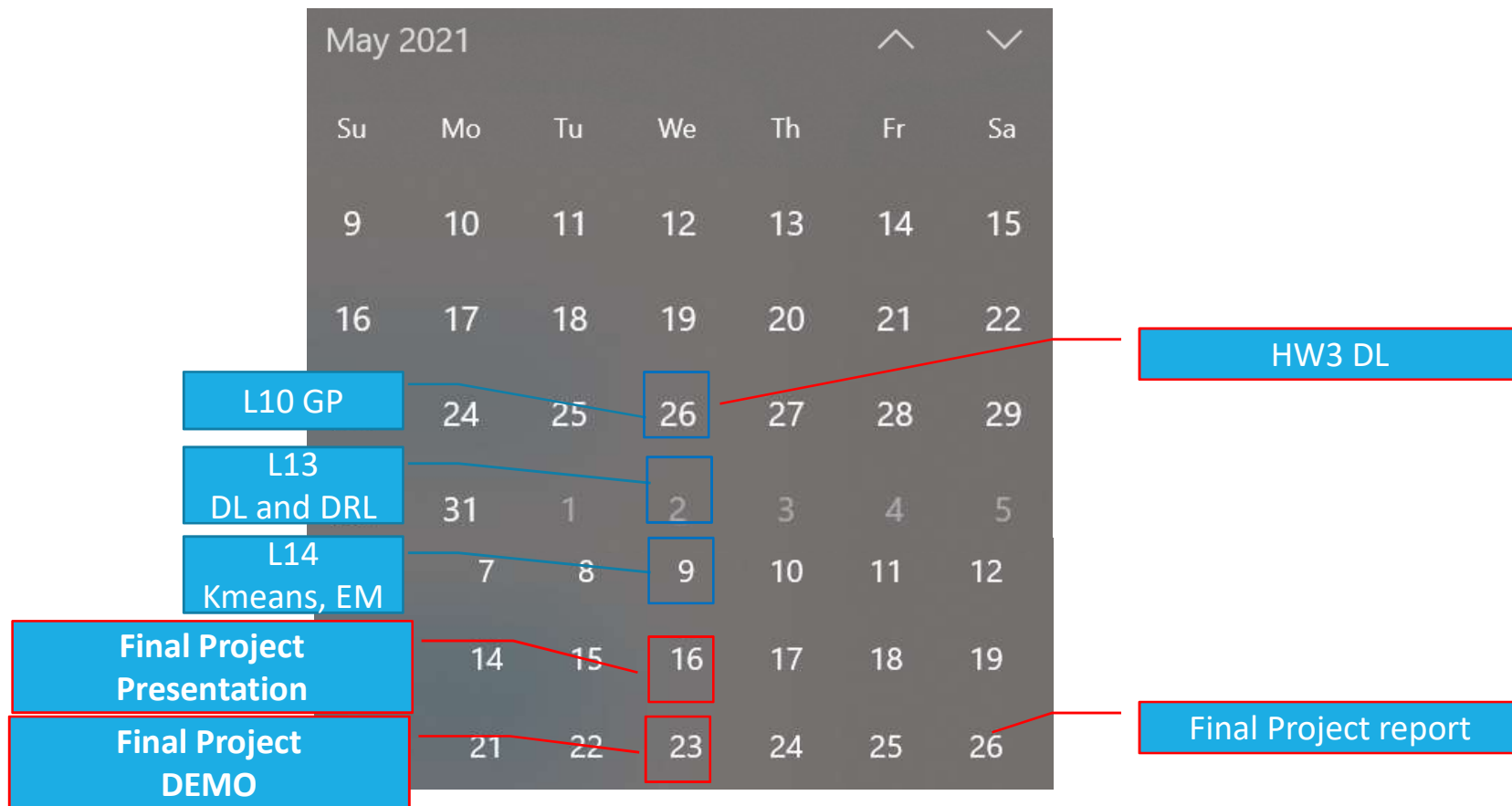# Gaussian Processes (GP) & Locally Weighted Project Regression(LWPR)

KUO-SHIH TSENG (曾國師)
DEPARTMENT OF MATHEMATICS
NATIONAL CENTRAL UNIVERSITY, TAIWAN

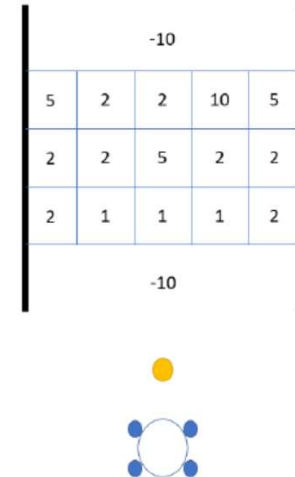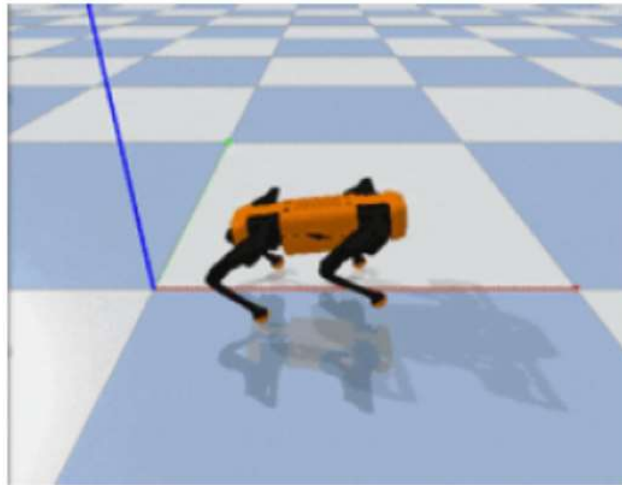2021/05/26

# Course Announcement



L10 GP

L13
DL and DRL

L14
Kmeans, EM

**Final Project
Presentation**

**Final Project
DEMO**

HW3 DL

Final Project report

# Course Announcement

- Final project proposal:
  - Check your proposal is
    - feasible?
    - Clear?

# Course Announcement

**M e s**                                    Sun, May 2, 5:33 PM  ☆  ↩  ⋮
to me ▾

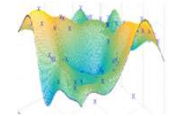文A  Chinese (Traditional) ▾  >  English ▾  Translate message    Turn off for: Chinese (Traditional)  ✕
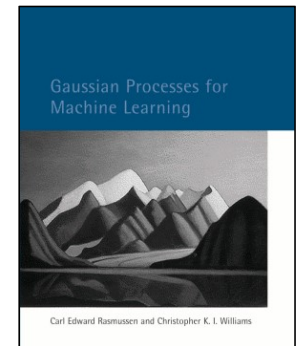
老師我聽朋友說做機器學習會一直用到高微和泛函，這是真的嗎XD 感覺上到現在都沒有看過。

Gaussian Processes: 我這不就來了嗎~

# Outline

- Need

- Gaussian Distribution

- MAP of Bayesian Linear Regression

- Kernel Functions

- Gaussian Processes (GP)
  ◦ EX: Coverage approximation

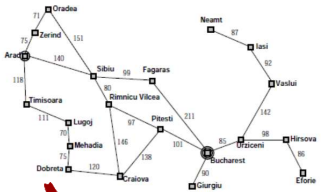- Locally Weighted Project Regression(LWPR)
  ◦ EX: Coverage approximation

[6] Carl Edward Rasmussen and Christopher K. I. Williams, " Gaussian Processes for Machine Learning," MIT Press, 2006.

# Outline

[Problem solving]

↓

Search problems
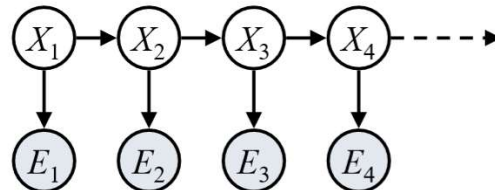
↓

Adversarial Search



[Perception and Uncertainty]

↓

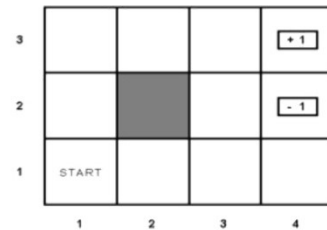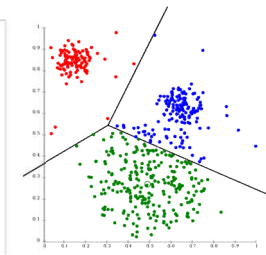Bayes Theorem

↓

Bayes Filter and Smoothing



[Learning and Decision-making]

| Supervised learning | Unsupervised learning | Reinforcement learning |



GP

# Need

- Parametric approximation needs good features. However, theses features are based on domain know-how. If we can learn the model <u>without features</u>, it could be applied to different domains.

$\{(\mathbf{x}_i, y_i)\}, i = 1 \sim N$

$y : measuremnt$

$\mathbf{x} : training\ sample$

$\boxed{f_i, i = 1 \sim m \\ f : features}$

Algorithm

$Q(\mathbf{x}) = \sum_i w_i f_i(\mathbf{x})$

Parametric approximation

$\boxed{\{(\mathbf{x}_i, y_i)\}, i = 1 \sim N \\ y : measuremnt \\ \mathbf{x} : training\ sample}$

Algorithm

$Q(\mathbf{x}) = ?$

Nonparametric approximation

# Need

- Function approximation for these cases:

*Given* :
$\{(\mathbf{x}_i, y_i)\}, i = 1 \sim 6$
$\mathbf{x}_7$
*Find* :
$y_7 = ? \operatorname{cov}(y_7) = ?$

*Given* :
$\{(\mathbf{x}_i, y_i)\}, i = 1 \sim N$
$\mathbf{x}_*$
*Find* :
$y_* = ? \operatorname{cov}(y_*) = ?$

No more feature information!

$P(\mathbf{y}_* \mid \mathbf{x}, \mathbf{y}, \mathbf{x}_*) = ?$

# Need

- Nonparametric approximation approaches include Gaussian processes (GP)[1], Locally Weighted Project Regression (LWPR) and neural network (NN).

- These approaches enable robots to learn some complex tasks[2].





[1] http://www.gaussianprocess.org/gpml/
[2] http://www-clmc.usc.edu/~sschaal/

# Gaussian Distribution

- Mean & variance → p(x)

- Probability density function (PDF)

$$p(x) = \frac{1}{\sqrt{2\pi}\,\sigma} \exp\left( -\frac{(x-\mu)^2}{2\sigma^2} \right)$$



https://en.wikipedia.org/wiki/Normal_distribution

# Gaussian Distribution

- Bivariate

$$\mathbf{X} = \begin{pmatrix} X \\ Y \end{pmatrix} \qquad X \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$p(\mathbf{x}) = \frac{1}{2\pi \|\boldsymbol{\Sigma}\|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \sigma^2_x & \sigma_{xy} \\ \sigma_{xy} & \sigma^2_y \end{pmatrix}$$



**Height**   **Weight**

# Gaussian Distribution

**Statistical Data Mining Tutorials by** **Andrew Moore**

# Gaussian Distribution



In this example, x and y are almost independent

# Gaussian Distribution



density values: 0.05 <= density < 0.11
density <= 0.05  0.11 < density

| | mean | cov | |
|---|---|---|---|
| x | -0.0579042 | 1.02654 | 1.06236 |
| x+y | -0.0885454 | 1.06236 | 2.0324 |

In this example, x and "x+y" are clearly not independent

# Gaussian Distribution



In this example, x and "20x+y" are clearly not independent

# Gaussian Distribution



The correlation coefficient $\rho_{X,Y}$ between two random variables $X$ and $Y$ with expected values $\mu_X$ and $\mu_Y$ and standard deviations $\sigma_X$ and $\sigma_Y$ is defined as:

$$\rho_{X,Y} = \frac{\mathrm{cov}(X,Y)}{\sigma_X \sigma_Y} = \frac{E((X - \mu_X)(Y - \mu_Y))}{\sigma_X \sigma_Y},$$

Source: WIKI

| Correlation | Negative | Positive |
|---|---|---|
| Small | −0.3 to −0.1 | 0.1 to 0.3 |
| Medium | −0.5 to −0.3 | 0.3 to 0.5 |
| Large | −1.0 to −0.5 | 0.5 to 1.0 |

# Gaussian Distribution

$$\rho_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y} = \frac{E((X - \mu_X)(Y - \mu_Y))}{\sigma_X \sigma_Y},$$

$$\rho_{X,Y} = \frac{\Sigma_{XY}}{\sigma_x \sigma_y}$$

$$= \frac{\Sigma_{XY}}{\Sigma_x^{0.5} \Sigma_y^{0.5}}$$

$$= \frac{20.5665}{\sqrt{1.02654}\sqrt{412.982}}$$

$$= 0.998$$

# Gaussian Distribution

## Error Ellipses in Action

**Kai Arras**
Social Robotics Lab, University of Freiburg

April 2010    Social Robotics Laboratory

# Gaussian Distribution

- Multivariate

$$\mathbf{X} = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_m \end{pmatrix} \qquad X \sim N(\mathbf{\mu}, \mathbf{\Sigma})$$

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{m/2} \|\mathbf{\Sigma}\|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}-\mathbf{\mu})^T \mathbf{\Sigma}^{-1}(\mathbf{x}-\mathbf{\mu})\right)$$

$$\mathbf{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_m \end{pmatrix} \qquad \mathbf{\Sigma} = \begin{pmatrix} \sigma^2_1 & \sigma_{12} & \cdots & \sigma_{1m} \\ \sigma_{12} & \sigma^2_2 & \cdots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1m} & \sigma_{2m} & \cdots & \sigma^2_m \end{pmatrix}$$

# Gaussian Distribution

- **Calculations of Gaussian**

- Summation

$$x \sim N(\mu_x, \Sigma_x)$$
$$y \sim N(\mu_y, \Sigma_y)$$

$\longrightarrow$

$$x + y \sim N(\mu_x + \mu_y, \Sigma_x + \Sigma_y)$$

- Multiplication

$$x \sim N(\mu_x, \Sigma_x)$$

$\longrightarrow$

$$bx \sim N(b\mu_x, b^T \Sigma_x b)$$

- Conditional Gaussian

$$x \sim N(\mu_x, \Sigma_{xc})$$
$$z \sim N(\mu_z, \Sigma_{zz})$$

$\longrightarrow$

$$\begin{cases} \mu_{x|z} = \mu_x + \Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z} - \mu_z) \\ \Sigma_{xx|z} = \Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx} \end{cases}$$

# Gaussian Distribution

- Conditional Gaussian

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}, P(\mathbf{x}) \sim N(\mu_x, \Sigma_{xx}) = \frac{1}{\sqrt{(2\pi)^n |\Sigma_{xx}|}} e^{-\frac{1}{2}(\mathbf{x}-\mu_x)^T \Sigma_{xx}^{-1}(\mathbf{x}-\mu_x)}$$

$$\mathbf{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix}, \mathbf{x} \text{ and } \mathbf{z} \text{ are jointly Gaussian, } P(\mathbf{z}) \sim N(\mu_z, \Sigma_{zz})$$

Reminder: We want to find $P(\mathbf{y}_* \mid \mathbf{x}, \mathbf{y}, \mathbf{x}_*) = ?$

$Find\ P(\mathbf{x}|\mathbf{z}) = ?$

$$\begin{cases} \mu_{x|z} = \mu_x + \Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z} - \mu_z) \\ \Sigma_{xx|z} = \Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx} \end{cases}$$

See the appendix for the proof.

Go to Appendix

# MAP of Bayesian Linear Regression

- Gaussian processes proof flow:

**1**

$$y = \mathbf{x}^T \mathbf{w} + \varepsilon, \quad \varepsilon \sim N(0, \sigma_n^2)$$

Projection

$y - \mathbf{x}^T \mathbf{w} = 0$

$y$

$x_2$

$(x_*, y_* = ?)$

**?**

$x_1$

**3**

$$y = \Phi(\mathbf{x})^T \mathbf{w} + \varepsilon, \quad \varepsilon \sim N(0, \sigma_n^2)$$

$y - \Phi(\mathbf{x})^T \mathbf{w} = 0$

$y$

$x'_2$

$(\Phi(x_*), y_* = ?)$

**?**

$x'_1$

**2** MAP :

$$P(\mathbf{w} \mid \mathbf{y}, X) \propto P(\mathbf{y} \mid X, \mathbf{w}) P(\mathbf{w})$$

Prediction :

$$P(y_* \mid x_*, y, X)$$

**4** MAP :

$$P(\mathbf{w} \mid \mathbf{y}, X) \propto P(\mathbf{y} \mid X, \mathbf{w}) P(\mathbf{w})$$

Prediction :

$$P(y_* \mid x_*, y, X)$$

# MAP of Bayesian Linear Regression

$\{(\mathbf{x}_i, y_i)\}, i = 1 \sim N$

$y : measuremnt$

$\mathbf{x} : traing\ sample$

$f_i, i = 1 \sim m$

$f : features$

Algorithm

$Q(\mathbf{x}) = \sum_i w_i f_i(\mathbf{x})$

Parametric approximation

*Training data D*
$D = \{(\boldsymbol{x_i}, y_i) | i = 1, \ldots N\} = \{X, \boldsymbol{y}\}$

$$y = \begin{bmatrix} x_{\{1,1\}} & \cdots & x_{\{1,m\}} \\ \vdots & \ddots & \vdots \\ x_{\{N,1\}} & \cdots & x_{\{N,m\}} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix} + \epsilon$$

$$y = \sum_i w_i f_i(\mathbf{x}) + \varepsilon, \quad \varepsilon \sim N(0, \sigma_n^2)$$

$$\Rightarrow y = \mathbf{x}^T \mathbf{w} + \varepsilon$$

Let's <u>take x vector</u> as **features**.
If we can find w vector, we got a linear approximation model!
We can use online/offline least square methods to find w.

# MAP of Bayesian Linear Regression

- Let's look at this problem from probabilistic perspective.

$$y = \mathbf{x}^T\mathbf{w} + \varepsilon, \quad \varepsilon \sim N(0, \sigma_n^2)$$

$$y = w_0 + w_1 x + \varepsilon, \quad \varepsilon \sim N(0, \sigma_n^2)$$

$$y\text{-}(w_0 + w_1 x) = \varepsilon$$

**Linear algebra perspective**

$$P\left[y\text{-}(w_0 + w_1 x)\right]$$

**Probabilistic perspective**

$$y = w_1 x + w_0$$

y

x

$$y\text{-}(w_0 + w_1 x)$$

$$-\sigma_n \; 0 \; \sigma_n$$

# MAP of Bayesian Linear Regression

- Let's look at this problem from probabilistic perspective.

$$y = \mathbf{x}^T \mathbf{w} + \varepsilon, \quad \varepsilon \sim N(0, \sigma_n^2)$$

$$y = w_0 + w_1 x_1 + w_2 x_2 + \varepsilon, \quad \varepsilon \sim N(0, \sigma_n^2)$$

$$y\text{-}(w_0 + w_1 x_1 + w_2 x_2) = \varepsilon$$

$$y = w_1 x + w_0$$

$$P[y - (w_0 + w_1 x_1 + w_2 x_2)]$$

$$y = (w_0 + w_1 x_1 + w_2 x_2)$$

# MAP of Bayesian Linear Regression

- Let's look at this problem in probabilistic domain and try to find **maximum a posterior** (MAP) estimation of w.

$$y = \mathbf{x}^T \mathbf{w} + \varepsilon, \quad \varepsilon \sim N(0, \sigma_n^2)$$

$$P(A \mid B) = \frac{P(B \mid A)P(A)}{P(B)}$$

$$P(\mathbf{w} \mid \mathbf{y}, X) = \frac{P(\mathbf{y}, X \mid \mathbf{w})P(\mathbf{w})}{P(\mathbf{y}, X)}$$

$$P(\mathbf{w} \mid \mathbf{y}, X) = \frac{P(\mathbf{y} \mid X, \mathbf{w})P(\mathbf{w})}{P(\mathbf{y} \mid X)}$$

$$= \frac{P(\mathbf{y} \mid X, \mathbf{w})P(X \mid \mathbf{w})P(\mathbf{w})}{P(\mathbf{y} \mid X)P(X)} \ldots \because P(X \mid \mathbf{w}) = P(X)$$

$$P(\mathbf{w} \mid \mathbf{y}, X) \propto P(\mathbf{y} \mid X, \mathbf{w})P(\mathbf{w})$$

$$= \frac{P(\mathbf{y} \mid X, \mathbf{w})P(\mathbf{w})}{P(\mathbf{y} \mid X)}$$

Let's find w vector with **MAP** estimation.

$\mathbf{x}^T : x$ vector/matrix

$X : $ input data

$\mathbf{y} : output\ data$

# MAP of Bayesian Linear Regression

$$P(\mathbf{w} \mid y, X) \propto P(y \mid X, \mathbf{w})P(\mathbf{w})$$

$$\mathbf{y} = \mathbf{x}^T\mathbf{w} + \varepsilon, \quad \varepsilon \sim N(0, \sigma_n^2)$$

$$P(\mathbf{y} \mid X, \mathbf{w})$$

$$= \prod_{i=1}^{n} p(y_i \mid \mathbf{x}_i, \mathbf{w}) \quad \because \text{iid}$$

$$= \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{(y_i - \mathbf{x}_i^T\mathbf{w})^2}{2\sigma_n^2}\right)$$

$$= \frac{1}{(2\pi\sigma_n)^{n/2}} \exp\left(-\frac{1}{2\sigma_n^2}\left|\mathbf{y} - X^T\mathbf{w}\right|^2\right)$$

$$P(\mathbf{y} \mid X, \mathbf{w}) \sim N(\quad X^T\mathbf{w}, \sigma_n^2)$$

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{m/2} \|\mathbf{\Sigma}\|^{1/2}} \exp\left(-\tfrac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

Assume $P(\mathbf{w}) \sim N(0, \Sigma_p)$

$$P(\mathbf{w}) = \frac{1}{(2\pi)^{n/2}\|\Sigma_p\|^{1/2}} \exp\left(-\frac{1}{2}\mathbf{w}^T\Sigma_p^{-1}\mathbf{w}\right)$$

**Weighting space perspective**

y

x

The slope (W)
has uncertainty

# MAP of Bayesian Linear Regression

$$P(\mathbf{y} \mid X, \mathbf{w}) \sim N(\quad X^T\mathbf{w}, \sigma_n^2)$$

$$\underline{P(\mathbf{y} \mid X, \mathbf{w})}$$

$$= \frac{1}{(2\pi\sigma)^{n/2}} \exp\left(-\frac{1}{2\sigma_n^2}\left|\mathbf{y} - X^T\mathbf{w}\right|^2\right)$$

$$P(\mathbf{w}) \sim N(0, \Sigma_p)$$

$$\underline{P(\mathbf{w})} = \frac{1}{(2\pi)^{n/2}\left\|\Sigma_p\right\|^{1/2}} \exp\left(-\frac{1}{2}\mathbf{w}^T \Sigma_p^{-1} \mathbf{w}\right)$$

$$P(\mathbf{w} \mid \mathbf{y}, X) \propto \underline{P(\mathbf{y} \mid X, \mathbf{w})}\,\underline{P(\mathbf{w})}$$

$$\propto \frac{1}{(2\pi\sigma)^{n/2}} \exp\left(-\frac{1}{2\sigma_n^2}\left|\mathbf{y} - X^T\mathbf{w}\right|^2\right) \frac{1}{(2\pi)^{n/2}\left\|\Sigma_p\right\|^{1/2}} \exp\left(-\frac{1}{2}\mathbf{w}^T \Sigma_p^{-1} \mathbf{w}\right)$$

$$\propto \exp\left(-\frac{1}{2\sigma_n^2}\left|\mathbf{y} - X^T\mathbf{w}\right|^2\right) \exp\left(-\frac{1}{2}\mathbf{w}^T \Sigma_p^{-1} \mathbf{w}\right)$$

# MAP of Bayesian Linear Regression

$$P(\mathbf{w} \mid \mathbf{y}, X) \propto P(\mathbf{y} \mid X, \mathbf{w}) P(\mathbf{w})$$

$$\propto \exp\left(-\frac{1}{2\sigma_n^2}\left|\mathbf{y} - X^T\mathbf{w}\right|^2\right) \exp\left(-\frac{1}{2}\mathbf{w}^T \Sigma_p^{-1} \mathbf{w}\right)$$

$$\propto \exp\left(-\frac{1}{2}(\mathbf{w} - \overline{\mathbf{w}})^T \left(\frac{1}{\sigma_n^2} XX^T + \Sigma_p^{-1}\right)(\mathbf{w} - \overline{\mathbf{w}})\right)$$

where $\overline{\mathbf{w}} = \sigma_n^{-2}\left(\sigma_n^{-2} XX^T + \Sigma_p^{-1}\right)^{-1} X\mathbf{y}$

$$P(\mathbf{w} \mid \mathbf{y}, X) \sim N\left(\overline{\mathbf{w}} = \frac{1}{\sigma_n^2} A^{-1} X\mathbf{y}, A^{-1}\right)$$

where $A = \sigma_n^{-2} XX^T + \Sigma_p^{-1}$
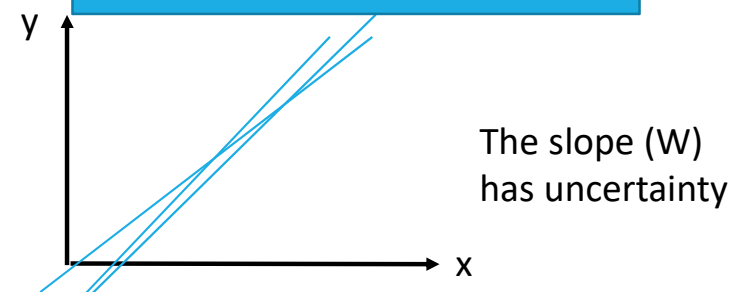
$$p(\mathbf{x}) = \frac{1}{(2\pi)^{m/2} \|\Sigma\|^{1/2}} \exp\left(-\tfrac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

Expand the summation term to Gaussian form.

$$P(\mathbf{y} \mid X, \mathbf{w}) \sim N(\quad X^T\mathbf{w}, \sigma_n^2)$$

$$P(\mathbf{w}) \sim N(0, \Sigma_p)$$

We got a new distribution and can use it to predict yi.

# MAP of Bayesian Linear Regression

- Prediction based on a Gaussian distribution

$$P(\mathbf{w} \mid y, X) \sim N\left( \overline{\mathbf{w}} = \frac{1}{\sigma_n^2} A^{-1} X\mathbf{y}, A^{-1} \right)$$

**Prediction**

(Total probability & independence)   (Scaling a Gaussian to another Gaussian)

$$P(y_* \mid x_*, y, X) = \int P(y_* \mid x_*, \mathbf{w}) P(\mathbf{w} \mid y, X) d\mathbf{w} = \int x_*^T \cdots P(\mathbf{w} \mid y, X) d\mathbf{w}$$

$$= N\left( \frac{1}{\sigma_n^2} x_*^T A^{-1} X\mathbf{y}, x_*^T A^{-1} x_* \right)$$

$$x \sim N(\mu_x, \Sigma_x)$$
$$bx \sim N(b\mu_x, b^T \Sigma_x b)$$

where $A = \sigma_n^{-2} XX^T + \Sigma_p^{-1}$

$$\begin{cases} \overline{y}_* = \dfrac{1}{\sigma_n^2} x_*^T A^{-1} X\mathbf{y} \\ \Sigma_{y_*} = x_*^T A^{-1} x_* \end{cases}$$

Where is **W**?

$$\begin{cases} \mu_{x|z} = \mu_x + \Sigma_{xz} \Sigma_{zz}^{-1} \left( \mathbf{z} - \mu_z \right) \\ \Sigma_{xx|z} = \Sigma_{xx} - \Sigma_{xz} \Sigma_{zz}^{-1} \Sigma_{zx} \end{cases}$$

# MAP of Bayesian Linear Regression

- Prediction based on a Gaussian distribution $P(\mathbf{w} \mid \mathbf{y}, X) \propto P(\mathbf{y} \mid X, \mathbf{w}) P(\mathbf{w})$

$P(w) \sim N(0, I)$

$y = w_1 + w_2 x + \epsilon$
$P(y_* \mid x_*. X, y)$

$P(y \mid X, w)$
Likelihood

$P(w \mid y, X)$
Posterior

# MAP of Bayesian Linear Regression

- Prediction based on a Gaussian distribution

$Given:$

$\{(\mathbf{x}_i, y_i)\}, i = 1 \sim 6$

$\mathbf{x}_7$

$Find:$

$y_7 = ? \operatorname{cov}(y_7) = ?$

$$\begin{cases} \overline{y}_* = \dfrac{1}{\sigma_n^2} x_*^T A^{-1} X \mathbf{y} \\ \sum_{y_*} = x_*^T A^{-1} x_* \end{cases}$$
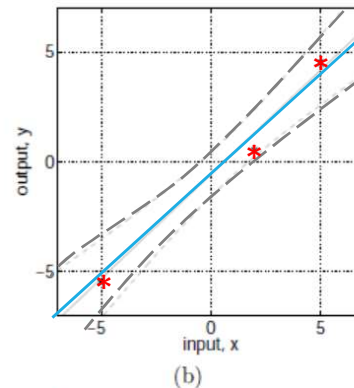
We can use it to predict yi based on the training data!

No!
It only works for linear data.

It works well?

Kuo-Shih

Math students

# Kernel functions

- The prediction of linear model could not work very well.

$$y = \mathbf{x}^T \mathbf{w} + \varepsilon, \quad \varepsilon \sim N(0, \sigma_n^2)$$

- In this model, it's a special case of feature base approximation. We adopt "X" as features.

We can project X to feature space

Why do you teach us this?

Kuo-Shih

Math students

# Kernel functions

- Let's project data (X) to feature space.

$$y = \mathbf{x}^T \mathbf{w} + \varepsilon, \quad \varepsilon \sim N(0, \sigma_n^2)$$

$$y - (w_0 + w_1 x_1 + w_2 x_2) = \varepsilon$$

$$y = \Phi(\mathbf{x})^T \mathbf{w} + \varepsilon, \quad \varepsilon \sim N(0, \sigma_n^2)$$

$$y - \mathbf{x}^T \mathbf{w} = 0$$

$y$

$x_2$

$x_1$

$$y - \Phi(\mathbf{x})^T \mathbf{w} = 0$$

$y$

$x'_2$

$x'_1$

What's the difference
Between GP and NN?

$\Phi(X)^T \mathbf{W}$

$\Phi(X^T \mathbf{W})$

# Kernel functions

- Let's project data (X) to feature space.

$$\mathbf{x}$$
$D \times 1$ vector

$$\Phi(.)$$

$$\Phi(\mathbf{x})$$
$N \times 1$ vector

$$y = \mathbf{x}^T \mathbf{w} + \varepsilon$$

$$y = \Phi(\mathbf{x})^T \mathbf{w} + \varepsilon$$

This is still a feature based regression!!

If $\Phi(x) = \left(1, x, x^2, x^3, ....\right)^T$,

it's a polynomial regression!

Math students

# Kernel functions

- Let's project data (X) to feature space.

$$\mathbf{x} \xrightarrow{\quad} \boxed{\Phi(.)} \xrightarrow{\quad} \Phi(\mathbf{x})$$

$D \times 1$ vector $\qquad\qquad\qquad\qquad N \times 1$ vector

$$y = \mathbf{x}^T \mathbf{w} + \varepsilon \qquad\qquad\qquad y = \Phi(\mathbf{x})^T \mathbf{w} + \varepsilon$$

$$P(y_* \mid x_*, y, X) = N\left( \frac{1}{\sigma_n^2} x_*^T A^{-1} X\mathbf{y}, x_*^T A^{-1} x_* \right) \qquad P(y_* \mid x_*, y, X) = N\left( \frac{1}{\sigma_n^2} \Phi(x_*)^T A^{-1} X\mathbf{y}, \Phi(x_*)^T A^{-1} \Phi(x_*) \right)$$

where $A = \sigma_n^{-2} XX^T + \sum_p^{-1}$ $\qquad\qquad\qquad$ where $A = \sigma_n^{-2} \Phi\Phi^T + \sum_p^{-1}$

$$\begin{cases} \bar{y}_* = \dfrac{1}{\sigma_n^2} x_*^T A^{-1} X\mathbf{y} \\ \sum_{y_*} = x_*^T A^{-1} x_* \end{cases} \qquad\qquad \begin{cases} \bar{y}_* = \dfrac{1}{\sigma_n^2} \Phi(x_*)^T A^{-1} \Phi\mathbf{y} \\ \sum_{y_*} = \Phi(x_*)^T A^{-1} \Phi(x_*) \end{cases}$$

# Kernel functions

$$P(y_* \mid x_*, y, X) = N\left(\frac{1}{\sigma_n^2} \Phi(x_*)^T A^{-1} X\mathbf{y}, \Phi(x_*)^T A^{-1} \Phi(x_*)\right)$$

where $A = \sigma_n^{-2}\Phi\Phi^T + \Sigma_p^{-1} = \sigma_n^{-2}K + \Sigma_p^{-1}$

$$\left\{ \begin{array}{l} \bar{y}_* = \dfrac{1}{\sigma_n^2} \Phi(x_*)^T A^{-1}\Phi\mathbf{y} \\[2mm] \Sigma_{y_*} = \Phi(x_*)^T A^{-1}\Phi(x_*) \end{array} \right.$$

Let $Z^{-1} = \Sigma_p, W^{-1} = \sigma_n^2 I, U = V = \Phi$

$$A^{-1} = \left(\sigma_n^{-2}\Phi\Phi^T + \Sigma_p^{-1}\right)^{-1}$$

$$= \Sigma_p - \Sigma_p \Phi\left(\sigma_n^2 I + \Phi\Sigma_p\Phi^T\right)\Phi^T \Sigma_p$$

$$\Sigma_{y_*} = \Phi(x_*)^T A^{-1}\Phi(x_*) = \Phi_*^T \Sigma_p \Phi_* - \Phi_*^T \Sigma_p \Phi \left(K + \sigma_n^2 I\right)^{-1}\Phi^T \Sigma_p \Phi_*$$

Matrix inversion lemma :
$$\left(Z + UWV^T\right)^{-1}$$
$$= Z^{-1} - Z^{-1}U\left(W^{-1} + V^T Z^{-1} U\right)^{-1} V^T Z^{-1}$$

$$\left\{ \begin{array}{l} \mu_{x|z} = \mu_x + \Sigma_{xz}\Sigma_{zz}^{-1}\left(\mathbf{z} - \mu_z\right) \\[2mm] \Sigma_{xx|z} = \Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx} \end{array} \right.$$

# Kernel functions

$$\Sigma_{y_*} = \Phi_*^T \Sigma_p \Phi_* - \Phi_*^T \Sigma_p \Phi \left(K + \sigma_n^2 I\right)^{-1} \Phi^T \Sigma_p \Phi_*$$

$$\mu_x = \mathbf{0}, \mu_z = \mathbf{0}$$

$$\Sigma_{xz} = \Phi_*^T \Sigma_p \Phi$$

$$\Sigma_{zz} = K + \sigma_n^2 I$$

$$\therefore \bar{y}_* = \Phi_*^T \Sigma_p \Phi \left(K + \sigma_n^2 I\right)^{-1} \mathbf{y}$$

$$\begin{cases} \mu_{x|z} = \mu_x + \Sigma_{xz} \Sigma_{zz}^{-1} \left(\mathbf{z} - \mu_z\right) \\ \Sigma_{xx|z} = \Sigma_{xx} - \Sigma_{xz} \Sigma_{zz}^{-1} \Sigma_{zx} \end{cases}$$

The Gaussian distribution is as follows:

$$\begin{cases} \bar{y}_* = \Phi_*^T \Sigma_p \Phi \left(K + \sigma_n^2 I\right)^{-1} \mathbf{y} \\ \Sigma_{y_*} = \Phi_*^T \Sigma_p \Phi_* - \Phi_*^T \Sigma_p \Phi \left(K + \sigma_n^2 I\right)^{-1} \Phi^T \Sigma_p \Phi_* \end{cases}$$

# Kernel functions

- Let's rewrite it.

$$y = \Phi(\mathbf{x})^T \mathbf{w} + \varepsilon$$

$$P(y_* \mid x_*, y, X) = N\left(\Phi_*^T \Sigma_p \Phi \left(K + \sigma_n^2 I\right)^{-1} \mathbf{y}, \Phi_*^T \Sigma_p \Phi_* - \Phi_*^T \Sigma_p \Phi \left(K + \sigma_n^2 I\right)^{-1} \Phi^T \Sigma_p \Phi_*\right)$$

where $K = \Phi^T \Sigma_p \Phi$

let's define $k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Sigma_p \Phi(\mathbf{x}')$

$$k_* = k(\mathbf{x}, \mathbf{x}_*)$$

$$K = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \cdots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \cdots & k(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_n, x_1) & k(x_n, x_2) & \cdots & k(x_n, x_n) \end{bmatrix}$$

$$k_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$$

$$P(y_* \mid x_*, y, X) = N\left(k_*^T \left(K + \sigma_n^2 I\right)^{-1} \mathbf{y}, k_{**} - k_*^T \left(K + \sigma_n^2 I\right)^{-1} k_*\right)$$

$$\begin{cases} \mu_{x\mid z} = \mu_x + \Sigma_{xz} \Sigma_{zz}^{-1} (\mathbf{z} - \mu_z) \\ \Sigma_{xx\mid z} = \Sigma_{xx} - \Sigma_{xz} \Sigma_{zz}^{-1} \Sigma_{zx} \end{cases}$$

We call "K" kernel function. It can project data to feature space.
How to choose a good kernel?

# Kernel functions



Figure 1.1: Examples of structures expressible by some basic kernels.

https://github.com/duvenaud/phd-thesis

# Kernel functions

- The squared exponential kernel

$$k(x, x') = \sigma_f^2 \exp\left( -\frac{1}{2l^2}(x - x')^2 \right)$$

We call them "Hyperparameters"

You use two parameters!!

Kuo-Shih

Math students

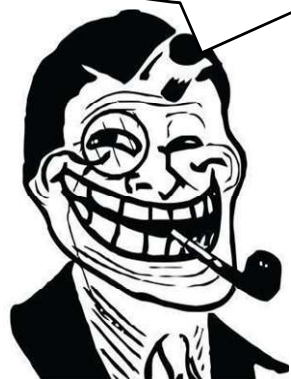# Kernel functions

- The squared exponential kernel

$$k(x, x') = \sigma_f^2 \exp\left(-\frac{1}{2l^2}(x - x')^2\right)$$

Hyperparameters can be learned from data (X)!

What's the difference?

Kuo-Shih

Math students

# Gaussian Processes (GP)

Definition : A Gaussian process is a collection of random variables **(Function space view)**

with a joint Gaussian distribution.

$$y = \Phi(\mathbf{x})^T \mathbf{w} + \varepsilon, \ \ \varepsilon \sim N(0, \sigma_n^2)$$

let's define $f(\mathbf{x}) = \Phi(\mathbf{x})^T \mathbf{w}, \mathbf{w} \sim N(0, \sum_P)$

Mean function : $m(\mathbf{x}) = E[f(\mathbf{x})]$

Covariance function : $k(\mathbf{x}, \mathbf{x}') = E[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x})')]$

$$\Rightarrow f(\mathbf{x}) \sim GP(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

$$m(\mathbf{x}) = E[f(\mathbf{x})] = \Phi(\mathbf{x})^T E[\mathbf{w}] = 0$$

Assume $P(\mathbf{w}) \sim N(0, \sum_p)$

$$E[f(\mathbf{x})f(\mathbf{x}')] = \Phi(\mathbf{x})^T E[\mathbf{w}\mathbf{w}^T]\Phi(\mathbf{x}') = \Phi(\mathbf{x})^T \sum_p \Phi(\mathbf{x}')$$

$$\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \sim N\left( 0, \begin{bmatrix} K(X,X) + \sigma_n^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix} \right)$$

# Gaussian Processes (GP)



**Linear algebra perspective**

$$y = w_1 x + w_0$$

$$P[y\text{-}(w_0 + w_1 x)]$$

**Probabilistic perspective**

$$y - \Phi(\mathbf{x})^T \mathbf{w} = 0$$

$y\text{-}(w_0 + w_1 x)$

— : functions

▨ : 95% bound

**+** : sampling data

(a), prior

(b), posterior

**Function space perspective**

$$P(\mathbf{w} \mid \mathbf{y}, X) \propto P(\mathbf{y} \mid X, \mathbf{w}) P(\mathbf{w})$$

$$y = \Psi(\boldsymbol{x})^T \boldsymbol{w} + \epsilon$$
$$\blacktriangleright y = f(\boldsymbol{x}) + \epsilon$$

# Gaussian Processes (GP)

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}, P(\mathbf{x}) \sim N(\mu_x, \Sigma_{xx}) = \frac{1}{\sqrt{(2\pi)^n |\Sigma_{xx}|}} e^{-\frac{1}{2}(\mathbf{x}-\mu_x)^T \Sigma_{xx}^{-1}(\mathbf{x}-\mu_x)}$$

$$\mathbf{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix}, \mathbf{x} \text{ and } \mathbf{z} \text{ are jointly Gaussian, } P(\mathbf{z}) \sim N(\mu_z, \Sigma_{zz})$$

*Find* $P(\mathbf{x}|\mathbf{z}) = ?$

$$\begin{cases} \mu_{x|z} = \mu_x + \Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z} - \mu_z) \\ \Sigma_{xx|z} = \Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx} \end{cases}$$

# Gaussian Processes (GP)

$$\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \sim N\left( \mathbf{0}, \begin{bmatrix} K(X,X)+\sigma_n^2 I & K(X,X_*) \\ K(X_*,X) & K(X_*,X_*) \end{bmatrix} \right)$$
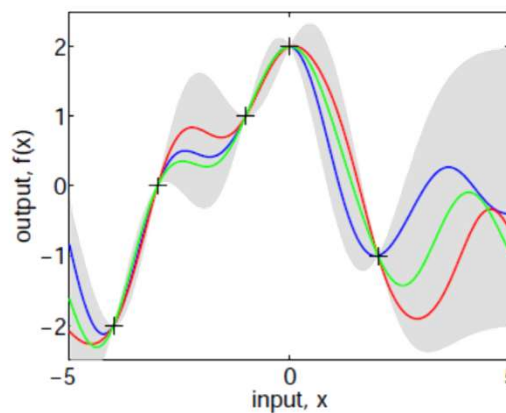
$$\begin{cases} \mu_{x|z} = \mu_x + \Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z}-\mu_z) \\ \Sigma_{xx|z} = \Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx} \end{cases}$$

$$\begin{cases} \mu_{y_*|y} = \mu_{y_*} + \Sigma_{y_*y}\Sigma_{yy}^{-1}(\mathbf{y}-\mathbf{0}) = K(X_*,X)\left(K(X,X)+\sigma_n^2 I\right)^{-1}\mathbf{y} \\ \Sigma_{y_*y_*|y} = \Sigma_{y_*y_*} - \Sigma_{y_*y}\Sigma_{yy}^{-1}\Sigma_{yy_*} = K(X_*,X_*) - K(X_*,X)\left(K(X,X)+\sigma_n^2 I\right)^{-1}K(X,X_*) \end{cases}$$

$$\begin{cases} \bar{y}_* = K_*^T\left(K+\sigma_n^2 I\right)^{-1}\mathbf{y} \\ \Sigma_{y_*} = K_{**} - K_*^T\left(K+\sigma_n^2 I\right)^{-1}K_* \end{cases}$$

# Gaussian Processes (GP)

- Prediction algorithm

$[Input]$

$X$ : data input, $\sigma_n$ : $noise$, $\mathbf{y}$ : data label

$k$ : covariance function

$X_*$ : test input

$1. \bar{y}_* = K_*^T \left( K + \sigma_n^2 I \right)^{-1} \mathbf{y}$

$2. \Sigma_{y_*} = K_{**} - K_*^T \left( K + \sigma_n^2 I \right)^{-1} K_*$

Return

$$k(x, x') = \sigma_f^2 \exp\left( -\frac{1}{2l^2} (x - x')^2 \right)$$

[Training]

$O\left( N^3 \right)$     Matrix invresion

[Prediction]

$O(N^2)$

# Gaussian Processes (GP)
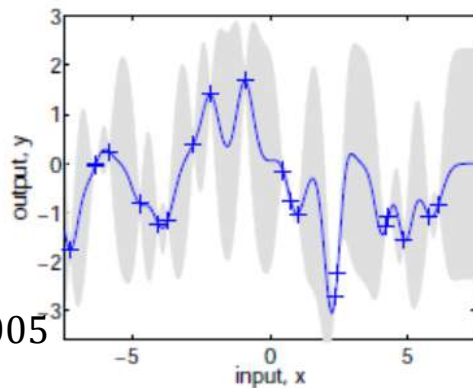


(a), $\ell = 1$
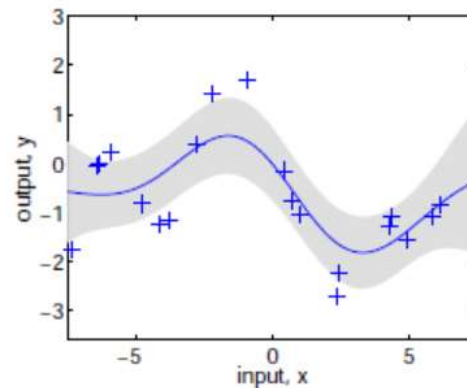
$l = 1$
$\sigma_f = 1$
$\sigma_n = 0.1$

$$k(x, x') = \sigma_f^2 \exp\left( -\frac{1}{2l^2}(x - x')^2 \right)$$

(b), $\ell = 0.3$

$l = 0.3$
$\sigma_f = 1.08$
$\sigma_n = 0.00005$

(c), $\ell = 3$

$l = 3$
$\sigma_f = 1.16$
$\sigma_n = 0.89$

How to tune **hyperparameters**?

# Gaussian Processes (GP)

- Learning Hyperparameters

$$\log P(\mathbf{y} \mid X, \mathbf{w})$$

$$= \frac{1}{2}\mathbf{y}^T K^{-1}\mathbf{y} - \frac{1}{2}\log|K| - \frac{n}{2}\log 2\pi$$

$$\frac{\partial \log P(\mathbf{y} \mid X, \mathbf{w})}{\partial w_j}$$

$$= \frac{1}{2}\mathbf{y}^T K^{-1}\frac{\partial K}{\partial w_j}\mathbf{y} - \frac{1}{2}tr\left(K^{-1}\frac{\partial K}{\partial w_j}\right)$$

$$= \frac{1}{2}tr\left((K^{-1}\mathbf{y}\mathbf{y}^T K^{-1^T} - K^{-1})\frac{\partial K}{\partial w_j}\right)$$

$$w_j = w_j + \alpha\frac{\partial \log P(\mathbf{y} \mid X, \mathbf{w})}{\partial w_j}$$

$$k(x, x') = \sigma_f^2 \exp\left(-\frac{1}{2l^2}(x - x')^2\right)$$

$$P(\boldsymbol{y}|\mathrm{X}, \boldsymbol{W}) = \frac{1}{2\pi\sigma^{n/2}}\exp(-\frac{1}{2\sigma_n^2}|y - \mathrm{X}^{\mathrm{T}}\boldsymbol{w}|^2)$$

$$P(\boldsymbol{y}|\mathrm{X}, \boldsymbol{W}) = \frac{1}{2\pi\sigma^{n/2}}\exp(-\frac{1}{2\sigma_n^2}|y - \Phi(\mathrm{X})^{\mathrm{T}}\boldsymbol{w}|^2)$$

# Gaussian Processes (GP)

- GP Library:

- scikit-learn (Python)
  ◦ https://scikit-learn.org/stable/modules/gaussian_process.html

- SheffieldML (C++)
  ◦ https://github.com/SheffieldML/GPc

- GP ML (MATLAB)
  ◦ http://www.gaussianprocess.org/gpml/code/matlab/doc/

- Gpstuff (MATLAB)
  ◦ https://github.com/gpstuff-dev/gpstuff

# GP — EX:

*Given* :

$\{(\mathbf{x}_i, y_i)\}, i = 1 \sim 6$

*Find* :

$y_7 = ? \; \mathrm{cov}(y_7) = ?$



Example source: Mark Ebden, "Gaussian Processes: A Quick Introduction," arXiv, 2015

# GP — EX: Coverage Approximation

- 1D case (X: robot's Y position)



Map1



Map2

# GP — EX: Coverage Approximation

- 2D case (X1: robot1's Y position, X2: robot2's Y position, )

# Gaussian Processes (GP)

- ☺
  ◦ Don't need hand crafting (features)
  ◦ Provide uncertainty information (mean and covariance)
  ◦ Easy to predict (O(N^2))

- ☹
  ◦ Need to select kernels
  ◦ Computational complexity of training: O(N^3), <15,000 samples

# Locally Weighted Project Regression(LWPR)

- Fastest & Scalable
  ◦ O(N^2) for training

- Input space is high-dimension, data lies on low-dimension manifold.



Schaal, S.;Atkeson, C. G.;Vijayakumar, S.,"Scalable techniques from nonparameteric statistics for real-time robot learning," Applied Intelligence, 2002

# Locally Weighted Project Regression(LWPR)

Least square(LS):

$$y = \beta X$$

$$\beta = \left(X^T X\right)^{-1} X^T y$$

Weighted LS(WLS):

$$y = \beta X$$

$$w_{ii} = \frac{1}{\sigma^2}$$

$$\beta = \left(X^T W X\right)^{-1} X^T W y$$

Locally Weighted Regression(LWR):

$$y = \beta X$$

$$w_{ii} = \exp\left(-\frac{1}{2}\left(x_i - x_q\right)^T D\left(x_i - x_q\right)\right)$$

$$\beta = \left(X^T W X\right)^{-1} X^T W y$$

LWPR=LWR + PLS (partial least square)

# LWPR – EX: CTG

Random 50 points → **LWPR** →

The fitted function: nMSE=1.524

1

The true function

The fitted function: nMSE=0.254

10

The fitted function: nMSE=0.041

20

http://wcms.inf.ed.ac.uk/ipab/slmc/research/software-lwpr

# Conclusions

- Nonparametric approximation :

- GP generates the most accurate approximation but its complexity is O(N^3).

- LWPR is the fastest method.

- If you have a few data (N<15,000), use GP.

- If you want to be fast, use LWPR.

- If you want to make a trade-off, use LGP, a mixture of GP and LWPR.

- After 2012, researchers start to adapt deep learning.

- During 2018~2020, researchers indicate that GP outperforms DL in many cases.

Nguyen-Tuong, D.; Seeger, M.; Peters, "*Model Learning with Local Gaussian Process Regression*, Advanced Robotics," 2009.

# Conditional Gaussian

$\mathbf{x}^{MMSE}$ in Gaussian distribution.

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}, P(\mathbf{x}) \sim N(\mu_x, \Sigma_{xx}) = \frac{1}{\sqrt{(2\pi)^n |\Sigma_{xx}|}} e^{-\frac{1}{2}(\mathbf{x}-\mu_x)^T \Sigma_{xx}^{-1}(\mathbf{x}-\mu_x)}$$

$$\mathbf{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix}, \mathbf{x} \text{ and } \mathbf{z} \text{ are jointly Gaussian}, P(\mathbf{z}) \sim N(\mu_z, \Sigma_{zz})$$

(multivariate normal distribution)

$Find\ P(\mathbf{x}|\mathbf{z}) = ?$

# Conditional Gaussian

$$P(\mathbf{x}) \sim N(\mu_x, \Sigma_{xx}) \text{ and } P(\mathbf{z}) \sim N(\mu_z, \Sigma_{zz}), P(\mathbf{x}, \mathbf{z}) = P(\mathbf{y}) = N(\mu_y, \Sigma_{yy})$$

$$\text{Mean}: \mathbf{y} = \begin{bmatrix} \mathbf{x}_{(m,1)} \\ \mathbf{z}_{(n,1)} \end{bmatrix}$$

$$\mu_y = \mathbf{E}[\mathbf{y}] = \mathbf{E}\begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix} = \begin{bmatrix} \mathbf{E}[\mathbf{x}] \\ \mathbf{E}[\mathbf{z}] \end{bmatrix} = \begin{bmatrix} \mu_x \\ \mu_z \end{bmatrix}$$

$$\text{Covariance}: \Sigma_{yy} = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xz} \\ \Sigma_{zx} & \Sigma_{zz} \end{bmatrix}$$

$$\Sigma_{xx} = \mathbf{E}\left[(\mathbf{x} - \mu_x)(\mathbf{x} - \mu_x)^T\right], \Sigma_{zz} = \mathbf{E}\left[(\mathbf{z} - \mu_z)(\mathbf{z} - \mu_z)^T\right]$$

$$\Sigma_{xz} = \mathbf{E}\left[(\mathbf{x} - \mu_x)^T(\mathbf{z} - \mu_z)^T\right], \Sigma_{zx} = \Sigma_{xz}^T$$

# Conditional Gaussian

$$P(\mathbf{x}|\mathbf{z}) = \frac{P(\mathbf{x},\mathbf{z})}{P(\mathbf{z})} = \frac{\dfrac{1}{\sqrt{(2\pi)^{m+n}|\Sigma_{yy}|}}e^{-\frac{1}{2}(\mathbf{y}-\mu_y)^T\Sigma_{yy}^{-1}(\mathbf{y}-\mu_y)}}{\dfrac{1}{\sqrt{(2\pi)^{m}|\Sigma_{zz}|}}e^{-\frac{1}{2}(\mathbf{z}-\mu_z)^T\Sigma_{zz}^{-1}(\mathbf{z}-\mu_z)}}$$

$$= \frac{1}{\sqrt{(2\pi)^{n}|\Sigma_{yy}|/|\Sigma_{zz}|}}e^{-\frac{1}{2}\left[(\mathbf{y}-\mu_y)^T\Sigma_{yy}^{-1}(\mathbf{y}-\mu_y)-(\mathbf{z}-\mu_z)^T\Sigma_{zz}^{-1}(\mathbf{z}-\mu_z)\right]}$$

$$= \frac{1}{\sqrt{(2\pi)^{n}\beta}}e^{-\frac{1}{2}[\alpha]}$$

If we can find $\alpha$ and $\beta$, we know the mean and covariance.

# Conditional Gaussian

$$\alpha = \left(\mathbf{y} - \mu_y\right)^T \Sigma_{yy}^{-1}\left(\mathbf{y} - \mu_y\right) - \left(\mathbf{z} - \mu_z\right)^T \Sigma_{zz}^{-1}\left(\mathbf{z} - \mu_z\right)$$

$$\begin{cases} \widetilde{\mathbf{x}} = \mathbf{x} - \mu_x \\ \widetilde{\mathbf{z}} = \mathbf{z} - \mu_z \end{cases}$$

$$= \begin{bmatrix} \widetilde{\mathbf{x}} \\ \widetilde{\mathbf{z}} \end{bmatrix}^T \begin{bmatrix} \Sigma_{xx} & \Sigma_{xz} \\ \Sigma_{zx} & \Sigma_{zz} \end{bmatrix}^{-1} \begin{bmatrix} \widetilde{\mathbf{x}} \\ \widetilde{\mathbf{z}} \end{bmatrix} - \widetilde{\mathbf{z}}^T \Sigma_{zz}^{-1} \widetilde{\mathbf{z}}$$

$$= \begin{bmatrix} \widetilde{\mathbf{x}} \\ \widetilde{\mathbf{z}} \end{bmatrix}^T \begin{bmatrix} I_{xx} & I_{xz} \\ I_{zx} & I_{zz} \end{bmatrix} \begin{bmatrix} \widetilde{\mathbf{x}} \\ \widetilde{\mathbf{z}} \end{bmatrix} - \widetilde{\mathbf{z}}^T \Sigma_{zz}^{-1} \widetilde{\mathbf{z}}$$

$$= \widetilde{\mathbf{x}}^T I_{xx} \widetilde{\mathbf{x}} + \widetilde{\mathbf{x}}^T I_{xz} \widetilde{\mathbf{z}} + \widetilde{\mathbf{z}}^T I_{zx} \widetilde{\mathbf{x}} + \widetilde{\mathbf{z}}^T I_{zz} \widetilde{\mathbf{z}} - \widetilde{\mathbf{z}}^T \Sigma_{zz}^{-1} \widetilde{\mathbf{z}}$$

# Conditional Gaussian

$$\alpha = \widetilde{\mathbf{x}}^T I_{xx} \widetilde{\mathbf{x}} + \widetilde{\mathbf{x}}^T I_{xz} \widetilde{\mathbf{z}} + \widetilde{\mathbf{z}}^T I_{zx} \widetilde{\mathbf{x}} + \widetilde{\mathbf{z}}^T I_{zz} \widetilde{\mathbf{z}} - \widetilde{\mathbf{z}}^T \Sigma_{zz}^{-1} \widetilde{\mathbf{z}}$$

$$= \widetilde{\mathbf{x}}^T I_{xx} \widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^T I_{xx} \Sigma_{xz} \Sigma_{zz}^{-1} \widetilde{\mathbf{z}} - \widetilde{\mathbf{z}}^T \Sigma_{zz}^{-1} \Sigma_{zx} I_{xx} \widetilde{\mathbf{x}}$$

$$+ \widetilde{\mathbf{z}}^T \left( \Sigma_{zz}^{-1} + \Sigma_{zz}^{-1} \Sigma_{zx} I_{xx} \Sigma_{xz} \Sigma_{zz}^{-1} \right) \widetilde{\mathbf{z}} - \widetilde{\mathbf{z}}^T \Sigma_{zz}^{-1} \widetilde{\mathbf{z}}$$

$$= \widetilde{\mathbf{x}}^T I_{xx} \widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^T I_{xx} \Sigma_{xz} \Sigma_{zz}^{-1} \widetilde{\mathbf{z}} - \widetilde{\mathbf{z}}^T \Sigma_{zz}^{-1} \Sigma_{zx} I_{xx} \widetilde{\mathbf{x}} + \widetilde{\mathbf{z}}^T \Sigma_{zz}^{-1} \Sigma_{zx} I_{xx} \Sigma_{xz} \Sigma_{zz}^{-1} \widetilde{\mathbf{z}}$$

$$= \left( \widetilde{\mathbf{x}} - \Sigma_{xz} \Sigma_{zz}^{-1} \widetilde{\mathbf{z}} \right)^T I_{xx} \left( \widetilde{\mathbf{x}} - \Sigma_{xz} \Sigma_{zz}^{-1} \widetilde{\mathbf{z}} \right)$$

See Appendix 2

$$\begin{cases} I_{xx} = \left( \Sigma_{xx} - \Sigma_{xz} \Sigma_{zz}^{-1} \Sigma_{zx} \right)^{-1} \\ I_{xz} = -I_{xx} \Sigma_{xz} \Sigma_{zz}^{-1} \\ I_{zx} = -\Sigma_{zz}^{-1} \Sigma_{zx} I_{xx} \\ I_{zz} = \Sigma_{zz}^{-1} + \Sigma_{zz}^{-1} \Sigma_{zx} I_{xx} \Sigma_{xz} \Sigma_{zz}^{-1} \end{cases}$$

# Conditional Gaussian

$$\alpha = \left(\widetilde{\mathbf{x}} - \Sigma_{xz}\Sigma_{zz}^{-1}\widetilde{\mathbf{z}}\right)^T I_{xx}\left(\widetilde{\mathbf{x}} - \Sigma_{xz}\Sigma_{zz}^{-1}\widetilde{\mathbf{z}}\right)$$

$$I_{xx} = \left(\Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx}\right)^{-1}$$

$$= \left(\widetilde{\mathbf{x}} - \Sigma_{xz}\Sigma_{zz}^{-1}\widetilde{\mathbf{z}}\right)^T \left(\Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx}\right)^{-1}\left(\widetilde{\mathbf{x}} - \Sigma_{xz}\Sigma_{zz}^{-1}\widetilde{\mathbf{z}}\right)$$

$$\begin{cases}\widetilde{\mathbf{x}} = \mathbf{x} - \mu_x \\ \widetilde{\mathbf{z}} = \mathbf{z} - \mu_z\end{cases}$$

$$= \left[\mathbf{x} - \left(\mu_x + \Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z} - \mu_z)\right)\right]^T \left(\Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx}\right)^{-1}\left[\mathbf{x} - \left(\mu_x + \Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z} - \mu_z)\right)\right]$$

$$P(\mathbf{x}|\mathbf{z}) = \frac{P(\mathbf{x}, \mathbf{z})}{P(\mathbf{z})} = \frac{1}{\sqrt{(2\pi)^n|\Sigma_{yy}|/|\Sigma_{zz}|}}e^{-\frac{1}{2}\left[(\mathbf{y}-\mu_y)^T\Sigma_{yy}^{-1}(\mathbf{y}-\mu_y) - (\mathbf{z}-\mu_z)^T\Sigma_{zz}^{-1}(\mathbf{z}-\mu_z)\right]}$$

$$= \frac{1}{\sqrt{(2\pi)^n|\Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx}|}}e^{-\frac{1}{2}\left\{\mathbf{x}-\left[\mu_x+\Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z}-\mu_z)\right]\right\}^T\left(\Sigma_{xx}-\Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx}\right)^{-1}\left\{\mathbf{x}-\left[\mu_x+\Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z}-\mu_z)\right]\right\}}$$

$$\frac{|\Sigma_{yy}|}{|\Sigma_{zz}|} = |\Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx}|$$

# Conditional Gaussian

$$P(\mathbf{x}|\mathbf{z}) = \frac{1}{\sqrt{(2\pi)^n \left| \Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx} \right|}} e^{-\frac{1}{2}\left\{\mathbf{x} - \left[\mu_x + \Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z}-\mu_z)\right]\right\}^T \left(\Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx}\right)^{-1}\left\{\mathbf{x} - \left[\mu_x + \Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z}-\mu_z)\right]\right\}}$$

$$\begin{cases} \mu_{x|z} = \mu_x + \Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z} - \mu_z) \\ \Sigma_{xx|z} = \Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx} \end{cases}$$



- Now, we know how to fuse two data based on jointly Gaussian distribution.

# Sensor Fusion

- Example:
- Location estimation

$x \sim N(10,2)$

$z \sim N(12,1), z = 12.5$

$\text{case}(i) \Sigma_{xz} = 0.1$

$\text{case}(ii) \Sigma_{xz} = 0.5$



$L = ?$

$x^{MMSE} = x^{MAP}$

$z$

$x_0$

# Sensor Fusion

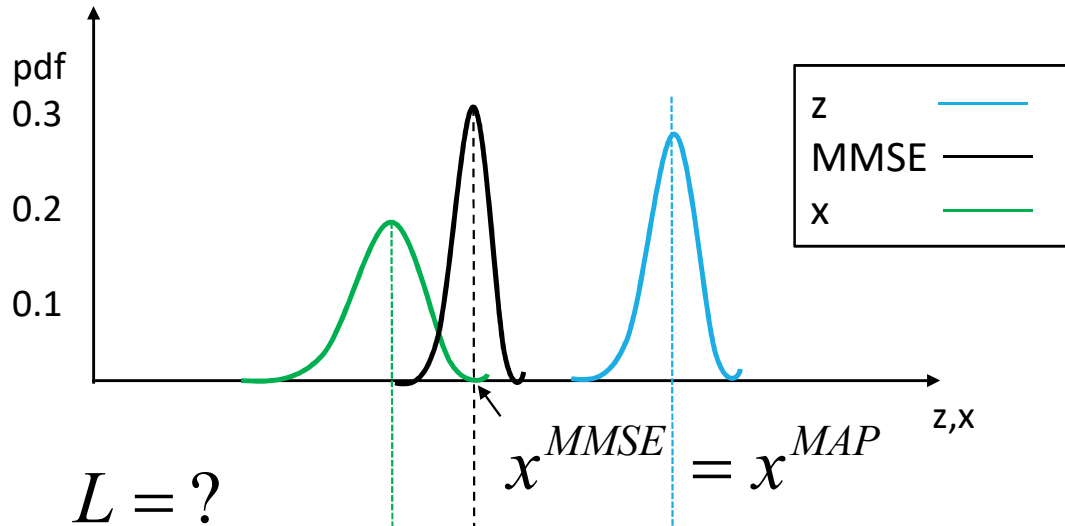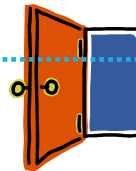$\text{case}(i)\, z = 12.5, x \sim N(10,2), z \sim N(12,1), \underline{\Sigma_{xz}} = 0.1$

$$\begin{cases} \mu_{x|z} = \mu_x + \Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z} - \mu_z) \\ \Sigma_{xx|z} = \Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx} \end{cases}$$

$\mu_{x|z} = 10 + (0.1)1^{-1}(12.5 - 12) = 10.05$

$\Sigma_{xx|z} = 2 - (0.1)1^{-1}(0.1) = 1.99$

$\text{case}(ii)\, z = 12.5, x \sim N(10,2), z \sim N(12,1), \underline{\Sigma_{xz}} = 0.5$

$$\begin{cases} \mu_{x|z} = \mu_x + \Sigma_{xz}\Sigma_{zz}^{-1}(\mathbf{z} - \mu_z) \\ \Sigma_{xx|z} = \Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx} \end{cases}$$
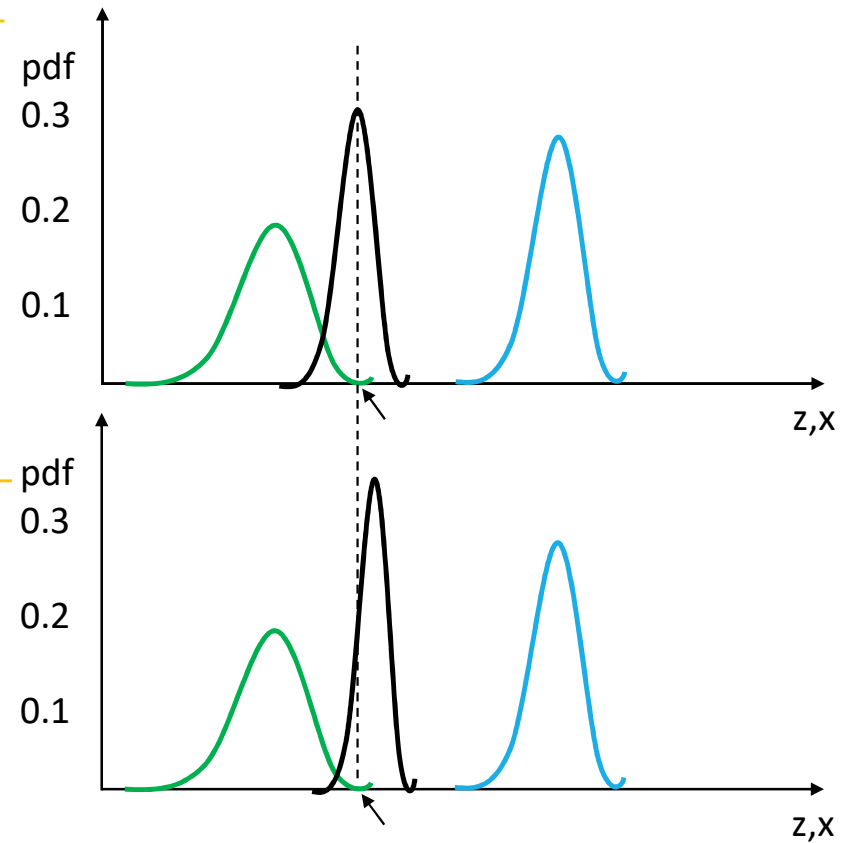
$\mu_{x|z} = 10 + (0.5)1^{-1}(12.5 - 12) = 10.25$

$\Sigma_{xx|z} = 2 - (0.5)1^{-1}(0.5) = 1.75$

# Appendix 1:

Determinant of conditional covariance of jointly Gaussian variables

$$\frac{\left|\Sigma_{yy}\right|}{\left|\Sigma_{zz}\right|} = \left|\Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx}\right|$$

$$\Sigma_{yy} = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xz} \\ \Sigma_{zx} & \Sigma_{zz} \end{bmatrix} = \begin{bmatrix} C & \Sigma_{xz} \\ 0 & \Sigma_{zz} \end{bmatrix}\begin{bmatrix} I_{n\times n} & 0 \\ D & I_{m\times m} \end{bmatrix}$$

$$\Sigma_{xx} = C + \Sigma_{xz}D$$

$$\Sigma_{zx} = \Sigma_{zz}D$$

$$D = \Sigma_{zz}^{-1}\Sigma_{zx}$$

$$C = \Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx}$$

$$\because \left|AB\right| = \left|A\right|\left|B\right|$$

$$\left|\Sigma_{yy}\right| = \left(\left|C\right|\left|\Sigma_{zz}\right|\right)\left(\left|I_{n\times n}\right|\left|I_{m\times m}\right|\right) = \left|\Sigma_{xx} - \Sigma_{xz}\Sigma_{zz}^{-1}\Sigma_{zx}\right|\left|\Sigma_{zz}\right|$$

# Appendix 2:

## Inversion of a Partitioned Matrix

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} E & F \\ G & H \end{bmatrix}$$

*Where*

$$E = A^{-1} + A^{-1}BHCA^{-1} = \left(A - BD^{-1}C\right)^{-1}$$

$$F = -A^{-1}BH = -EBD^{-1}$$

$$G = -HCA^{-1} = -D^{-1}CE$$

$$H = \left(D - CA^{-1}B\right)^{-1} = D^{-1} + D^{-1}CEBD^{-1}$$