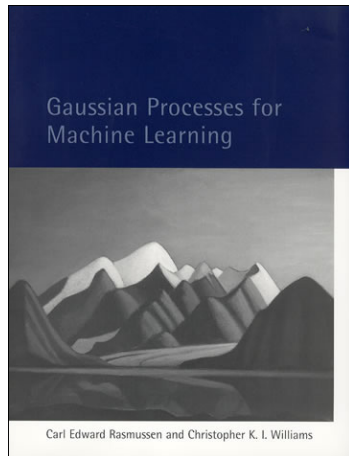


Introduction to Gaussian Processes

Neil D. Lawrence

GPMC
6th February 2017





Outline

The Gaussian Density

Covariance from Basis Functions

Outline

The Gaussian Density

Covariance from Basis Functions

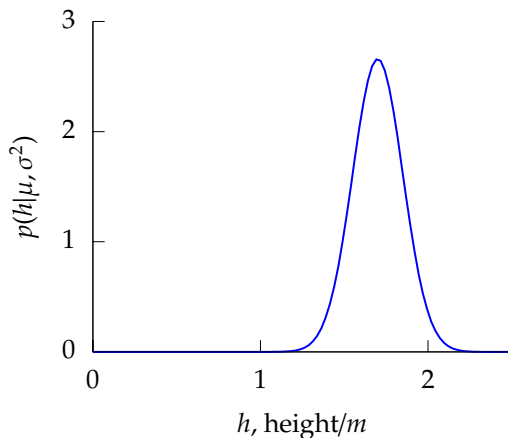
The Gaussian Density

- ▶ Perhaps the most common probability density.

$$\begin{aligned} p(y|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right) \\ &\triangleq \mathcal{N}(y|\mu, \sigma^2) \end{aligned}$$

- ▶ The Gaussian density.

Gaussian Density



The Gaussian PDF with $\mu = 1.7$ and variance $\sigma^2 = 0.0225$. Mean shown as red line. It could represent the heights of a population of students.

$$\mathcal{N}(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right)$$

σ^2 is the variance of the density and μ is the mean.

Two Important Gaussian Properties

Sum of Gaussians

- ▶ Sum of Gaussian variables is also Gaussian.

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

Two Important Gaussian Properties

Sum of Gaussians

- Sum of Gaussian variables is also Gaussian.

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

And the sum is distributed as

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

Two Important Gaussian Properties

Sum of Gaussians

- Sum of Gaussian variables is also Gaussian.

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

And the sum is distributed as

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

(*Aside:* As sum increases, sum of non-Gaussian, finite variance variables is also Gaussian [central limit theorem].)

Two Important Gaussian Properties

Sum of Gaussians

- Sum of Gaussian variables is also Gaussian.

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

And the sum is distributed as

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

(*Aside:* As sum increases, sum of non-Gaussian, finite variance variables is also Gaussian [central limit theorem].)

Two Important Gaussian Properties

Scaling a Gaussian

- ▶ Scaling a Gaussian leads to a Gaussian.

Two Important Gaussian Properties

Scaling a Gaussian

- ▶ Scaling a Gaussian leads to a Gaussian.

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

Two Important Gaussian Properties

Scaling a Gaussian

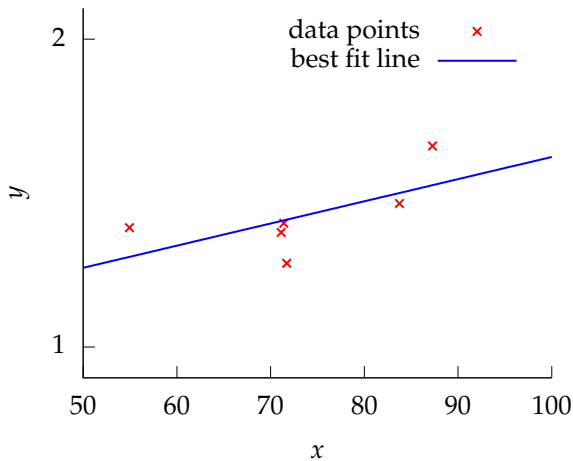
- ▶ Scaling a Gaussian leads to a Gaussian.

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

And the scaled density is distributed as

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

Linear Function

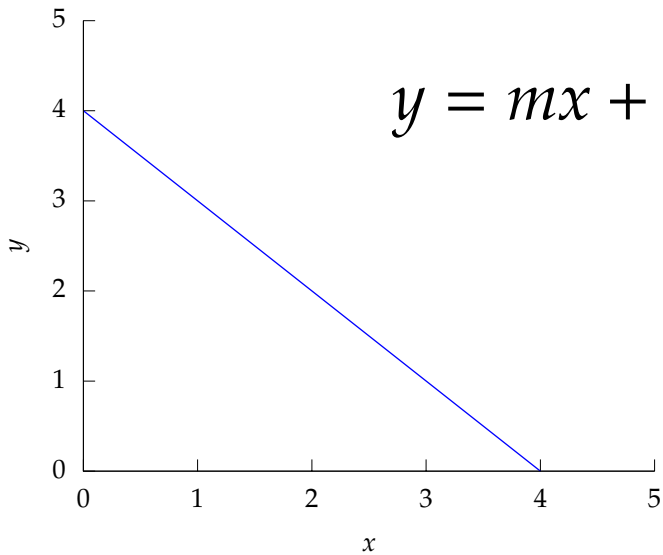


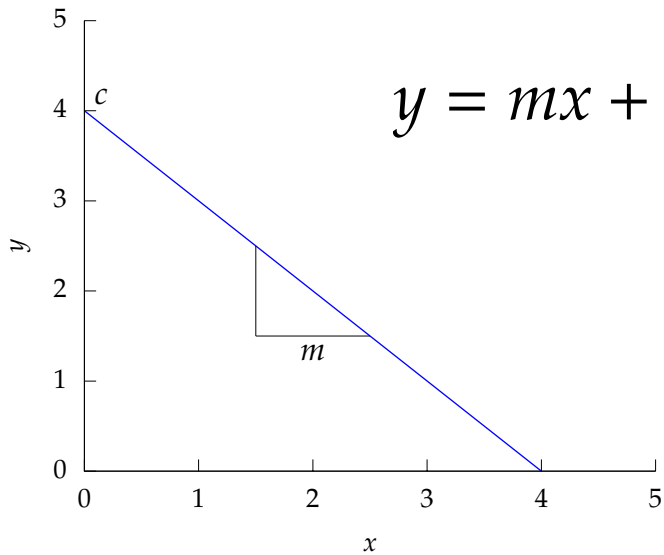
A linear regression between x and y .

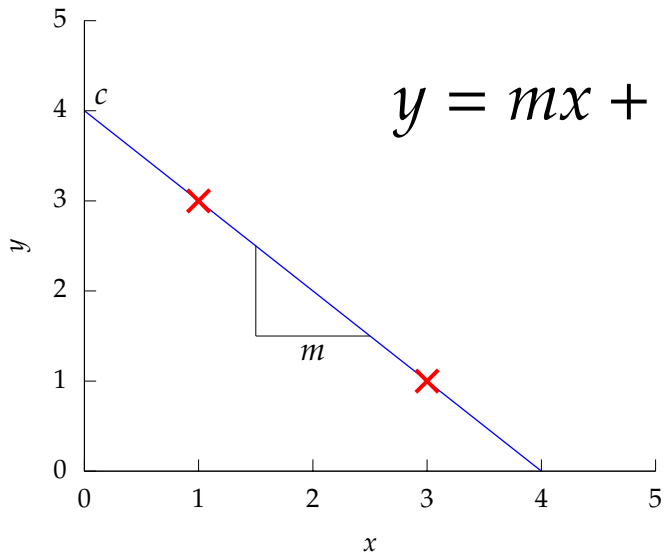
Regression Examples

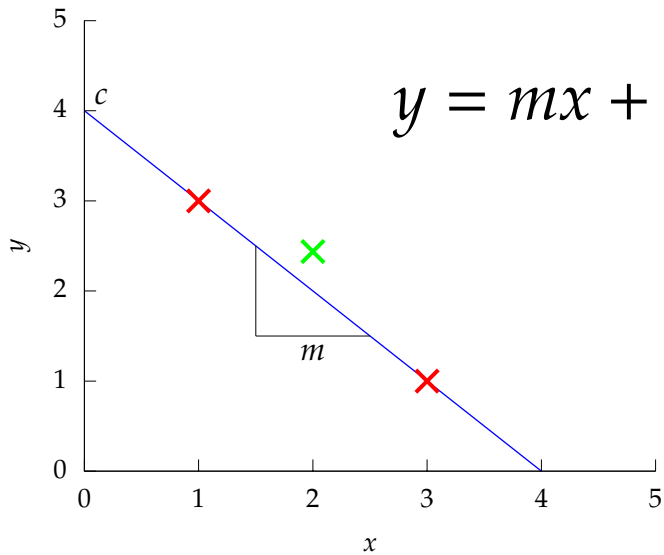
- ▶ Predict a real value, y_i given some inputs x_i .
- ▶ Predict quality of meat given spectral measurements (Tecator data).
- ▶ Radiocarbon dating, the C14 calibration curve: predict age given quantity of C14 isotope.
- ▶ Predict quality of different Go or Backgammon moves given expert rated training data.

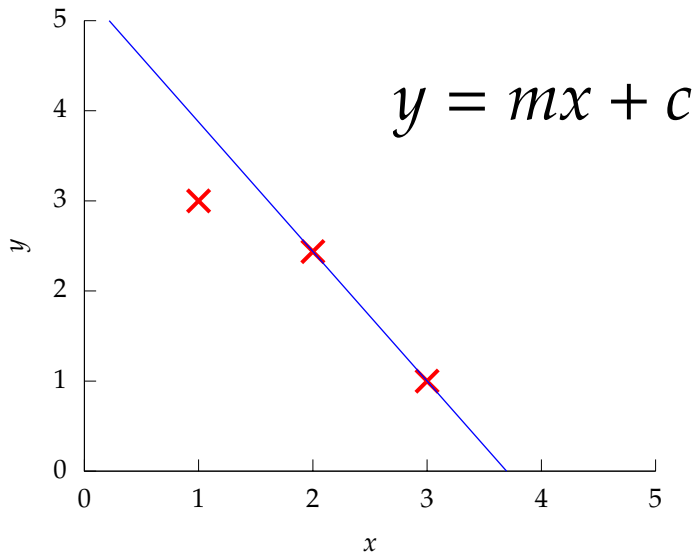
$$y = mx + c$$

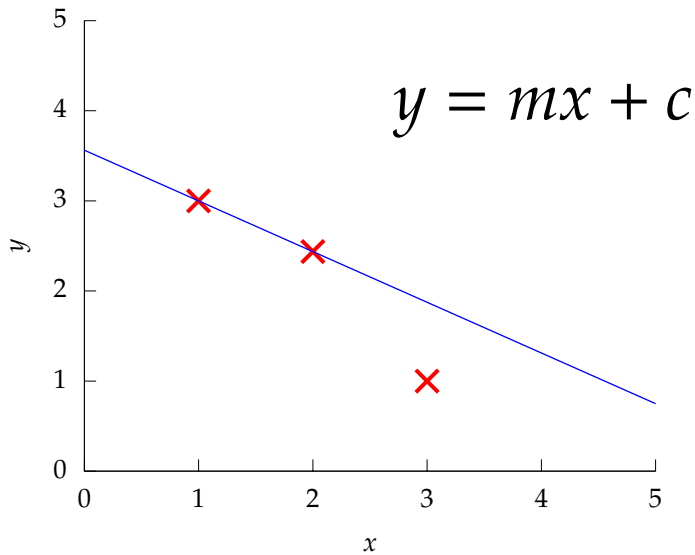


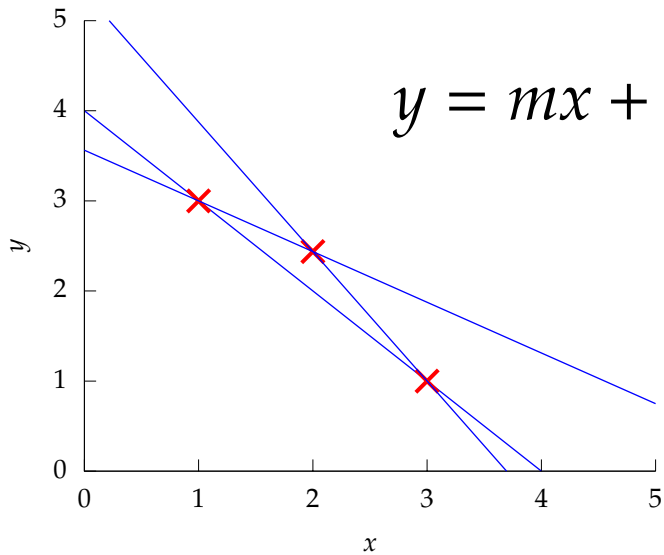












$$y = mx + c$$

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'elle-même 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 parvenu 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

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

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: $x = 3, y = 1$

$$1 = 3m + c + \epsilon_2$$

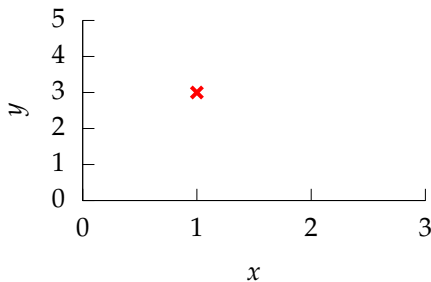
point 3: $x = 2, y = 2.5$

$$2.5 = 2m + c + \epsilon_3$$

Underdetermined System

What about two unknowns and *one* observation?

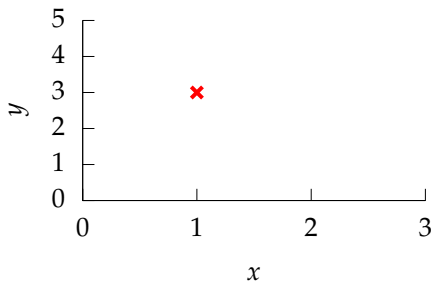
$$y_1 = mx_1 + c$$



Underdetermined System

Can compute m given c .

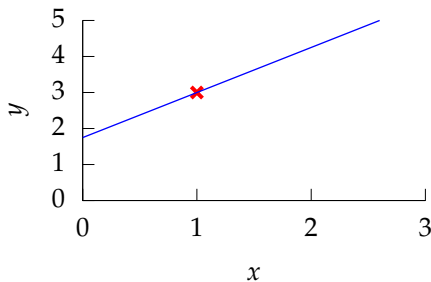
$$m = \frac{y_1 - c}{x}$$



Underdetermined System

Can compute m given c .

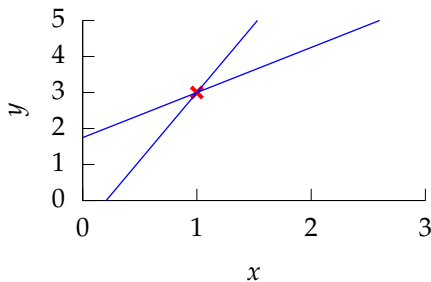
$$c = 1.75 \Rightarrow m = 1.25$$



Underdetermined System

Can compute m given c .

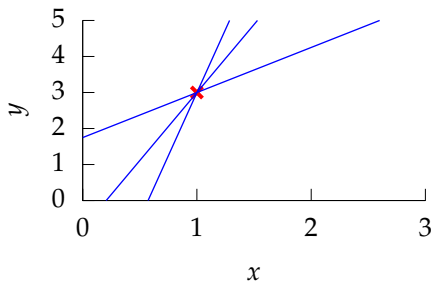
$$c = -0.777 \Rightarrow m = 3.78$$



Underdetermined System

Can compute m given c .

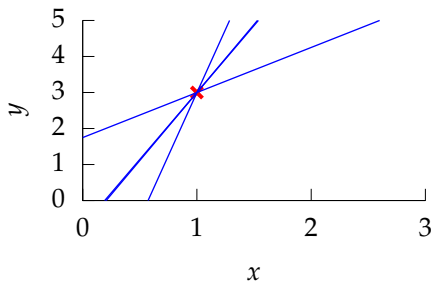
$$c = -4.01 \Rightarrow m = 7.01$$



Underdetermined System

Can compute m given c .

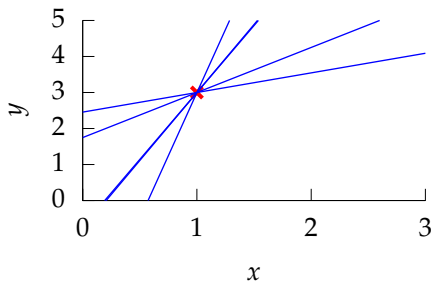
$$c = -0.718 \implies m = 3.72$$



Underdetermined System

Can compute m given c .

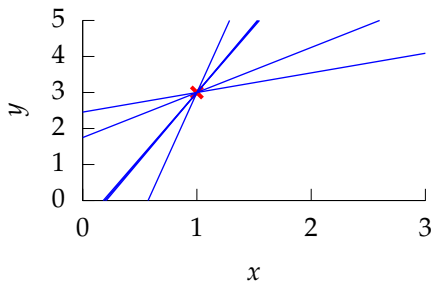
$$c = 2.45 \implies m = 0.545$$



Underdetermined System

Can compute m given c .

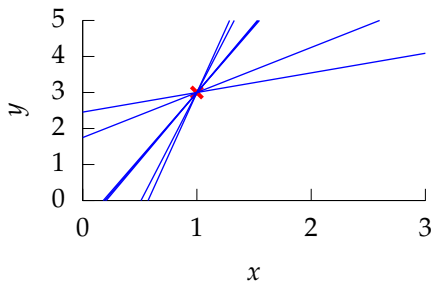
$$c = -0.657 \implies m = 3.66$$



Underdetermined System

Can compute m given c .

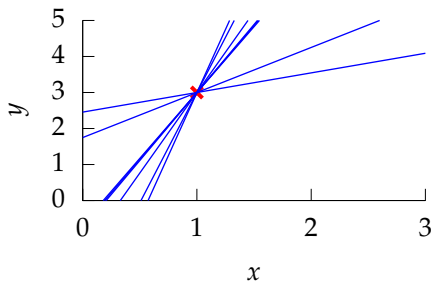
$$c = -3.13 \Rightarrow m = 6.13$$



Underdetermined System

Can compute m given c .

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



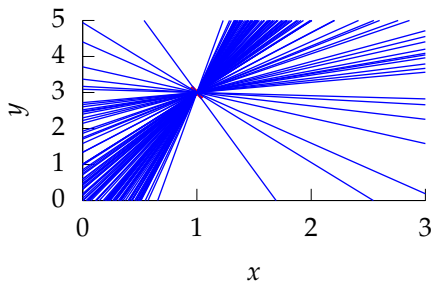
Underdetermined System

Can compute m given c .

Assume

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

we find a distribution of solutions.



Probability for Under- and Overdetermined

- ▶ To deal with overdetermined introduced probability distribution for 'variable', ϵ_i .
- ▶ For underdetermined system introduced probability distribution for 'parameter', c .
- ▶ This is known as a Bayesian treatment.

Multivariate Prior Distributions

- ▶ For general Bayesian inference need multivariate priors.
- ▶ E.g. for multivariate linear regression:

$$y_i = \sum_j w_j x_{i,j} + \epsilon_i$$

(where we've dropped c for convenience), we need a prior over \mathbf{w} .

- ▶ This motivates a *multivariate* Gaussian density.
- ▶ We will use the multivariate Gaussian to put a prior *directly* on the function (a Gaussian process).

Multivariate Prior Distributions

- ▶ For general Bayesian inference need multivariate priors.
- ▶ E.g. for multivariate linear regression:

$$y_i = \mathbf{w}^\top \mathbf{x}_{i,:} + \epsilon_i$$

(where we've dropped c for convenience), we need a prior over \mathbf{w} .

- ▶ This motivates a *multivariate* Gaussian density.
- ▶ We will use the multivariate Gaussian to put a prior *directly* on the function (a Gaussian process).

Prior Distribution

- ▶ Bayesian inference requires a prior on the parameters.
- ▶ The prior represents your belief *before* you see the data of the likely value of the parameters.
- ▶ For linear regression, consider a Gaussian prior on the intercept:

$$c \sim \mathcal{N}(0, \alpha_1)$$

Posterior Distribution

- ▶ Posterior distribution is found by combining the prior with the likelihood.
- ▶ Posterior distribution is your belief *after* you see the data of the likely value of the parameters.
- ▶ The posterior is found through **Bayes' Rule**

$$p(c|y) = \frac{p(y|c)p(c)}{p(y)}$$

Bayes Update

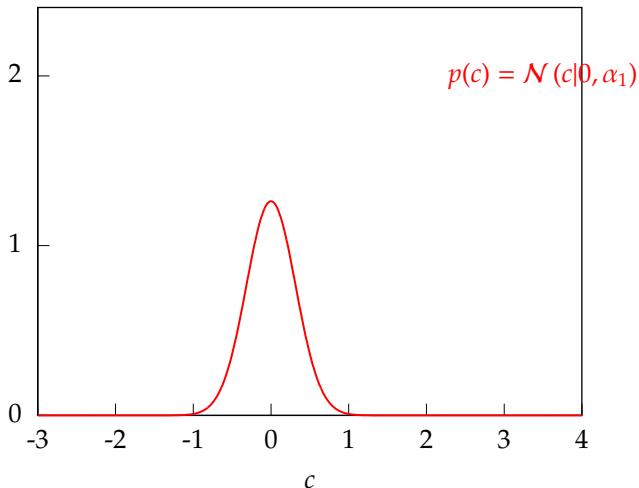


Figure: A Gaussian prior combines with a Gaussian likelihood for a Gaussian posterior.

Bayes Update

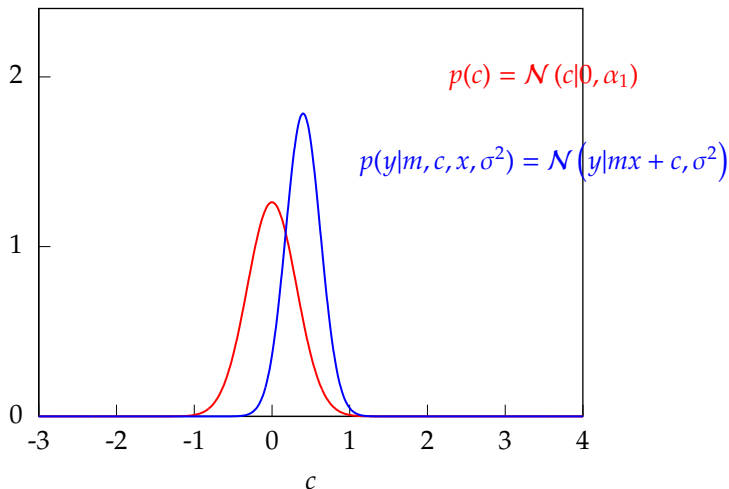


Figure: A Gaussian prior combines with a Gaussian likelihood for a Gaussian posterior.

Bayes Update

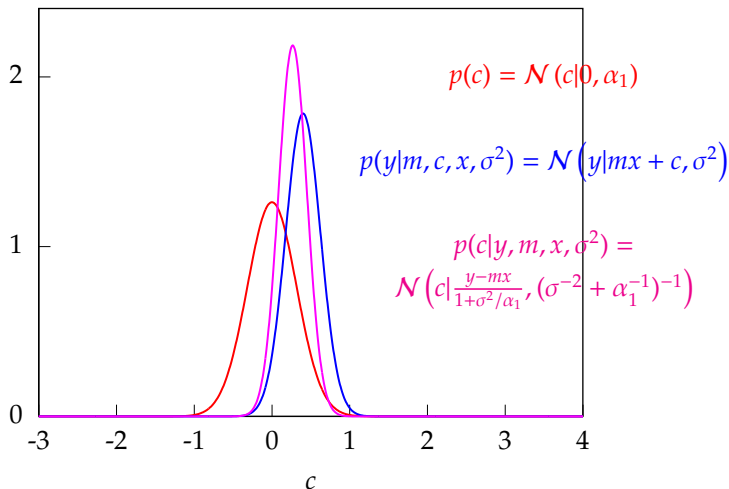


Figure: A Gaussian prior combines with a Gaussian likelihood for a Gaussian posterior.

Stages to Derivation of the Posterior

- ▶ Multiply likelihood by prior
 - ▶ they are “exponentiated quadratics”, the answer is always also an exponentiated quadratic because $\exp(a^2) \exp(b^2) = \exp(a^2 + b^2)$.
- ▶ Complete the square to get the resulting density in the form of a Gaussian.
- ▶ Recognise the mean and (co)variance of the Gaussian. This is the estimate of the posterior.

Multivariate Regression Likelihood

- ▶ Noise corrupted data point

$$y_i = \mathbf{w}^\top \mathbf{x}_{i,:} + \epsilon_i$$

Multivariate Regression Likelihood

- ▶ Noise corrupted data point

$$y_i = \mathbf{w}^\top \mathbf{x}_{i,:} + \epsilon_i$$

- ▶ Multivariate regression likelihood:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_{i,:})^2\right)$$

Multivariate Regression Likelihood

- ▶ Noise corrupted data point

$$y_i = \mathbf{w}^\top \mathbf{x}_{i,:} + \epsilon_i$$

- ▶ Multivariate regression likelihood:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_{i,:})^2\right)$$

- ▶ Now use a multivariate Gaussian prior:

$$p(\mathbf{w}) = \frac{1}{(2\pi\alpha)^{\frac{p}{2}}} \exp\left(-\frac{1}{2\alpha} \mathbf{w}^\top \mathbf{w}\right)$$

Two Dimensional Gaussian

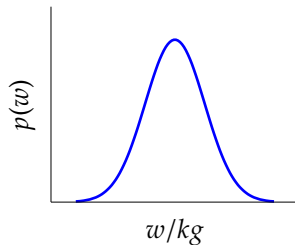
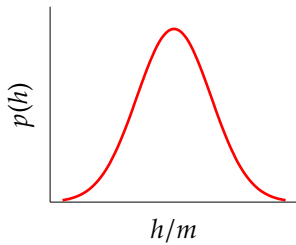
- ▶ Consider height, h/m and weight, w/kg .
- ▶ Could sample height from a distribution:

$$p(h) \sim \mathcal{N}(1.7, 0.0225)$$

- ▶ And similarly weight:

$$p(w) \sim \mathcal{N}(75, 36)$$

Height and Weight Models

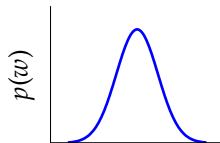
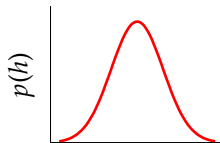
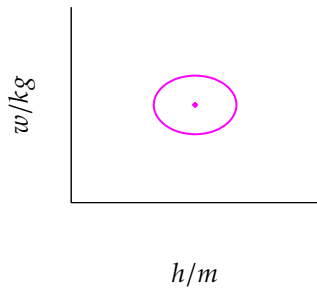


Gaussian distributions for height and weight.

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

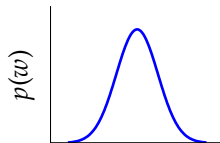
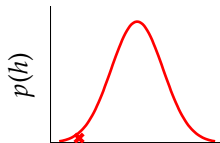
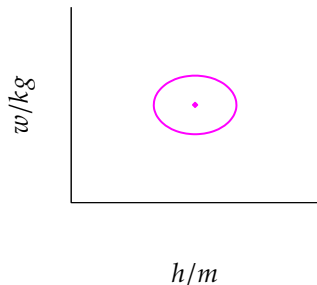


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

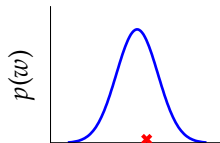
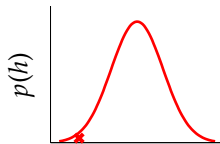
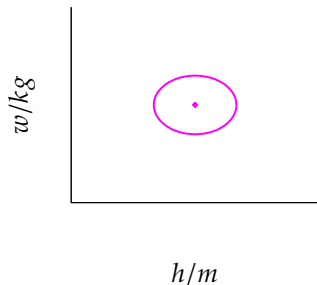


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

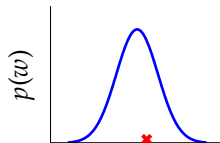
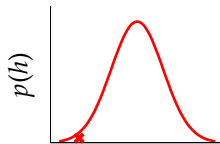
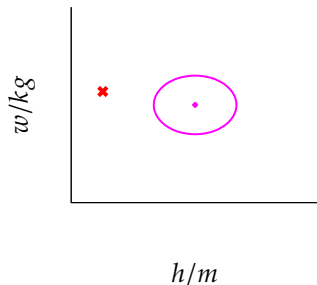


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

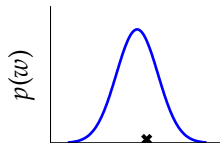
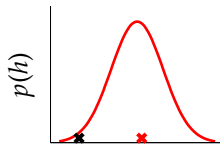
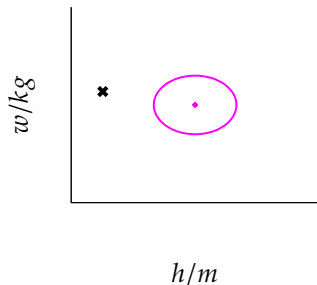


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

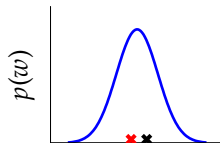
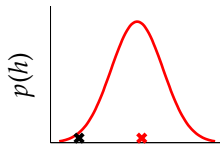
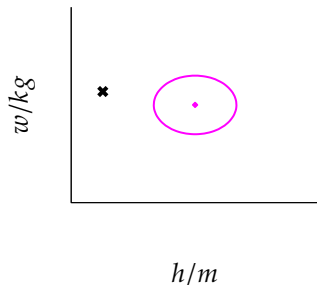


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

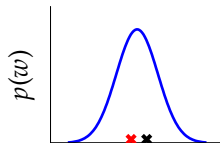
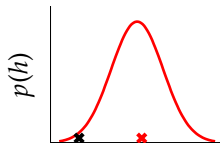
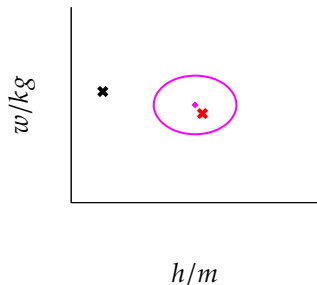


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

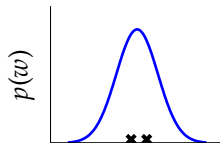
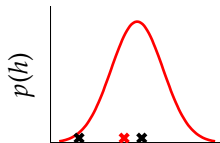
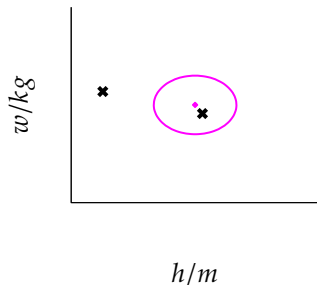


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

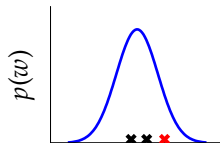
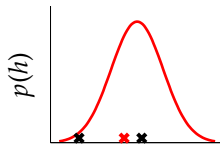
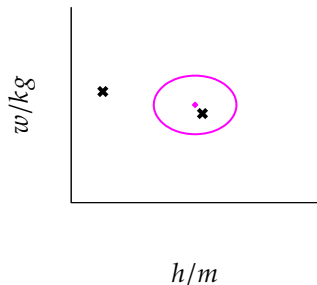


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

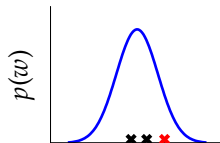
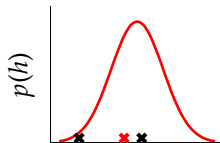
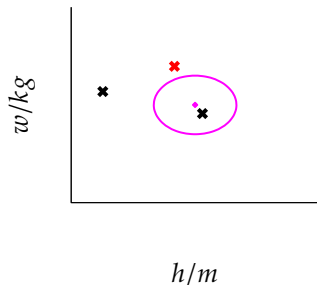


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

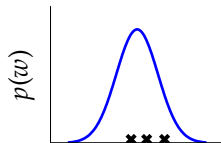
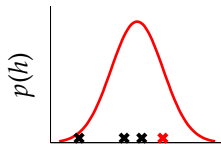
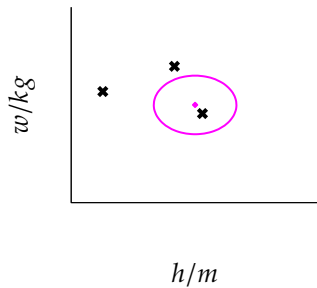


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

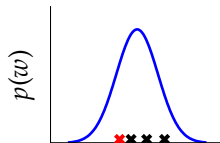
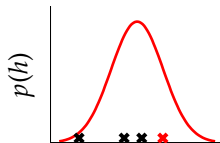
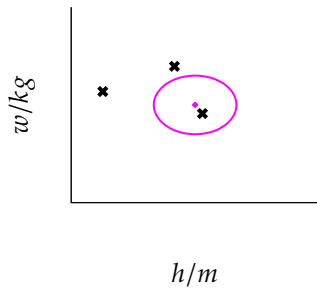


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

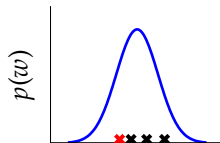
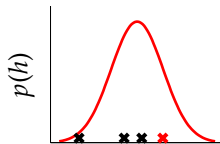
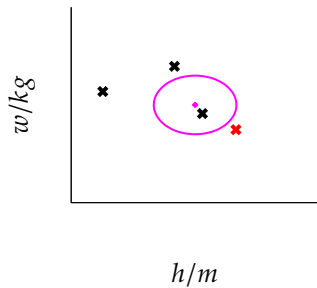


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

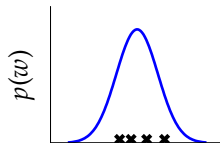
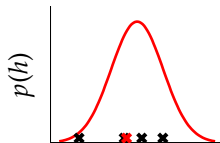
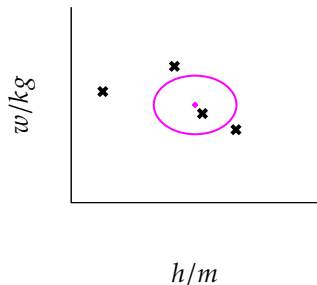


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

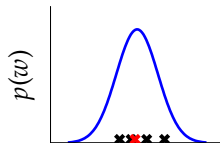
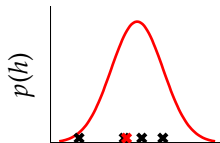
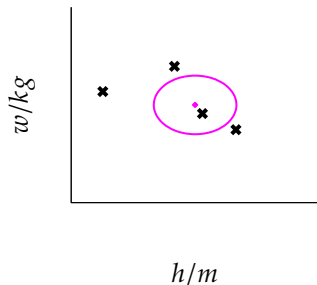


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

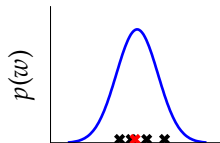
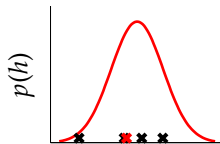
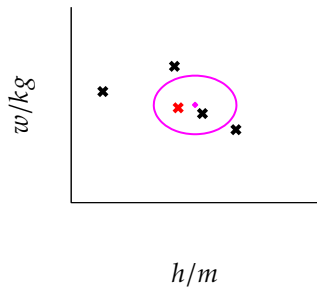


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

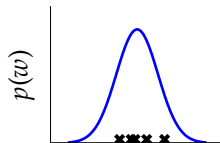
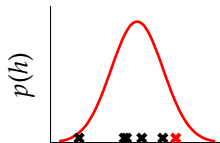
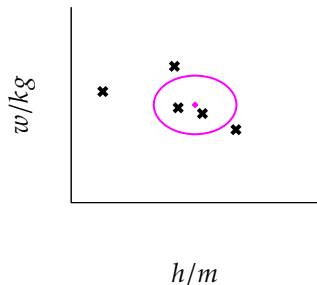


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

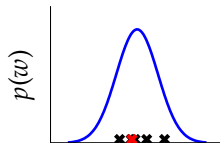
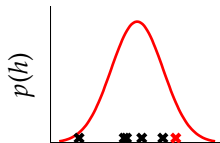
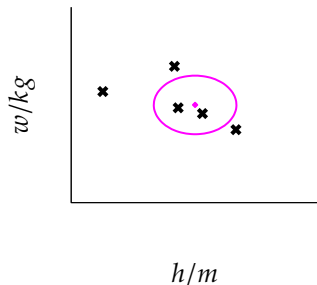


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

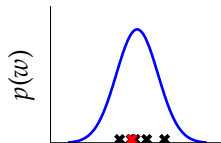
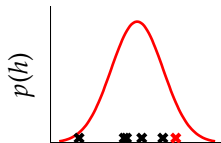
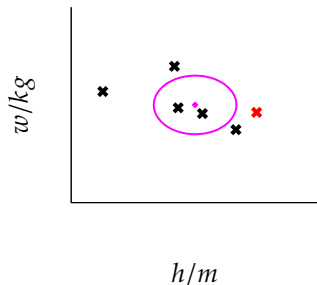


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

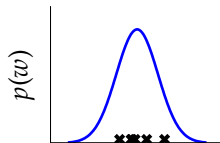
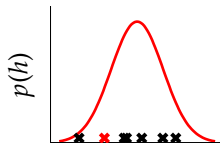
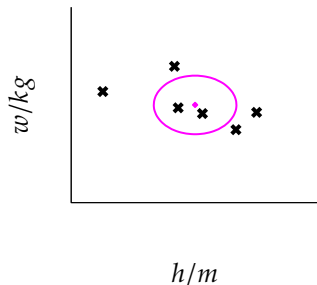


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

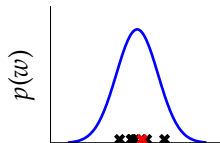
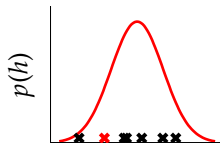
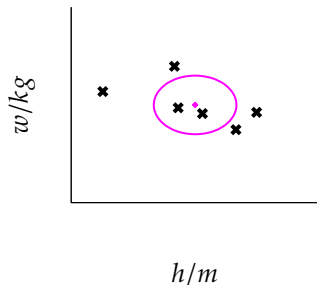


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

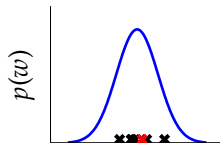
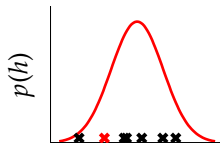
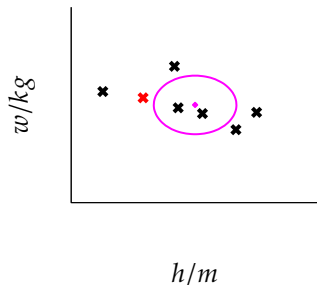


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution

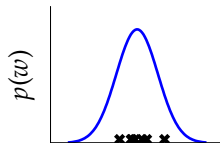
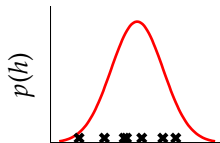
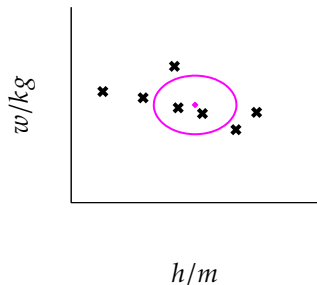


Samples of height and weight

Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution



Samples of height and weight

Independence Assumption

- ▶ This assumes height and weight are independent.

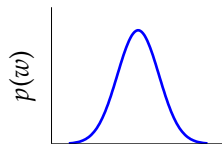
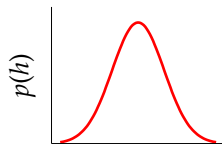
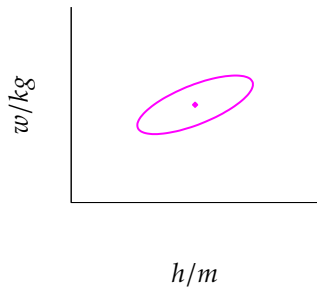
$$p(h, w) = p(h)p(w)$$

- ▶ In reality they are dependent (body mass index) = $\frac{w}{h^2}$.

Sampling Two Dimensional Variables

Marginal Distributions

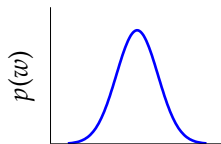
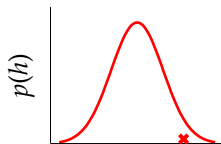
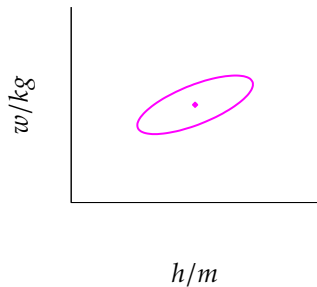
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

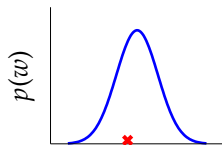
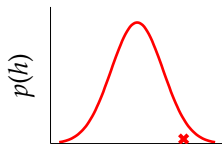
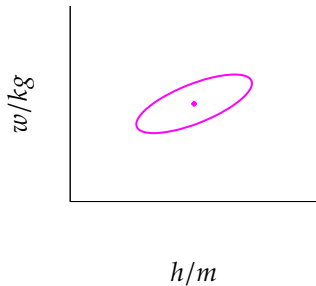
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

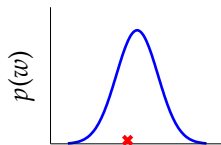
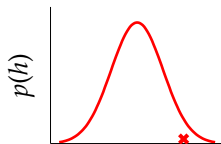
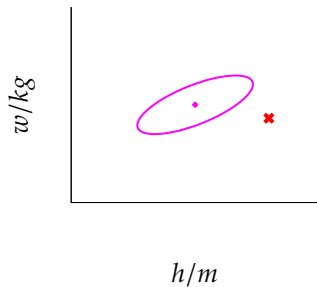
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

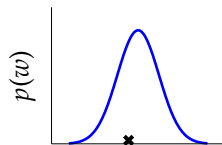
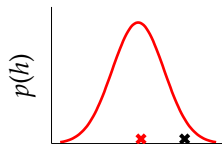
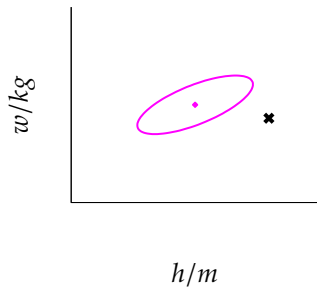
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

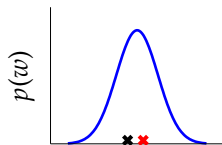
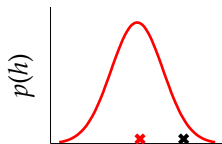
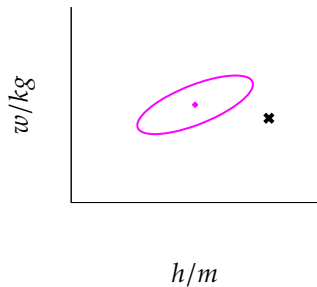
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

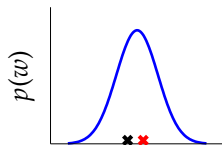
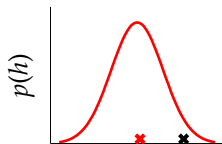
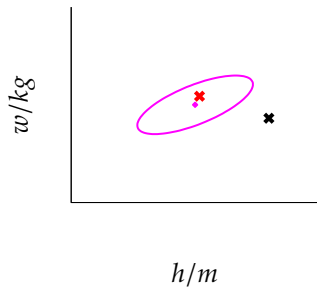
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

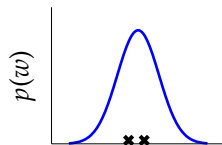
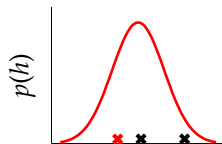
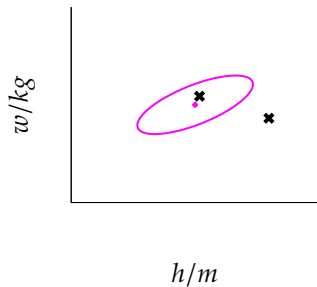
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

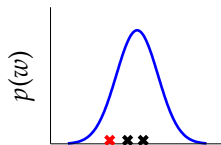
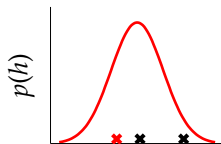
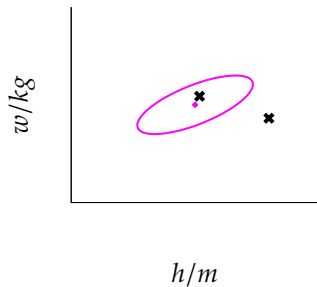
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

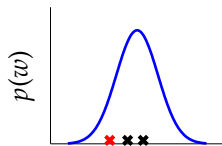
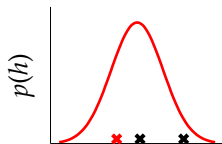
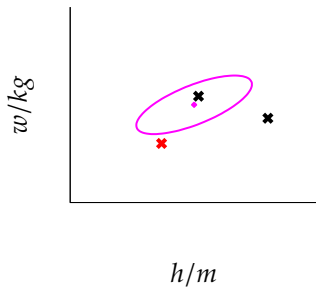
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

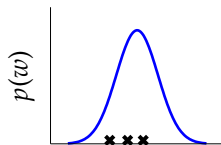
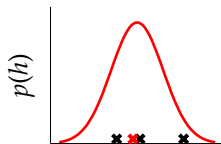
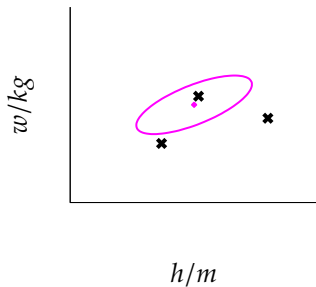
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

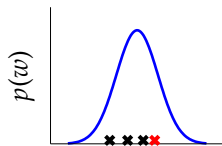
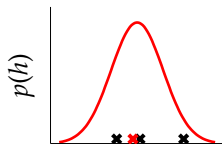
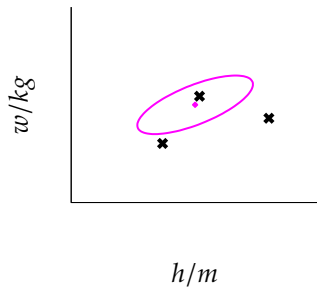
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

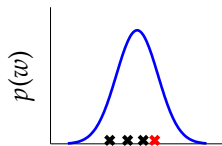
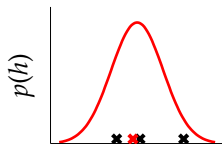
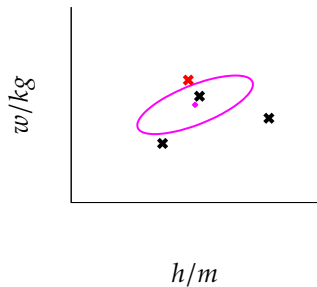
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

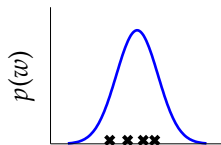
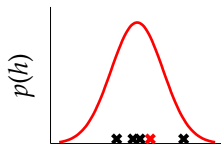
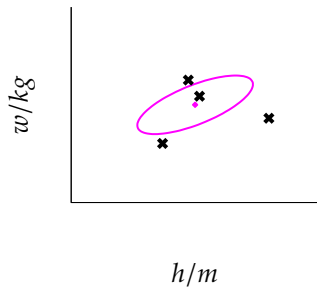
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

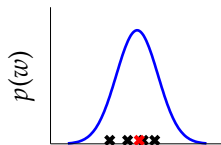
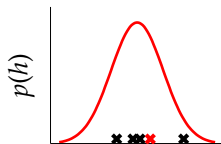
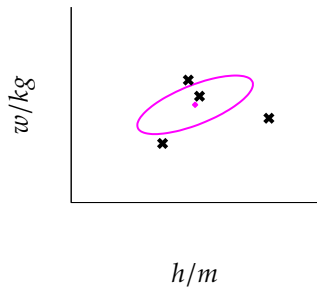
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

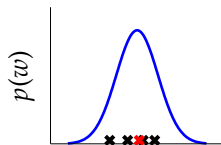
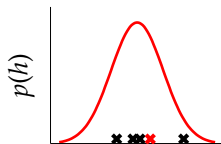
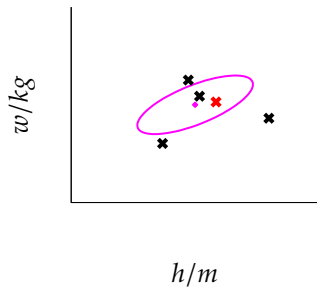
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

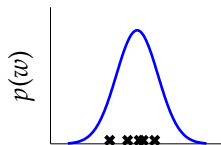
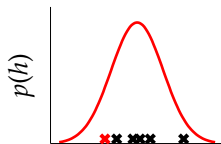
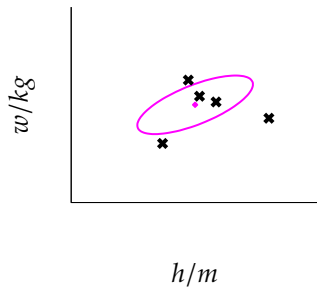
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

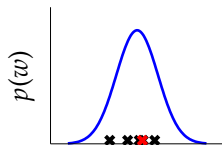
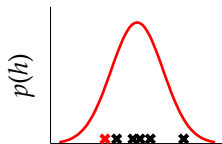
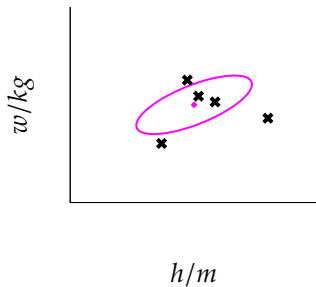
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

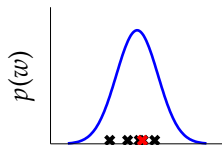
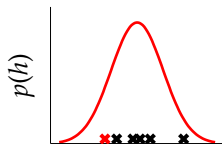
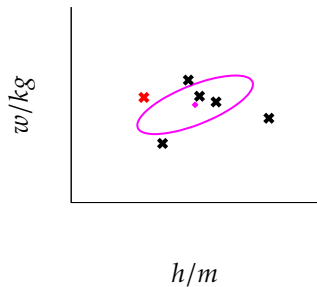
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

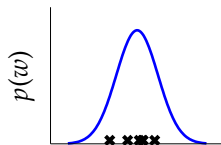
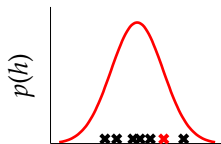
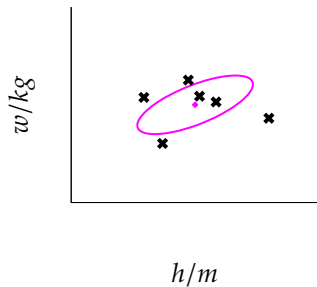
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

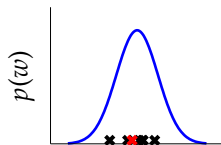
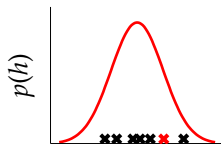
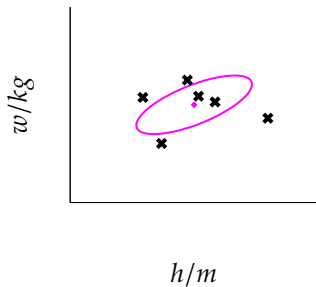
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

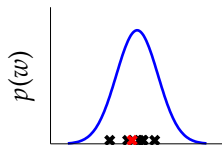
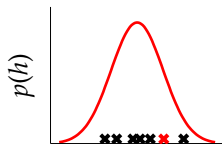
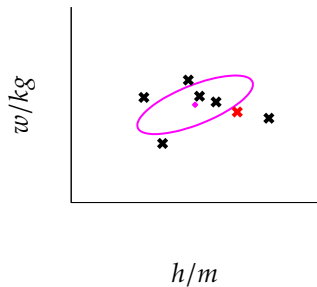
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

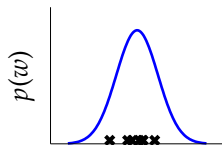
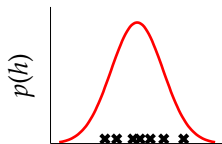
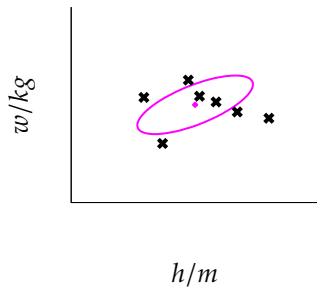
Joint Distribution



Sampling Two Dimensional Variables

Marginal Distributions

Joint Distribution



Independent Gaussians

$$p(w, h) = p(w)p(h)$$

Independent Gaussians

$$p(w, h) = \frac{1}{\sqrt{2\pi\sigma_1^2}\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{1}{2}\left(\frac{(w - \mu_1)^2}{\sigma_1^2} + \frac{(h - \mu_2)^2}{\sigma_2^2}\right)\right)$$

Independent Gaussians

$$p(w, h) = \frac{1}{\sqrt{2\pi\sigma_1^2} \sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{1}{2} \left(\begin{bmatrix} w \\ h \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}\right)^\top \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}^{-1} \left(\begin{bmatrix} w \\ h \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}\right)\right)$$

Independent Gaussians

$$p(\mathbf{y}) = \frac{1}{|2\pi\mathbf{D}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^\top \mathbf{D}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right)$$

Correlated Gaussian

Form correlated from original by rotating the data space using matrix \mathbf{R} .

$$p(\mathbf{y}) = \frac{1}{|2\pi\mathbf{D}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^{\top}\mathbf{D}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right)$$

Correlated Gaussian

Form correlated from original by rotating the data space using matrix \mathbf{R} .

$$p(\mathbf{y}) = \frac{1}{|2\pi\mathbf{D}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{R}^\top \mathbf{y} - \mathbf{R}^\top \boldsymbol{\mu})^\top \mathbf{D}^{-1}(\mathbf{R}^\top \mathbf{y} - \mathbf{R}^\top \boldsymbol{\mu})\right)$$

Correlated Gaussian

Form correlated from original by rotating the data space using matrix \mathbf{R} .

$$p(\mathbf{y}) = \frac{1}{|2\pi\mathbf{D}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^\top \mathbf{R}\mathbf{D}^{-1}\mathbf{R}^\top(\mathbf{y} - \boldsymbol{\mu})\right)$$

this gives a covariance matrix:

$$\mathbf{C}^{-1} = \mathbf{R}\mathbf{D}^{-1}\mathbf{R}^\top$$

Correlated Gaussian

Form correlated from original by rotating the data space using matrix \mathbf{R} .

$$p(\mathbf{y}) = \frac{1}{|2\pi\mathbf{C}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^{\top}\mathbf{C}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right)$$

this gives a covariance matrix:

$$\mathbf{C} = \mathbf{R}\mathbf{D}\mathbf{R}^{\top}$$

Recall Univariate Gaussian Properties

1. Sum of Gaussian variables is also Gaussian.

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

Recall Univariate Gaussian Properties

1. Sum of Gaussian variables is also Gaussian.

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

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

Recall Univariate Gaussian Properties

1. Sum of Gaussian variables is also Gaussian.

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

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

2. Scaling a Gaussian leads to a Gaussian.

Recall Univariate Gaussian Properties

1. Sum of Gaussian variables is also Gaussian.

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

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

2. Scaling a Gaussian leads to a Gaussian.

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

Recall Univariate Gaussian Properties

1. Sum of Gaussian variables is also Gaussian.

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

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

2. Scaling a Gaussian leads to a Gaussian.

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

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

Multivariate Consequence

- If

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Multivariate Consequence

- ▶ If

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- ▶ And

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

Multivariate Consequence

- ▶ If

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- ▶ And

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

- ▶ Then

$$\mathbf{y} \sim \mathcal{N}(\mathbf{W}\boldsymbol{\mu}, \mathbf{W}\boldsymbol{\Sigma}\mathbf{W}^\top)$$

Sampling a Function

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.

Gaussian Distribution Sample

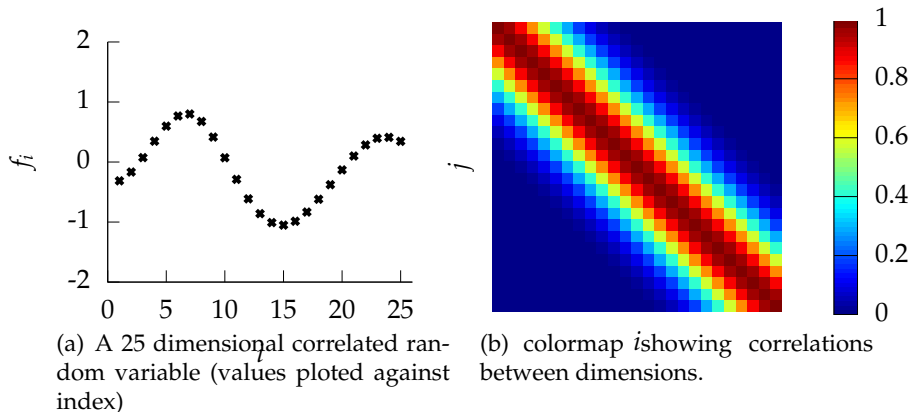


Figure: A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample

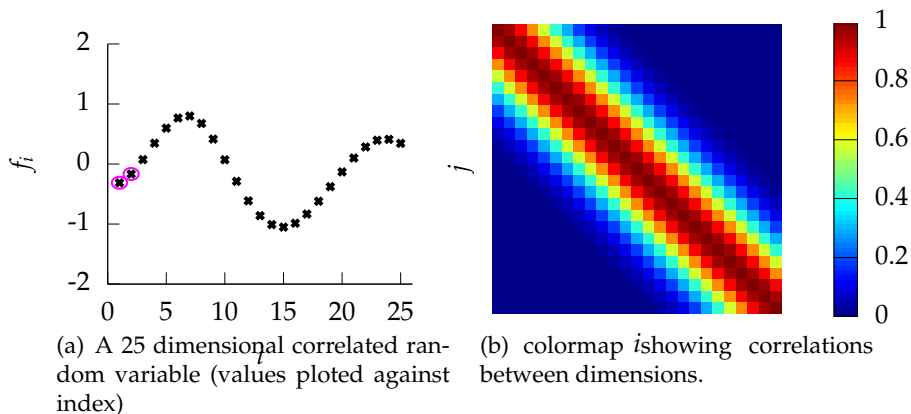


Figure: A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample

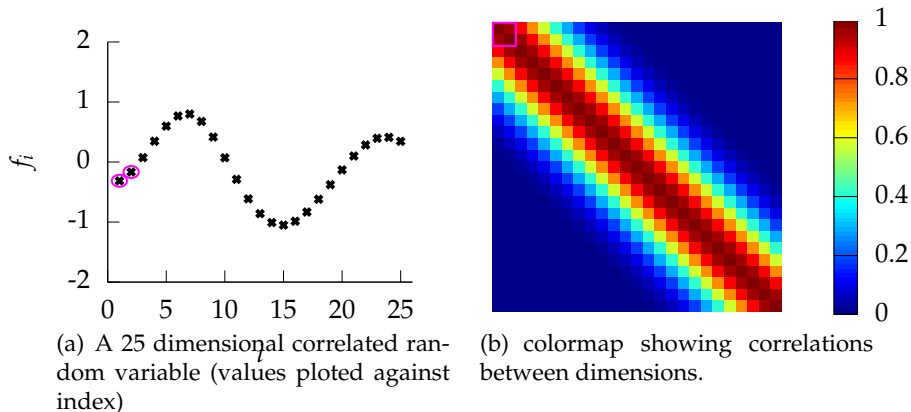


Figure: A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample

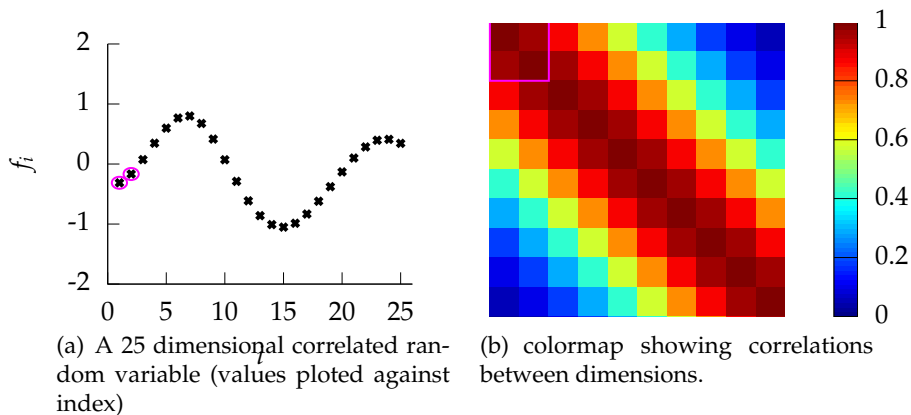


Figure: A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample

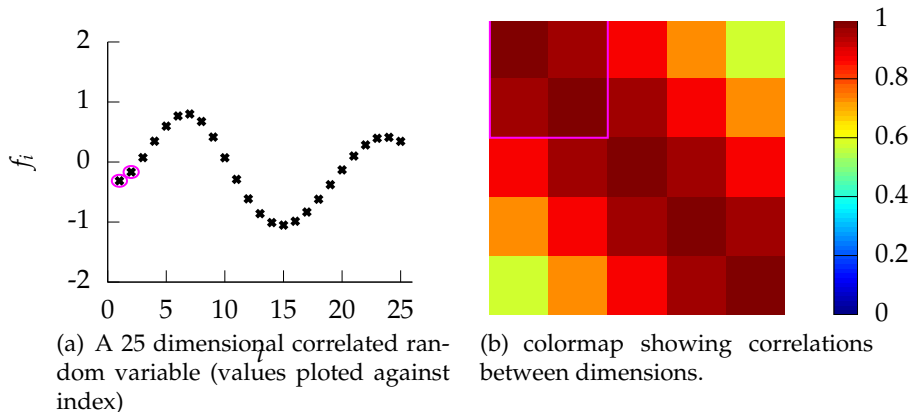


Figure: A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample

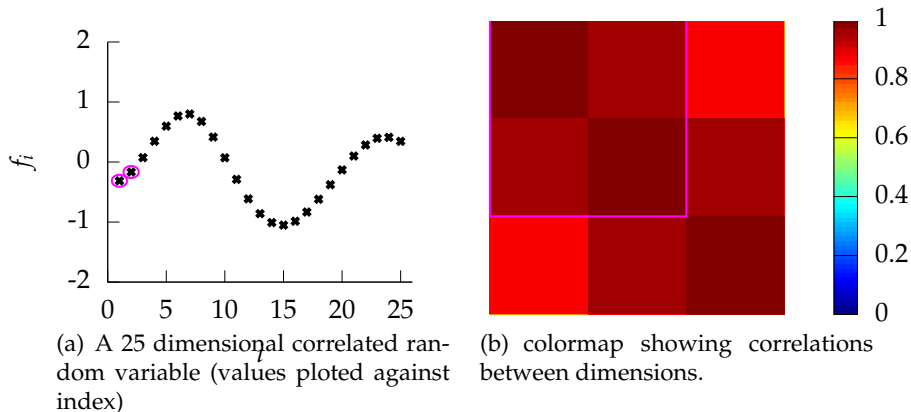


Figure: A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample

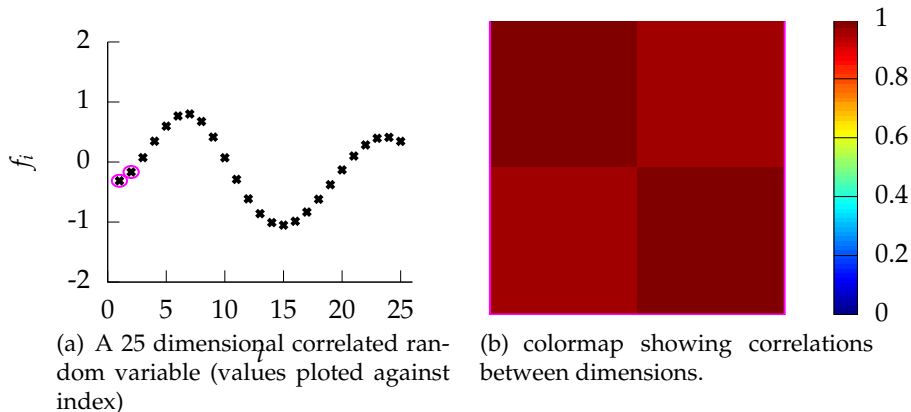
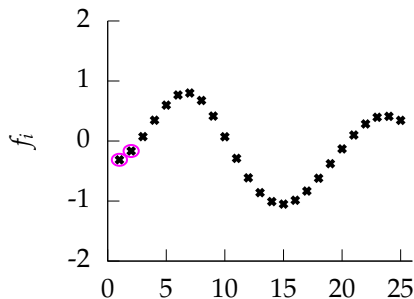


Figure: A sample from a 25 dimensional Gaussian distribution.

Gaussian Distribution Sample



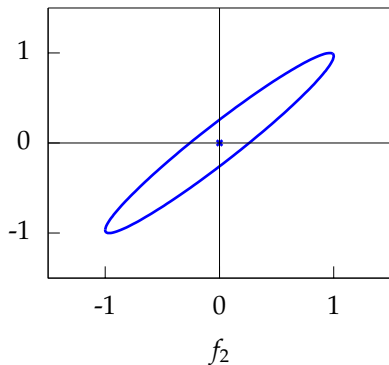
(a) A 25 dimensional correlated random variable (values plotted against index)

$$\begin{bmatrix} 1 & 0.96587 \\ 0.96587 & 1 \end{bmatrix}$$

(b) correlation between f_1 and f_2 .

Figure: A sample from a 25 dimensional Gaussian distribution.

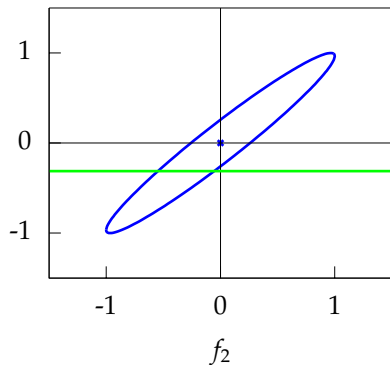
Prediction of f_2 from f_1



$$\begin{bmatrix} 1 & 0.96587 \\ 0.96587 & 1 \end{bmatrix}$$

- The single contour of the Gaussian density represents the joint distribution, $p(f_1, f_2)$.

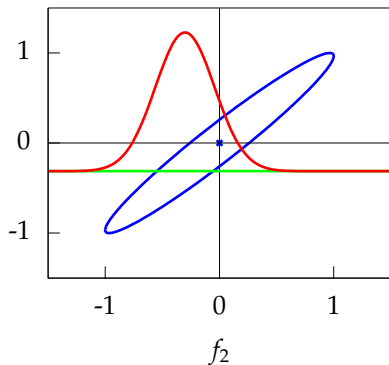
Prediction of f_2 from f_1



$$\begin{bmatrix} 1 & 0.96587 \\ 0.96587 & 1 \end{bmatrix}$$

- ▶ The single contour of the Gaussian density represents the **joint distribution**, $p(f_1, f_2)$.
- ▶ We observe that $f_1 = -0.313$.

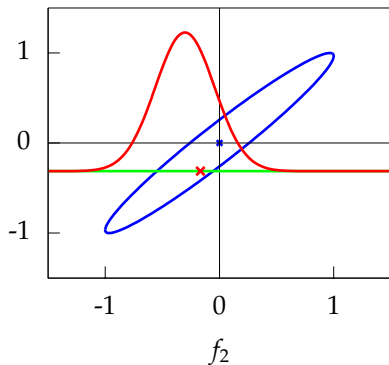
Prediction of f_2 from f_1



$$\begin{bmatrix} 1 & 0.96587 \\ 0.96587 & 1 \end{bmatrix}$$

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

Prediction of f_2 from f_1



$$\begin{bmatrix} 1 & 0.96587 \\ 0.96587 & 1 \end{bmatrix}$$

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

Prediction with Correlated Gaussians

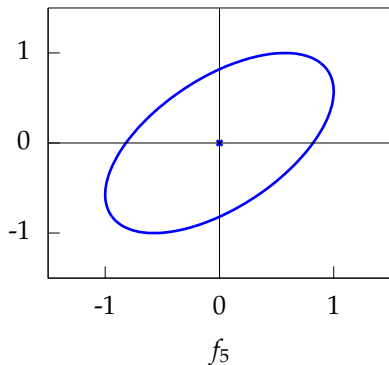
- ▶ Prediction of f_2 from f_1 requires *conditional density*.
- ▶ Conditional density is *also* Gaussian.

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

where covariance of joint density is given by

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

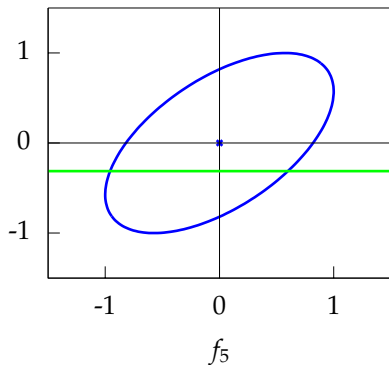
Prediction of f_5 from f_1



$$\begin{bmatrix} 1 & 0.57375 \\ 0.57375 & 1 \end{bmatrix}$$

- The single contour of the Gaussian density represents the joint distribution, $p(f_1, f_5)$.

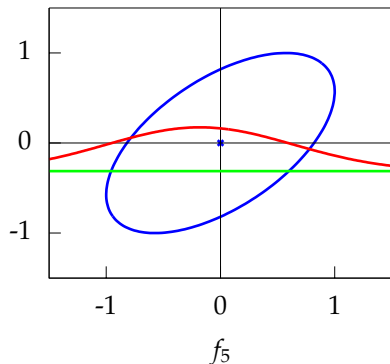
Prediction of f_5 from f_1



$$\begin{bmatrix} 1 & 0.57375 \\ 0.57375 & 1 \end{bmatrix}$$

- ▶ The single contour of the Gaussian density represents the **joint distribution, $p(f_1, f_5)$** .
- ▶ We observe that **$f_1 = -0.313$** .

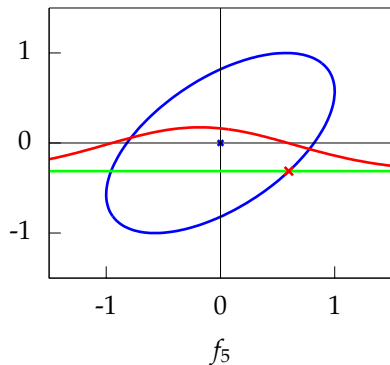
Prediction of f_5 from f_1



$$\begin{bmatrix} 1 & 0.57375 \\ 0.57375 & 1 \end{bmatrix}$$

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

Prediction of f_5 from f_1



$$\begin{bmatrix} 1 & 0.57375 \\ 0.57375 & 1 \end{bmatrix}$$

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

Prediction with Correlated Gaussians

- ▶ Prediction of \mathbf{f}_* from \mathbf{f} requires multivariate *conditional density*.
- ▶ Multivariate conditional density is *also* Gaussian.

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

- ▶ Here covariance of joint density is given by

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

Prediction with Correlated Gaussians

- ▶ Prediction of \mathbf{f}_* from \mathbf{f} requires multivariate *conditional density*.
- ▶ Multivariate conditional density is *also* Gaussian.

$$p(\mathbf{f}_*|\mathbf{f}) = \mathcal{N}(\mathbf{f}_*|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\mu} = \mathbf{K}_{*,\mathbf{f}}\mathbf{K}_{\mathbf{f},\mathbf{f}}^{-1}\mathbf{f}$$

$$\boldsymbol{\Sigma} = \mathbf{K}_{*,*} - \mathbf{K}_{*,\mathbf{f}}\mathbf{K}_{\mathbf{f},\mathbf{f}}^{-1}\mathbf{K}_{\mathbf{f},*}$$

- ▶ Here covariance of joint density is given by

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

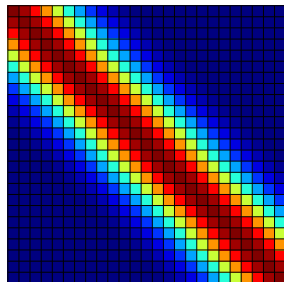
Covariance Functions

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 \mathbf{x} .
- ▶ For the example above it was based on Euclidean distance.
- ▶ The covariance function is also known as a kernel.



\mathbf{x}_1

\mathbf{x}_2

Covariance Functions

Where did this covariance matrix come from?

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

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

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

$$x_1 = -3.0, x_2 = 1.20, \text{ and } x_3 = 1.40 \text{ with } \ell = 2.00 \text{ and } \alpha = 1.00.$$

Outline

The Gaussian Density

Covariance from Basis Functions

Basis Function Form

Radial basis functions commonly have the form

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

- Basis function maps data into a “feature space” in which a linear sum is a non linear function.

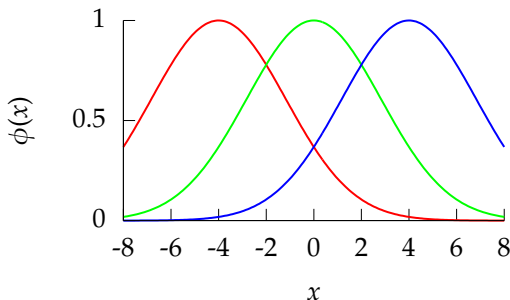


Figure: A set of radial basis functions with width $\ell = 2$ and location parameters $\boldsymbol{\mu} = [-4 \ 0 \ 4]^\top$.

Basis Function Representations

- Represent a function by a linear sum over a basis,

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

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

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

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

Random Functions

Functions derived
using:

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

where elements of \mathbf{w}
are independently
sampled from a
Gaussian density,

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

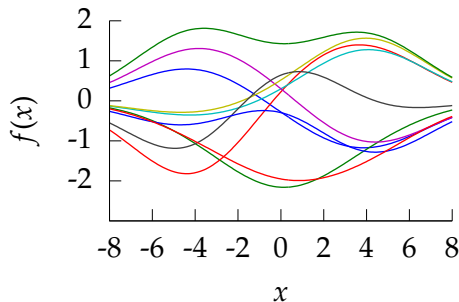


Figure: Functions sampled using the basis set from figure 3. Each line is a separate sample, generated by a weighted sum of the basis set. The weights, \mathbf{w} are sampled from a Gaussian density with variance $\alpha = 1$.

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

computed at training data gives a vector

$$\mathbf{f} = \mathbf{\Phi} \mathbf{w}.$$

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

computed at training data gives a vector

$$\mathbf{f} = \mathbf{\Phi} \mathbf{w}.$$

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

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

computed at training data gives a vector

$$\mathbf{f} = \mathbf{\Phi} \mathbf{w}.$$

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

\mathbf{w} and \mathbf{f} are only related by an *inner product*.

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

computed at training data gives a vector

$$\mathbf{f} = \Phi \mathbf{w}.$$

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

\mathbf{w} and \mathbf{f} are only related by an *inner product*.

$\Phi \in \mathbb{R}^{n \times p}$ is a *design matrix*

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

computed at training data gives a vector

$$\mathbf{f} = \mathbf{\Phi} \mathbf{w}.$$

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

\mathbf{w} and \mathbf{f} are only related by an *inner product*.

$\mathbf{\Phi} \in \mathbb{R}^{n \times p}$ is a *design matrix*

$\mathbf{\Phi}$ is fixed and non-stochastic for a given training set.

Direct Construction of Covariance Matrix

Use matrix notation to write function,

$$f(\mathbf{x}_i; \mathbf{w}) = \sum_{k=1}^m w_k \phi_k(\mathbf{x}_i)$$

computed at training data gives a vector

$$\mathbf{f} = \Phi \mathbf{w}.$$

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

\mathbf{w} and \mathbf{f} are only related by an *inner product*.

$\Phi \in \mathbb{R}^{n \times p}$ is a *design matrix*

Φ is fixed and non-stochastic for a given training set.

\mathbf{f} is Gaussian distributed.

Expectations

- ▶ We have

$$\langle \mathbf{f} \rangle = \Phi \langle \mathbf{w} \rangle .$$

We use $\langle \cdot \rangle$ to denote expectations under prior distributions.

Expectations

- ▶ We have

$$\langle \mathbf{f} \rangle = \Phi \langle \mathbf{w} \rangle .$$

- ▶ Prior mean of \mathbf{w} was zero giving

$$\langle \mathbf{f} \rangle = \mathbf{0} .$$

We use $\langle \cdot \rangle$ to denote expectations under prior distributions.

Expectations

- ▶ We have

$$\langle \mathbf{f} \rangle = \mathbf{\Phi} \langle \mathbf{w} \rangle .$$

- ▶ Prior mean of \mathbf{w} was zero giving

$$\langle \mathbf{f} \rangle = \mathbf{0} .$$

- ▶ Prior covariance of \mathbf{f} is

$$\mathbf{K} = \langle \mathbf{f} \mathbf{f}^\top \rangle - \langle \mathbf{f} \rangle \langle \mathbf{f} \rangle^\top$$

We use $\langle \cdot \rangle$ to denote expectations under prior distributions.

Expectations

- ▶ We have

$$\langle \mathbf{f} \rangle = \mathbf{\Phi} \langle \mathbf{w} \rangle.$$

- ▶ Prior mean of \mathbf{w} was zero giving

$$\langle \mathbf{f} \rangle = \mathbf{0}.$$

- ▶ Prior covariance of \mathbf{f} is

$$\mathbf{K} = \langle \mathbf{f} \mathbf{f}^\top \rangle - \langle \mathbf{f} \rangle \langle \mathbf{f} \rangle^\top$$

$$\langle \mathbf{f} \mathbf{f}^\top \rangle = \mathbf{\Phi} \langle \mathbf{w} \mathbf{w}^\top \rangle \mathbf{\Phi}^\top,$$

giving

$$\mathbf{K} = \alpha \mathbf{\Phi} \mathbf{\Phi}^\top.$$

We use $\langle \cdot \rangle$ to denote expectations under prior distributions.

Covariance between Two Points

- ▶ The prior covariance between two points \mathbf{x}_i and \mathbf{x}_j is

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

Covariance between Two Points

- The prior covariance between two points \mathbf{x}_i and \mathbf{x}_j is

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

or in sum notation

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \sum_{k=1}^m \phi_k(\mathbf{x}_i) \phi_k(\mathbf{x}_j)$$

Covariance between Two Points

- ▶ The prior covariance between two points \mathbf{x}_i and \mathbf{x}_j is

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

or in sum notation

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \sum_{k=1}^m \phi_k(\mathbf{x}_i) \phi_k(\mathbf{x}_j)$$

- ▶ For the radial basis used this gives

Covariance between Two Points

- ▶ The prior covariance between two points \mathbf{x}_i and \mathbf{x}_j is

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \phi^T(\mathbf{x}_i) \phi(\mathbf{x}_j),$$

or in sum notation

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \sum_{k=1}^m \phi_k(\mathbf{x}_i) \phi_k(\mathbf{x}_j)$$

- ▶ For the radial basis used this gives

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha \sum_{k=1}^m \exp\left(-\frac{|\mathbf{x}_i - \boldsymbol{\mu}_k|^2 + |\mathbf{x}_j - \boldsymbol{\mu}_k|^2}{2\ell^2}\right).$$

Covariance Functions

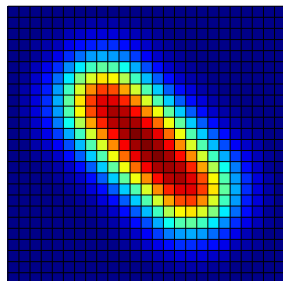
RBF Basis Functions

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

$$\phi_k(x) = \exp\left(-\frac{\|x - \mu_k\|_2^2}{\ell^2}\right)$$

$$\mu = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$$

x_1



x_2

Selecting Number and Location of Basis

- ▶ 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:

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

Selecting Number and Location of Basis

- ▶ 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:

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

Selecting Number and Location of Basis

- ▶ 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:

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

Selecting Number and Location of Basis

- ▶ 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:

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

Selecting Number and Location of Basis

- ▶ 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 $\Delta\mu$.

Infinite Basis Functions

- Take

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

Infinite Basis Functions

- ▶ Take

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

- ▶ This implies

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

Infinite Basis Functions

- Take

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

- This implies

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

and therefore

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

Infinite Basis Functions

- ▶ Take

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

- ▶ 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

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

- ▶ 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 $m \rightarrow \infty$

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

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

Result

- Performing the integration leads to

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

Result

- ▶ Performing the integration leads to

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

- ▶ Now take limit as $a \rightarrow -\infty$ and $b \rightarrow \infty$

Result

- ▶ Performing the integration leads to

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

- ▶ Now take limit as $a \rightarrow -\infty$ and $b \rightarrow \infty$

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

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

Infinite Feature Space

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

Infinite Feature Space

- ▶ 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 ??.

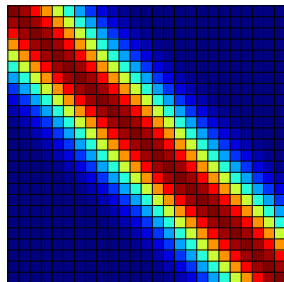
Covariance Functions

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 \mathbf{x} .
- ▶ For the example above it was based on Euclidean distance.
- ▶ The covariance function is also known as a kernel.



\mathbf{x}_1

\mathbf{x}_2

Covariance Functions

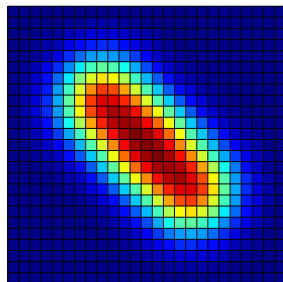
RBF Basis Functions

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

$$\phi_k(x) = \exp\left(-\frac{\|x - \mu_k\|_2^2}{\ell^2}\right)$$

$$\mu = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$$

x_1



x_2

References I

- P. S. Laplace. *Essai philosophique sur les probabilités*. Courcier, Paris, 2nd edition, 1814. Sixth edition of 1840 translated and reprinted (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.
- R. M. Neal. *Bayesian Learning for Neural Networks*. Springer, 1996. Lecture Notes in Statistics 118.
- C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, MA, 2006. [[Google Books](#)].
- C. K. I. Williams. Computation with infinite neural networks. *Neural Computation*, 10(5):1203–1216, 1998.