.
- A. G. Journel and C. J. Huijbregts. *Mining Geostatistics*. Academic Press, London, 1978. [Google Books].
- A. A. Kalaitzis and N. D. Lawrence. A simple approach to ranking differentially expressed gene expression time courses through Gaussian process regression. *BMC Bioinformatics*, 12(180), 2011. [DOI].
- A. Klami and S. Kaski. Local dependent components analysis. In Ghahramani (2007). [Google Books].
- A. Klami and S. Kaski. Probabilistic approach to detecting dependencies between data sets. *Neurocomputing*, 72:39–46, 2008.
- J. B. Kruskal. Multidimensional scaling by optimizing goodness-of-fit to a nonmetric hypothesis. *Psychometrika*, 29(1):1–28, 1964. [DOI].

References IV

- P. S. Laplace. *Essai philosophique sur les probabilités*. Courcier, Paris, 2nd edition, 1814. Sixth edition of 1840 translated and repreinted (1951) as *A Philosophical Essay on Probabilities*, New York: Dover; fifth edition of 1825 reprinted 1986 with notes by Bernard Bru, Paris: Christian Bourgois Éditeur, translated by Andrew Dale (1995) as *Philosophical Essay on Probabilities*, New York:Springer-Verlag.
- N. D. Lawrence. Gaussian process models for visualisation of high dimensional data. In S. Thrun, L. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems*, volume 16, pages 329–336, Cambridge, MA, 2004. MIT Press.
- N. D. Lawrence. Probabilistic non-linear principal component analysis with Gaussian process latent variable models. *Journal of Machine Learning Research*, 6:1783–1816, 11 2005.
- N. D. Lawrence and A. J. Moore. Hierarchical Gaussian process latent variable models. In Ghahramani (2007), pages 481–488. [Google Books] . [PDF].
- N. D. Lawrence and J. C. Platt. Learning to learn with the informative vector machine. In R. Greiner and D. Schuurmans, editors, *Proceedings of the International Conference in Machine Learning*, volume 21, pages 512–519. Omnipress, 2004. [PDF].

References V

- N. D. Lawrence and J. Quiñonero Candela. Local distance preservation in the GP-LVM through back constraints. In W. Cohen and A. Moore, editors, *Proceedings of the International Conference in Machine Learning*, volume 23, pages 513–520. Omnipress, 2006. [Google Books]. [PDF].
- G. Leen and C. Fyfe. A Gaussian process latent variable model formulation of canonical correlation analysis. Bruges (Belgium), 26-28 April 2006 2006.
- S. Levine, J. M. Wang, A. Haraux, Z. Popović, and V. Koltun. Continuous character control with low-dimensional embeddings. *ACM Transactions on Graphics (SIGGRAPH 2012)*, 31(4), 2012.
- D. Lowe and M. E. Tipping. Neuroscale: Novel topographic feature extraction with radial basis function networks. In M. C. Mozer, M. I. Jordan, and T. Petsche, editors, *Advances in Neural Information Processing Systems*, volume 9, pages 543–549, Cambridge, MA, 1997. MIT Press.
- C. Lu and X. Tang. Surpassing human-level face verification performance on LFW with GaussianFace. Technical report,
- D. J. C. MacKay. Bayesian neural networks and density networks. Nuclear Instruments and Methods in Physics Research, A, 354(1):73–80, 1995. [DOI].
- D. J. C. MacKay. *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, Cambridge, U.K., 2003. [Google Books].

References VI

- K. V. Mardia, J. T. Kent, and J. M. Bibby. *Multivariate analysis*. Academic Press, London, 1979. [Google Books].
- T. P. Minka and R. W. Picard. Learning how to learn is learning with point sets. Available on-line., 1997. [URL]. Revised 1999, available at http://www.stat.cmu.edu/~{}minka/.
- R. Navaratnam, A. Fitzgibbon, and R. Cipolla. The joint manifold model for semi-supervised multi-valued regression. In *IEEE International Conference on Computer Vision (ICCV)*. IEEE Computer Society Press, 2007.
- R. M. Neal. *Bayesian Learning for Neural Networks*. Springer, 1996. Lecture Notes in Statistics 118.
- J. Oakley and A. O'Hagan. Bayesian inference for the uncertainty distribution of computer model outputs. *Biometrika*, 89(4):769–784, 2002.
- V. Priacuriu and I. D. Reid. Nonlinear shape manifolds as shape priors in level set segmentation and tracking. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2011a.
- V. Priacuriu and I. D. Reid. Shared shape spaces. In *IEEE International Conference on Computer Vision (ICCV)*, 2011b.
- C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, MA, 2006. [Google Books].

References VII

- S. T. Roweis. EM algorithms for PCA and SPCA. In M. I. Jordan, M. J. Kearns, and S. A. Solla, editors, *Advances in Neural Information Processing Systems*, volume 10, pages 626–632, Cambridge, MA, 1998. MIT Press.
- S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000. [DOI].
- J. W. Sammon. A nonlinear mapping for data structure analysis. *IEEE Transactions on Computers*, C-18(5):401–409, 1969. [DOI].
- B. Schölkopf, A. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10:1299–1319, 1998. [DOI].
- M. Seeger and M. I. Jordan. Sparse Gaussian Process Classification With Multiple Classes. Technical Report 661, Department of Statistics, University of California at Berkeley,
- A. P. Shon, K. Grochow, A. Hertzmann, and R. P. N. Rao. Learning shared latent structure for image synthesis and robotic imitation. In Weiss et al. (2006).
- G. Skolidis and G. Sanguinetti. Bayesian multitask classification with Gaussian process priors. *IEEE Transactions on Neural Networks*, 22(12):2011 2021, 2011.
- M. L. Stein. *Interpolation of Spatial Data: Some Theory for Kriging*. Springer-Verlag, 1999. [Google Books].

References VIII

- Y. W. Teh, M. Seeger, and M. I. Jordan. Semiparametric latent factor models. In R. G. Cowell and Z. Ghahramani, editors, *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics*, pages 333–340, Barbados, 6-8 January 2005. Society for Artificial Intelligence and Statistics.
- J. B. Tenenbaum, V. de Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500): 2319–2323, 2000. [DOI].
- M. E. Tipping and C. M. Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society, B,* 6(3):611–622, 1999. [PDF]. [DOI].
- M. K. Titsias and N. D. Lawrence. Bayesian Gaussian process latent variable model. In Y. W. Teh and D. M. Titterington, editors, *Proceedings of the Thirteenth International Workshop on Artificial Intelligence and Statistics*, volume 9, pages 844–851, Chia Laguna Resort, Sardinia, Italy, 13-16 May 2010. JMLR W&CP 9. [PDF].
- L. R. Tucker. An inter-battery method of factor analysis. *Psychometrika*, 23(2): 111–136, 1958.
- R. Urtasun and T. Darrell. Discriminative Gaussian process latent variable model for classification. In Ghahramani (2007). [Google Books] .

References IX

- R. Urtasun, D. J. Fleet, and P. Fua. 3D people tracking with Gaussian process dynamical models. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 238–245, New York, U.S.A., 17–22 Jun. 2006. IEEE Computer Society Press.
- R. Urtasun, D. J. Fleet, A. Hertzmann, and P. Fua. Priors for people tracking from small training sets. In *IEEE International Conference on Computer Vision (ICCV)*, pages 403–410, Bejing, China, 17–21 Oct. 2005. IEEE Computer Society Press.
- S. Virtanen, A. Klami, and S. Kaski. Bayesian CCA via group sparsity. In L. Getoor and T. Scheffer, editors, *Proceedings of the International Conference in Machine Learning*, volume 28, 2011.
- H. Wackernagel. *Multivariate Geostatistics: An Introduction With Applications*. Springer-Verlag, 3rd edition, 2003. [Google Books].
- J. M. Wang, D. J. Fleet, and A. Hertzmann. Gaussian process dynamical models. In Weiss et al. (2006).
- J. M. Wang, D. J. Fleet, and A. Hertzmann. Gaussian process dynamical models for human motion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(2):283–298, 2008. ISSN 0162-8828. [DOI].
- Y. Weiss, B. Schölkopf, and J. C. Platt, editors. *Advances in Neural Information Processing Systems*, volume 18, Cambridge, MA, 2006. MIT Press.

- C. K. Williams and D. Barber. Bayesian Classification with Gaussian processes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20 (12):1342–1351, 1998.
- C. K. I. Williams. Computation with infinite neural networks. *Neural Computation*, 10(5):1203–1216, 1998.
- I. Woodward, M. R. Lomas, and R. A. Betts. Vegetation-climate feedbacks in a greenhouse world. *Philosophical Transactions: Biological Sciences*, 353(1365): 29–39, 1998.
- K. Yu, V. Tresp, and A. Schwaighofer. Learning Gaussian processes from multiple tasks. In *Proceedings of the 22nd International Conference on Machine Learning (ICML 2005)*, pages 1012–1019, 2005.