Statistical tools for analysing uncertainty in computer model predictions

Jeremy Oakley

School of Mathematics and Statistics



Jeremy Oakley (Sheffield)

Computer Models

January 2015 1 / 23

 \bullet Computer model represented by function f with inputs x and outputs y

$$y = f(x).$$

• Computer model represented by function f with inputs \boldsymbol{x} and outputs \boldsymbol{y}

$$y = f(x).$$

- f constructed from modeller's understanding of the process.
 - There may be no physical input-output data.

• Computer model represented by function f with inputs \boldsymbol{x} and outputs \boldsymbol{y}

$$y = f(x).$$

- f constructed from modeller's understanding of the process.
 - There may be no physical input-output data.
- f may be deterministic.

• Computer model represented by function f with inputs \boldsymbol{x} and outputs \boldsymbol{y}

$$y = f(x).$$

- f constructed from modeller's understanding of the process.
 - There may be no physical input-output data.
- f may be deterministic.
- Computer experiment: evaluating f at difference choices of x
 - A 'model run': evaluating f at a single choice of x.

- Model may be set up to accept 'controllable' inputs only.
- But there may be other parameters/coefficients/variables 'hard-wired' within the model.

- Model may be set up to accept 'controllable' inputs only.
- But there may be other parameters/coefficients/variables 'hard-wired' within the model.
- We define the input x to include these other numerical values used to calculate the outputs.

- Model may be set up to accept 'controllable' inputs only.
- But there may be other parameters/coefficients/variables 'hard-wired' within the model.
- We define the input x to include these other numerical values used to calculate the outputs.
- Suppose that there is a true input value, X, with at least some elements of X uncertain.

- Model may be set up to accept 'controllable' inputs only.
- But there may be other parameters/coefficients/variables 'hard-wired' within the model.
- We define the input x to include these other numerical values used to calculate the outputs.
- Suppose that there is a true input value, X, with at least some elements of X uncertain.
- What is our uncertainty about Y = f(X)?

- Model may be set up to accept 'controllable' inputs only.
- But there may be other parameters/coefficients/variables 'hard-wired' within the model.
- We define the input x to include these other numerical values used to calculate the outputs.
- Suppose that there is a true input value, X, with at least some elements of X uncertain.
- What is our uncertainty about Y = f(X)?
- We quantify uncertainty about X with a probability distribution p_X

- Model may be set up to accept 'controllable' inputs only.
- But there may be other parameters/coefficients/variables 'hard-wired' within the model.
- We define the input x to include these other numerical values used to calculate the outputs.
- Suppose that there is a true input value, X, with at least some elements of X uncertain.
- What is our uncertainty about Y = f(X)?
- We quantify uncertainty about X with a probability distribution p_X
- Then need to obtain the distribution p_Y .

- Model may be set up to accept 'controllable' inputs only.
- But there may be other parameters/coefficients/variables 'hard-wired' within the model.
- We define the input x to include these other numerical values used to calculate the outputs.
- Suppose that there is a true input value, X, with at least some elements of X uncertain.
- What is our uncertainty about Y = f(X)?
- We quantify uncertainty about X with a probability distribution p_X
- Then need to obtain the distribution p_Y .
- Can propagate uncertainty using Monte Carlo: sample X_1, \ldots, X_N from p_X and evaluate $f(X_1), \ldots, f(X_N)$

- Model may be set up to accept 'controllable' inputs only.
- But there may be other parameters/coefficients/variables 'hard-wired' within the model.
- We define the input x to include these other numerical values used to calculate the outputs.
- Suppose that there is a true input value, X, with at least some elements of X uncertain.
- What is our uncertainty about Y = f(X)?
- We quantify uncertainty about X with a probability distribution p_X
- Then need to obtain the distribution p_Y .
- Can propagate uncertainty using Monte Carlo: sample X_1, \ldots, X_N from p_X and evaluate $f(X_1), \ldots, f(X_N)$
- What do we do if f is computationally expensive?



• Want $f(x_1), \ldots, f(x_N)$, but can only evaluate $f(x_1), \ldots, f(x_n)$, for $n \ll N$.



- Want $f(x_1), \ldots, f(x_N)$, but can only evaluate $f(x_1), \ldots, f(x_n)$, for $n \ll N$.
- Could estimate f given $f(x_1), \ldots, f(x_n)$
 - but can we quantify uncertainty in the estimate?



- Want $f(x_1), \ldots, f(x_N)$, but can only evaluate $f(x_1), \ldots, f(x_n)$, for $n \ll N$.
- Could estimate f given $f(x_1), \ldots, f(x_n)$
 - but can we quantify uncertainty in the estimate?
- A statistical inference problem:
 - Treat f as an *uncertain* function
 - Derive a probability distribution for f given $f(x_1),\ldots,f(x_n)$ (an "emulator")



- Want $f(x_1), \ldots, f(x_N)$, but can only evaluate $f(x_1), \ldots, f(x_n)$, for $n \ll N$.
- Could estimate f given $f(x_1), \ldots, f(x_n)$
 - but can we quantify uncertainty in the estimate?
- A statistical inference problem:
 - Treat f as an *uncertain* function
 - Derive a probability distribution for f given $f(x_1),\ldots,f(x_n)$ (an "emulator")

• Popular technique: Gaussian process emulation (Sacks et al 1989)

Example: 18 input climate model, 255 model runs

Example: 18 input climate model, 255 model runs

Emulator means and 95% intervals



simulator output

• Interested in $Y=f(\mathbf{X}),$ where \mathbf{X} is uncertain with distribution $p_{\mathbf{X}}(\mathbf{x}).$

- Interested in $Y=f(\mathbf{X}),$ where \mathbf{X} is uncertain with distribution $p_{\mathbf{X}}(\mathbf{x}).$
- Sensitivity analysis: which elements in $\mathbf{X} = \{X_1, \dots, X_d\}$ are most responsible for the uncertainty in $Y = f(\mathbf{X})$?

- Interested in $Y = f(\mathbf{X})$, where \mathbf{X} is uncertain with distribution $p_{\mathbf{X}}(\mathbf{x})$.
- Sensitivity analysis: which elements in $\mathbf{X} = \{X_1, \dots, X_d\}$ are most responsible for the uncertainty in $Y = f(\mathbf{X})$?
- Write $\mathbf{X} = (\mathbf{X}_u, \mathbf{X}_{-u}).$

- Interested in $Y = f(\mathbf{X})$, where \mathbf{X} is uncertain with distribution $p_{\mathbf{X}}(\mathbf{x})$.
- Sensitivity analysis: which elements in $\mathbf{X} = \{X_1, \dots, X_d\}$ are most responsible for the uncertainty in $Y = f(\mathbf{X})$?
- Write $\mathbf{X} = (\mathbf{X}_u, \mathbf{X}_{-u})$. Consider 'importance' of \mathbf{X}_u via

$$Var_{\mathbf{X}_{u}}\{E_{\mathbf{X}_{-u}}(Y|\mathbf{X}_{u})\}$$

- Interested in $Y = f(\mathbf{X})$, where \mathbf{X} is uncertain with distribution $p_{\mathbf{X}}(\mathbf{x})$.
- Sensitivity analysis: which elements in $\mathbf{X} = \{X_1, \dots, X_d\}$ are most responsible for the uncertainty in $Y = f(\mathbf{X})$?
- Write $\mathbf{X} = (\mathbf{X}_u, \mathbf{X}_{-u})$. Consider 'importance' of \mathbf{X}_u via

$$Var_{\mathbf{X}_{u}}\{E_{\mathbf{X}_{-u}}(Y|\mathbf{X}_{u})\}$$

• The expected reduction in variance if value of \mathbf{X}_u is learnt, because

$$Var(Y) = Var_{\mathbf{X}_{u}} \{ E_{\mathbf{X}_{-u}}(Y|\mathbf{X}_{u}) \} + E_{\mathbf{X}_{u}} \{ Var_{\mathbf{X}_{-u}}(Y|\mathbf{X}_{u}) \}$$

- Interested in $Y = f(\mathbf{X})$, where \mathbf{X} is uncertain with distribution $p_{\mathbf{X}}(\mathbf{x})$.
- Sensitivity analysis: which elements in $\mathbf{X} = \{X_1, \dots, X_d\}$ are most responsible for the uncertainty in $Y = f(\mathbf{X})$?
- Write $\mathbf{X} = (\mathbf{X}_u, \mathbf{X}_{-u})$. Consider 'importance' of \mathbf{X}_u via

$$Var_{\mathbf{X}_{u}}\{E_{\mathbf{X}_{-u}}(Y|\mathbf{X}_{u})\}$$

• The expected reduction in variance if value of \mathbf{X}_u is learnt, because

$$Var(Y) = Var_{\mathbf{X}_u} \{ E_{\mathbf{X}_{-u}}(Y|\mathbf{X}_u) \} + E_{\mathbf{X}_u} \{ Var_{\mathbf{X}_{-u}}(Y|\mathbf{X}_u) \}$$

• Can speed up computation with emulator (Oakley & O'Hagan, 2004)

$$- p_{X_1}(x_1)$$
$$- E(Y|X_1 = x_1)$$

$$- p_{X_2}(x_2) - E(Y|X_2 = x_2)$$

Computer Models

January 2015 8 / 23

¹MUCM case study: analysis by John Paul Gosling, Hugo Maruri-Aguilar, Alexis Boukouvalas

 Model developed by GlaxoSmithKline. Predicts incidence of rotavirus in a population before and after a vaccine is administered to a proportion of the infant population

¹MUCM case study: analysis by John Paul Gosling, Hugo Maruri-Aguilar, Alexis Boukouvalas

Example¹: modelling Rotavirus

- Model developed by GlaxoSmithKline. Predicts incidence of rotavirus in a population before and after a vaccine is administered to a proportion of the infant population
- Deterministic compartmental model, 672 compartments (16 disease stages × 42 age classes)

¹MUCM case study: analysis by John Paul Gosling, Hugo Maruri-Aguilar, Alexis Boukouvalas

Example¹: modelling Rotavirus

- Model developed by GlaxoSmithKline. Predicts incidence of rotavirus in a population before and after a vaccine is administered to a proportion of the infant population
- Deterministic compartmental model, 672 compartments (16 disease stages \times 42 age classes)
- Inputs include transmission rates between age groups, reduction in risk following each infection

¹MUCM case study: analysis by John Paul Gosling, Hugo Maruri-Aguilar, Alexis Boukouvalas
Example¹: modelling Rotavirus

- Model developed by GlaxoSmithKline. Predicts incidence of rotavirus in a population before and after a vaccine is administered to a proportion of the infant population
- Deterministic compartmental model, 672 compartments (16 disease stages \times 42 age classes)
- Inputs include transmission rates between age groups, reduction in risk following each infection
- Outputs: time series of rotavirus incidence for six age groups following vaccination programme

¹MUCM case study: analysis by John Paul Gosling, Hugo Maruri-Aguilar, Alexis Boukouvalas

Jeremy Oakley (Sheffield)

Example¹: modelling Rotavirus

- Model developed by GlaxoSmithKline. Predicts incidence of rotavirus in a population before and after a vaccine is administered to a proportion of the infant population
- Deterministic compartmental model, 672 compartments (16 disease stages \times 42 age classes)
- Inputs include transmission rates between age groups, reduction in risk following each infection
- Outputs: time series of rotavirus incidence for six age groups following vaccination programme
- GSK analysis investigated sensitivity of output to 9 inputs, using 8200 model runs
- We consider sensitivity of output to 20 inputs, using 340 model runs

Jeremy Oakley (Sheffield)

¹MUCM case study: analysis by John Paul Gosling, Hugo Maruri-Aguilar, Alexis Boukouvalas

Variance based sensitivity analysis

Analysis for an individual output: no. of infections in 2-3 age group after 2 years



Need to think carefully about input distributions

Need to think carefully about input distributions

Consider

$$f(x) = \exp(-x),$$

with Y = f(X) and

 $X \sim U[0,b].$

Consider

$$f(x) = \exp(-x),$$

with $\boldsymbol{Y}=\boldsymbol{f}(\boldsymbol{X})$ and

 $X \sim U[0,b].$

In this case we have

$$Var(Y) = \frac{b - 2 + 4\exp(-b) - (b + 2)\exp(-2b)}{2b^2},$$

Consider

$$f(x) = \exp(-x),$$

with $\boldsymbol{Y}=\boldsymbol{f}(\boldsymbol{X})$ and

 $X \sim U[0,b].$

In this case we have

$$Var(Y) = \frac{b - 2 + 4\exp(-b) - (b + 2)\exp(-2b)}{2b^2},$$

Increasing b increases the variance of X but *decreases* the variance of Y



Sensitivity analysis for decision making



Sensitivity analysis for decision making



Jeremy Oakley (Sheffield)

• A Gaussian plume deposition model $f(x_{cont}, x_{calib})$ predicts deposition of radionuclides at a location x_{cont} following release of unknown concentration X_{calib} from point source

- A Gaussian plume deposition model $f(x_{cont}, x_{calib})$ predicts deposition of radionuclides at a location x_{cont} following release of unknown concentration X_{calib} from point source
- Measurements of the true deposition $z(x_{cont})$ at a limited number of locations x_{cont} available.

- A Gaussian plume deposition model $f(x_{cont}, x_{calib})$ predicts deposition of radionuclides at a location x_{cont} following release of unknown concentration X_{calib} from point source
- Measurements of the true deposition $z(x_{cont})$ at a limited number of locations x_{cont} available.
- Aim: to predict deposition at other locations using both data and model.

- A Gaussian plume deposition model $f(x_{cont}, x_{calib})$ predicts deposition of radionuclides at a location x_{cont} following release of unknown concentration X_{calib} from point source
- Measurements of the true deposition $z(x_{cont})$ at a limited number of locations x_{cont} available.
- Aim: to predict deposition at other locations using both data and model.
- What value of x_{calib} do we use?

- A Gaussian plume deposition model $f(x_{cont}, x_{calib})$ predicts deposition of radionuclides at a location x_{cont} following release of unknown concentration X_{calib} from point source
- Measurements of the true deposition $z(x_{cont})$ at a limited number of locations x_{cont} available.
- Aim: to predict deposition at other locations using both data and model.
- What value of x_{calib} do we use?
- And what happens if the model is wrong?

• Wish to estimate $X_{calib} = g$: acceleration due to Earth's gravity



- Wish to estimate $X_{calib} = g$: acceleration due to Earth's gravity
- I drop a tennis ball from my office window at height x_{cont} , and time its descent to the ground



- Wish to estimate $X_{calib} = g$: acceleration due to Earth's gravity
- I drop a tennis ball from my office window at height x_{cont} , and time its descent to the ground
- Estimate g via

$$t = \sqrt{2x_{cont}}g$$



- Wish to estimate $X_{calib} = g$: acceleration due to Earth's gravity
- I drop a tennis ball from my office window at height x_{cont} , and time its descent to the ground
- Estimate g via

$$t = \sqrt{2x_{cont}}g$$

• Will have error in measurements, so take replicates



- Wish to estimate $X_{calib} = g$: acceleration due to Earth's gravity
- I drop a tennis ball from my office window at height x_{cont}, and time its descent to the ground
- Estimate g via

$$t = \sqrt{2x_{cont}}g$$

- Will have error in measurements, so take replicates
- The more measurements I take, the more certain I become about the *wrong value*



- Wish to estimate $X_{calib} = g$: acceleration due to Earth's gravity
- I drop a tennis ball from my office window at height x_{cont}, and time its descent to the ground
- Estimate g via

$$t = \sqrt{2x_{cont}}g$$

- Will have error in measurements, so take replicates
- The more measurements I take, the more certain I become about the *wrong value*



$$z(x_{cont,i}) = f(x_{cont,i}, X_{calib}) + \delta(x_{cont,i}) + \varepsilon_i$$

 $\delta(x_{cont,i})$ is the discrepancy (bias) between model output and reality.

$$z(x_{cont,i}) = f(x_{cont,i}, X_{calib}) + \delta(x_{cont,i}) + \varepsilon_i$$

 $\delta(x_{cont,i}) \text{ is the discrepancy (bias) between model output and reality.}$ • Doesn't always go down well with modellers!

$$z(x_{cont,i}) = f(x_{cont,i}, X_{calib}) + \delta(x_{cont,i}) + \varepsilon_i$$

 $\delta(x_{cont,i}) \text{ is the discrepancy (bias) between model output and reality.}$ • Doesn't always go down well with modellers!

"I'm horrified! You should be improving your models with better physics!"

$$z(x_{cont,i}) = f(x_{cont,i}, X_{calib}) + \delta(x_{cont,i}) + \varepsilon_i$$

 $\delta(x_{cont,i})$ is the **discrepancy** (bias) between model output and reality. • Doesn't always go down well with modellers!

"I'm horrified! You should be improving your models with better physics!"

- Accounting for model discrepancy important, otherwise
 - can become certain about a 'wrong' input value
 - model predictions can be spuriously precise

$$z(x_{cont,i}) = f(x_{cont,i}, X_{calib}) + \delta(x_{cont,i}) + \varepsilon_i$$

 $\delta(x_{cont,i})$ is the discrepancy (bias) between model output and reality.

• Doesn't always go down well with modellers!

"I'm horrified! You should be improving your models with better physics!"

- Accounting for model discrepancy important, otherwise
 - can become certain about a 'wrong' input value
 - model predictions can be spuriously precise
- In some settings, can infer $\delta()$ from observations, but in others, (prior) judgements important

• Wish to find x such that $f(\boldsymbol{x})$ is 'close' to observation \boldsymbol{z}

- Wish to find x such that f(x) is 'close' to observation z
- Model is slow to run, so will use an emulator to explore input space

- Wish to find x such that f(x) is 'close' to observation z
- Model is slow to run, so will use an emulator to explore input space
- Emphasis now on discarding region of input space where the model can't fit the data

- Wish to find x such that f(x) is 'close' to observation z
- Model is slow to run, so will use an emulator to explore input space
- Emphasis now on discarding region of input space where the model can't fit the data
- First run the model and construct an emulator for f

- Wish to find x such that f(x) is 'close' to observation z
- Model is slow to run, so will use an emulator to explore input space
- Emphasis now on discarding region of input space where the model can't fit the data
- First run the model and construct an emulator for f
- Assess the "implausibility" of an input value \boldsymbol{x} via

$$I(x) = \frac{|z - E\{f(x)\}|}{\left[Var\{f(x)\} + Var(\varepsilon) + Var(\delta)\right]^{1/2}}.$$

- Wish to find x such that $f(\boldsymbol{x})$ is 'close' to observation \boldsymbol{z}
- Model is slow to run, so will use an emulator to explore input space
- Emphasis now on discarding region of input space where the model can't fit the data
- First run the model and construct an emulator for f
- \bullet Assess the "implausibility" of an input value x via

$$I(x) = \frac{|z - E\{f(x)\}]|}{\left[Var\{f(x)\} + Var(\varepsilon) + Var(\delta)\right]^{1/2}}.$$

• Do further runs in the non-implausible region and rebuild emulator

Toy example



Jeremy Oakley (Sheffield)

Computer Models

18 / 23

Toy example


- "Mukwano": a dynamic, stochastic, individual based model that simulates sexual partnerships and HIV transmission
- Births, deaths, partnership formation and dissolution and HIV transmission were modelled using time-dependent rates

Jeremy Oakley (Sheffield)

²I. Andrianakis, I. Vernon, N. McCreesh, T.J. McKinley, J.O., R. Nsubuga, M. Goldstein and R.G. White

- "Mukwano": a dynamic, stochastic, individual based model that simulates sexual partnerships and HIV transmission
- Births, deaths, partnership formation and dissolution and HIV transmission were modelled using time-dependent rates
- 22 inputs, e.g. proportions of men and women in "high sexual activity" groups, transition probability of HIV per sex act during primary stage of infection

Jeremy Oakley (Sheffield)

²I. Andrianakis, I. Vernon, N. McCreesh, T.J. McKinley, J.O., R. Nsubuga, M. Goldstein and R.G. White

- Calibration data were collected from a rural general population cohort in South-West Uganda. The cohort was established in 1989 and currently consists of the residents of 25 villages
 - 18 demographic, behavioural and epidemiological outputs, e.g., male and female HIV prevalences at three time points

- Calibration data were collected from a rural general population cohort in South-West Uganda. The cohort was established in 1989 and currently consists of the residents of 25 villages
 - 18 demographic, behavioural and epidemiological outputs, e.g., male and female HIV prevalences at three time points
- Model run on a high performance cluster with 240 nodes. The run time for a single simulation was around 10 minutes

- Calibration data were collected from a rural general population cohort in South-West Uganda. The cohort was established in 1989 and currently consists of the residents of 25 villages
 - 18 demographic, behavioural and epidemiological outputs, e.g., male and female HIV prevalences at three time points
- Model run on a high performance cluster with 240 nodes. The run time for a single simulation was around 10 minutes
- History matching iterated through 10 waves, 200-500 model runs per wave





January 2015

Jeremy Oakley (Sheffield)













Model runs after history matching



Summary

• Emulators for computationally expensive models. Helpful also in calibration.

- Emulators for computationally expensive models. Helpful also in calibration.
- Sensitivity analysis tools for investigating input uncertainty.

- Emulators for computationally expensive models. Helpful also in calibration.
- Sensitivity analysis tools for investigating input uncertainty.
- Model discrepancy essential for 'complete' uncertainty quantification.

- Emulators for computationally expensive models. Helpful also in calibration.
- Sensitivity analysis tools for investigating input uncertainty.
- Model discrepancy essential for 'complete' uncertainty quantification.

Further reading/papers at

```
jeremy-oakley.staff.shef.ac.uk
```