

Regression methods

(CIV6540 - Probabilistic Machine Learning for Civil Engineers)

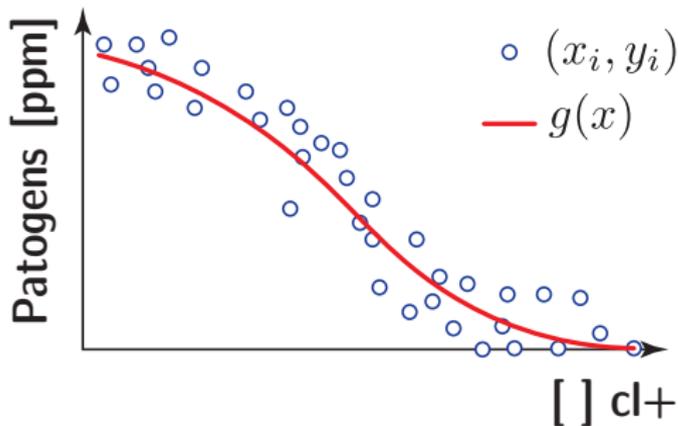
Professor: James-A. Goulet

Département des génies civil, géologique et des mines
Polytechnique Montréal



Chapter 8 – Goulet (2020)
Probabilistic Machine Learning for Civil Engineers
MIT Press

What is regression?



Data

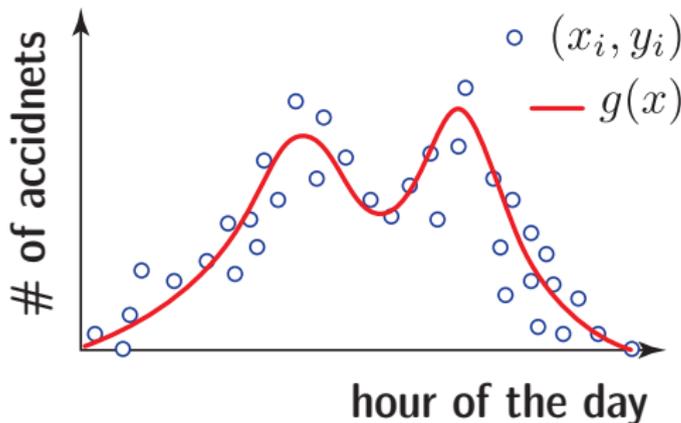
$$\mathcal{D} = \{(x_i, y_i), \forall i = 1 : D\}$$

$$x_i \in \mathbb{R} : \begin{cases} \text{Covariate} \\ \text{attribute} \\ \text{regressor} \end{cases}$$

$$y_i \in \mathbb{R} : \text{Observation}$$

Regression methods: mathematical models for $g(x)$

What is regression?



Data

$$\mathcal{D} = \{(x_i, y_i), \forall i = 1 : D\}$$

$$x_i \in \mathbb{R} : \begin{cases} \text{Covariate} \\ \text{attribute} \\ \text{regressor} \end{cases}$$

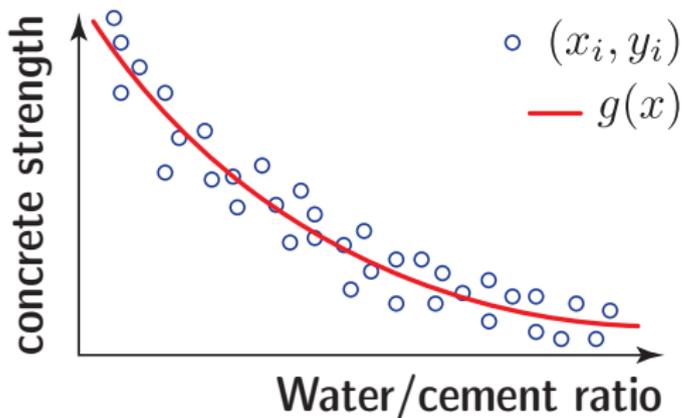
$$y_i \in \mathbb{R} : \text{Observation}$$

Model

$$g(x) \equiv \text{fct}(x)$$

Regression methods: mathematical models for $g(x)$

What is regression?



Data

$$\mathcal{D} = \{(x_i, y_i), \forall i = 1 : D\}$$

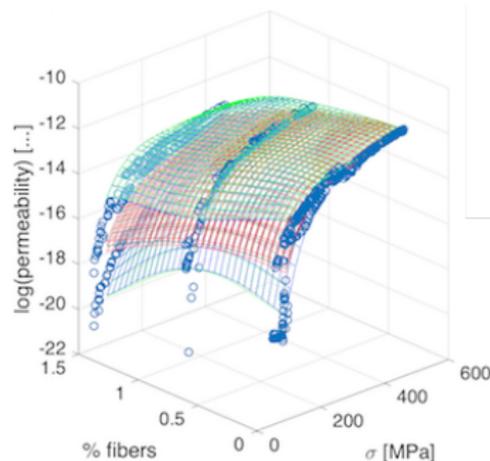
$$x_i \in \mathbb{R} : \begin{cases} \text{Covariate} \\ \text{attribute} \\ \text{regressor} \end{cases}$$

$$y_i \in \mathbb{R} : \text{Observation}$$

Model

$$g(x) \equiv \text{fct}(x)$$

Regression methods: mathematical models for $g(x)$

Regression in $> 1D$ **Data**

$$\mathcal{D} = \{(\mathbf{x}_i, y_i), \forall i = 1 : D\}$$

$$\mathbf{x}_i = [x_1, x_2]^T \in \mathbb{R}^2 : \begin{cases} \text{Covariates} \\ \text{attributes} \\ \text{regressors} \end{cases}$$

$$y_i \in \mathbb{R}^1 : \text{Observation}$$

Model

$$g(\mathbf{x}) \equiv \text{fct}(\mathbf{x}) \in \mathbb{R}^1$$

General case: $\mathbf{x}_i = [x_1, x_2, \dots, x_x]^T \in \mathbb{R}^x$, $y_i \in \mathbb{R}^1$, $g(\mathbf{x}_i) \in \mathbb{R}^1$

What is regression?

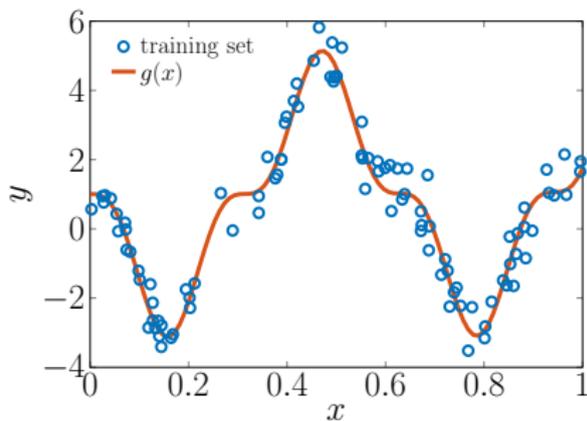
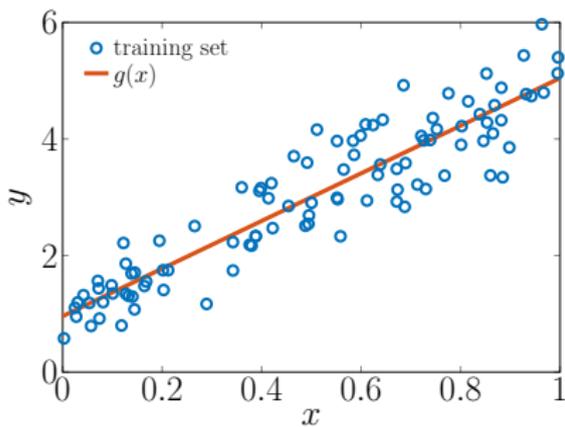
Module #5 Outline

Intro.
Linear Regression
Gaussian Process
Regression
Examples
Adv. topics
Neural Networks
S.

Topics organization

<i>Background</i>	<div style="display: flex; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px; margin-right: 5px;">1</div> <div style="margin-right: 5px;">}</div> <div style="margin-right: 5px;">Revision probability & linear algebra</div> <div style="font-size: 1em;">🌐</div> </div> <div style="display: flex; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px; margin-right: 5px;">2</div> <div style="margin-right: 5px;">}</div> <div style="margin-right: 5px;">Probability distributions</div> <div style="font-size: 1em;">📐📐</div> </div>
<i>Machine Learning Basics</i>	<div style="display: flex; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px; margin-right: 5px;">0</div> <div style="margin-right: 5px;">}</div> <div style="margin-right: 5px;">Introduction</div> <div style="font-size: 1em;">📖</div> </div> <div style="display: flex; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px; margin-right: 5px;">3</div> <div style="margin-right: 5px;">}</div> <div style="margin-right: 5px;">Bayesian Estimation</div> <div style="margin-right: 5px;">$P(A B) = \frac{P(B A)P(A)}{P(B)}$</div> </div> <div style="display: flex; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px; margin-right: 5px;">4</div> <div style="margin-right: 5px;">}</div> <div style="margin-right: 5px;">MCMC sampling & Newton</div> <div style="font-size: 1em;">📦🔴</div> </div>
<i>Supervised learning</i>	<div style="display: flex; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px; margin-right: 5px;">5</div> <div style="margin-right: 5px;">}</div> <div style="margin-right: 5px;">Regression</div> <div style="font-size: 1em;">📈</div> </div> <div style="display: flex; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px; margin-right: 5px;">6</div> <div style="margin-right: 5px;">}</div> <div style="margin-right: 5px;">Classification</div> <div style="font-size: 1em;">🚂🌪️</div> </div> <div style="display: flex; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px; margin-right: 5px;">7</div> <div style="margin-right: 5px;">}</div> <div style="margin-right: 5px;">LSTM networks for time series</div> <div style="font-size: 1em;">🚗</div> </div>
<i>Unsupervised learning</i>	<div style="display: flex; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px; margin-right: 5px;">7</div> <div style="margin-right: 5px;">}</div> <div style="margin-right: 5px;">State-space model for time-series</div> <div style="font-size: 1em;">📈</div> </div>
<i>Decision Making & RL</i>	<div style="display: flex; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px; margin-right: 5px;">8</div> <div style="margin-right: 5px;">}</div> <div style="margin-right: 5px;">Decision Theory</div> <div style="font-size: 1em;">📐</div> </div> <div style="display: flex; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px; margin-right: 5px;">9</div> <div style="margin-right: 5px;">}</div> <div style="margin-right: 5px;">AI & Sequential decision problems</div> <div style="font-size: 1em;">🧠</div> </div>

Which one is an example of linear regression?



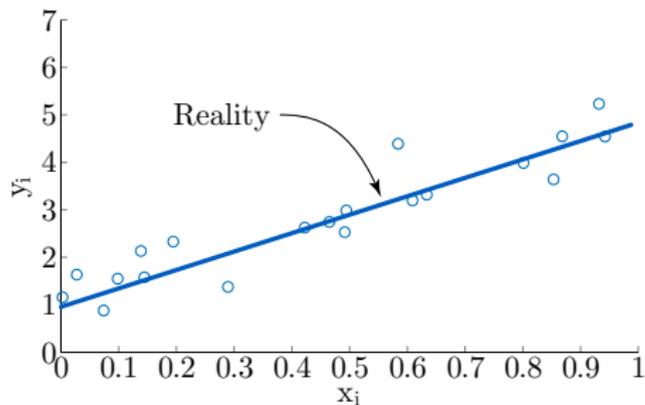
Both.

Section Outline

Linear Regression

- 2.1 Introduction
- 2.2 Mathematical formulation 1D
- 2.3 Cross-validation
- 2.4 Mathematical formulation $>1D$
- 2.5 Limitation

Nomenclature & hypotheses



Data

$$\mathcal{D} = \{(x_i, y_i), \forall i = 1 : D\}$$

$$x_i \in \mathbb{R} : \begin{cases} \text{Covariate} \\ \text{attribute} \\ \text{regressor} \end{cases}$$

$$y_i \in \mathbb{R} : \text{Observation}$$

Model

$$g(x) \equiv \text{fct}(x)$$

Hypothesis: $g(x) \equiv \text{reality}$, i.e. no variability

$$y = g(x) + v, \quad v : V \sim \mathcal{N}(v; 0, \sigma_V^2), \quad \overbrace{V_i \perp V_j, \forall i \neq j}^{\text{independent observation errors}}$$

Mathematical formulation – special case

Let us assume that our model takes the form

$$g(x) = b_0 + b_1x = [b_0 \ b_1] \begin{bmatrix} 1 \\ x \end{bmatrix}$$

$$y = g(x) + v, \quad v : V \sim \mathcal{N}(v; 0, \sigma_V^2)$$

For the entire dataset $\mathcal{D} = \{(x_i, y_i), \forall i = 1 : D\} = \{\mathcal{D}_y, \mathcal{D}_x\}$

$$\underbrace{\mathbf{b}}_{\theta} = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}, \quad \underbrace{\mathbf{y}}_{\mathcal{D}_y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_D \end{bmatrix}, \quad \underbrace{\mathbf{X}}_{\mathcal{D}_x} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_D \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_D \end{bmatrix}$$

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{v}, \quad \mathbf{v} : \mathbf{V} \sim \mathcal{N}(\mathbf{v}; \mathbf{0}, \sigma_V^2 \cdot \mathbf{I})$$

Joint likelihood of $\mathcal{D} = \{(x_i, y_i), \forall i = 1 : D\} = \{\mathcal{D}_y, \mathcal{D}_x\}$

$$\underbrace{\mathbf{y}}_{\mathcal{D}_y} = \underbrace{\mathbf{X}}_{\mathcal{D}_x} \underbrace{\mathbf{b}}_{\boldsymbol{\theta}} + \mathbf{v}, \quad \mathbf{v} : \mathbf{V} \sim \mathcal{N}(\mathbf{v}; \mathbf{0}, \sigma_V^2 \cdot \mathbf{I})$$

$$\begin{aligned} f(\mathcal{D}_y | \mathcal{D}_x, \boldsymbol{\theta}) &= f(\mathcal{D}_y = \mathbf{y} | \mathcal{D}_x = \mathbf{x}, \boldsymbol{\theta} = \mathbf{b}) \\ &= \prod_{i=1}^D \mathcal{N}(y_i; g(x_i), \sigma_V^2) \\ &\propto \prod_{i=1}^D \exp\left(-\frac{1}{2} \frac{(y_i - g(x_i))^2}{\sigma_V^2}\right) \\ &= \exp\left(-\frac{1}{2} \sum_{i=1}^D \frac{(y_i - g(x_i))^2}{\sigma_V^2}\right) \\ &= \exp\left(-\frac{1}{2\sigma_V^2} (\mathbf{y} - \mathbf{Xb})^\top (\mathbf{y} - \mathbf{Xb})\right) \end{aligned}$$

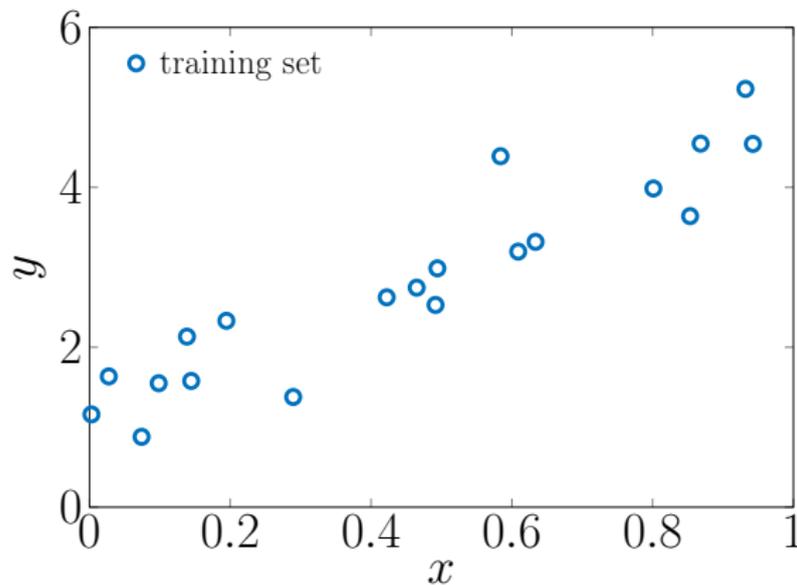
Estimating optimal parameters $\theta^* = \mathbf{b}^*$

The goal is to estimate optimal parameters $\theta^* = \mathbf{b}^*$ which maximize the likelihood $f(\mathcal{D}_y | \mathcal{D}_x, \theta)$, or equivalently which **maximize** the log-likelihood

$$\begin{aligned}
 \ln f(\mathcal{D}_y | \mathcal{D}_x, \theta) &= -\frac{1}{2\sigma_V^2}(\mathbf{y} - \mathbf{X}\mathbf{b})^T(\mathbf{y} - \mathbf{X}\mathbf{b}) \\
 &\propto -\frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{b})^T(\mathbf{y} - \mathbf{X}\mathbf{b}) \\
 &= -J(\mathbf{b}) \quad (\text{Loss function})
 \end{aligned}$$

In the context of linear regression, we **minimize** $J(\mathbf{b})$, which is **equivalent to maximizing** $f(\mathcal{D}_y | \mathcal{D}_x, \theta)$

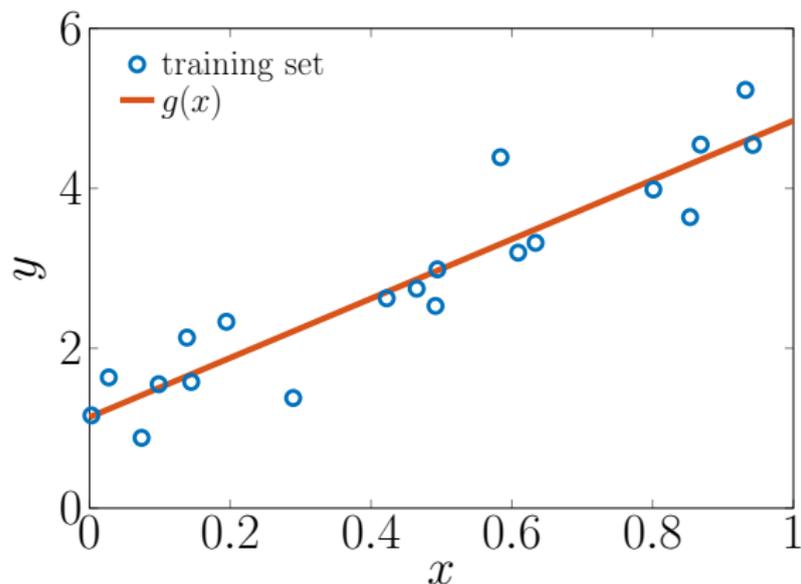
Example #1 [🐍]



$$g(x) = b_0 + b_1x$$

$$\mathbf{b}^* = [0.96, 4.09]^T$$

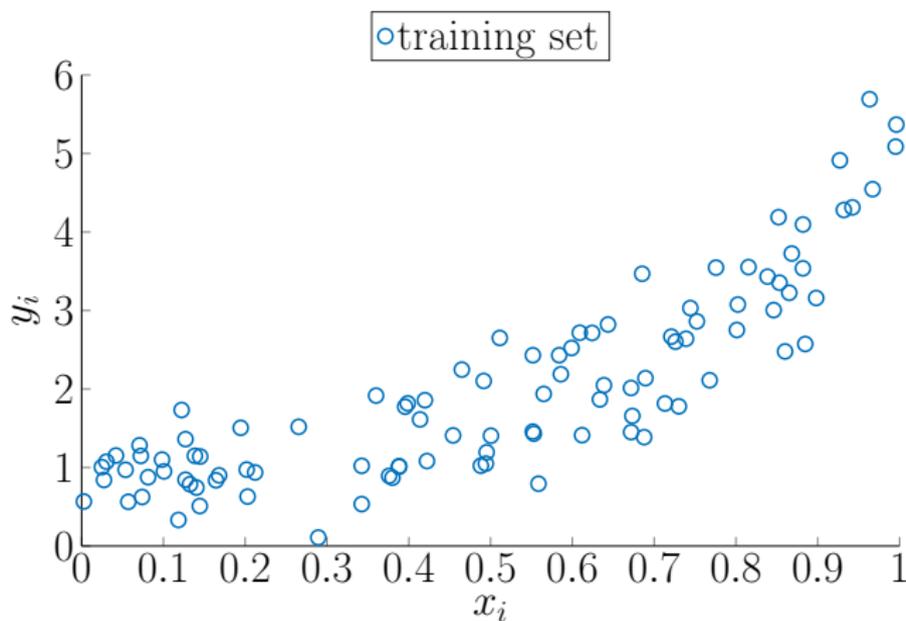
Example #1 [🐍]



$$g(x) = b_0 + b_1x$$

$$\mathbf{b}^* = [0.96, 4.09]^T$$

What if the data is non-linear?



Mathematical formulation – general case

$$g(x) = b_0 + b_1\phi_1(x) + b_2\phi_2(x) + \dots + b_B\phi_B(x) = \mathbf{X}\mathbf{b}$$

where $\phi_j(x)$ can be **any linear or non-linear function**,

$$\text{e.g.: } \phi_j(x) = x^2, \phi_j(x) = \sin(x), \dots$$

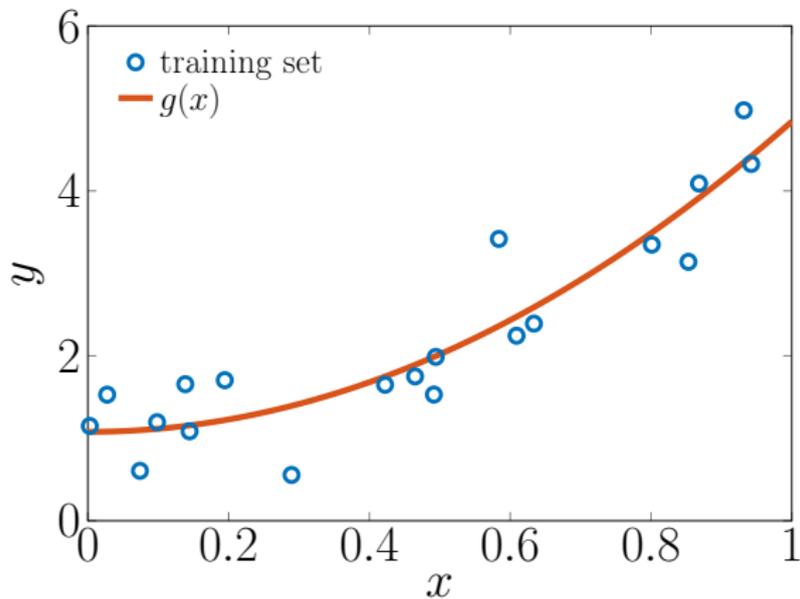
For the entire dataset

$$\mathbf{b} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_B \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_D \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & \phi_1(x_1) & \phi_2(x_1) & \dots & \phi_B(x_1) \\ 1 & \phi_1(x_2) & \phi_2(x_2) & \dots & \phi_B(x_2) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1(x_D) & \phi_2(x_D) & \dots & \phi_B(x_D) \end{bmatrix}$$

Linear regression: Linear with respect to $\phi_j(x)$ and \mathbf{b} not x

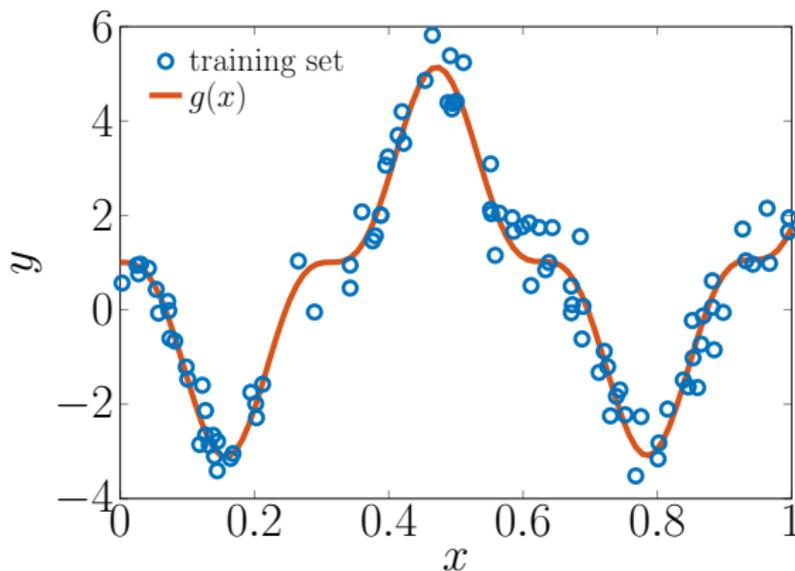
Example #2a [Python]

$$g(x) = b_0 + b_1x^2$$



Example #2b [Python]

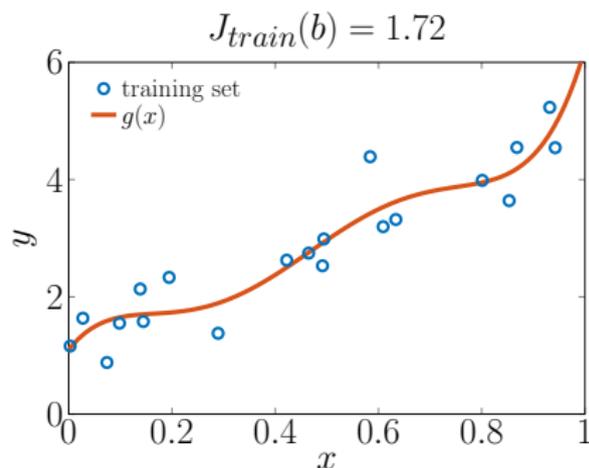
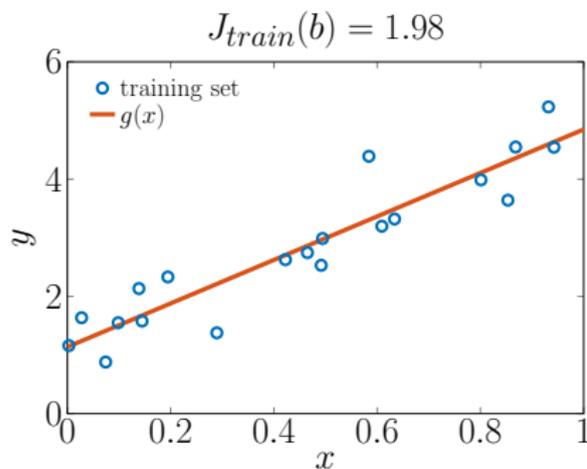
$$g(x) = b_0 + b_1x^3 + b_2 \sin(10x)$$



What model structure is best?

$$g(x) = b_0 + b_1x$$

$$g(x) = b_0 + b_1x + b_2x^2 + b_3x^3 + b_4x^4 + b_5x^5$$



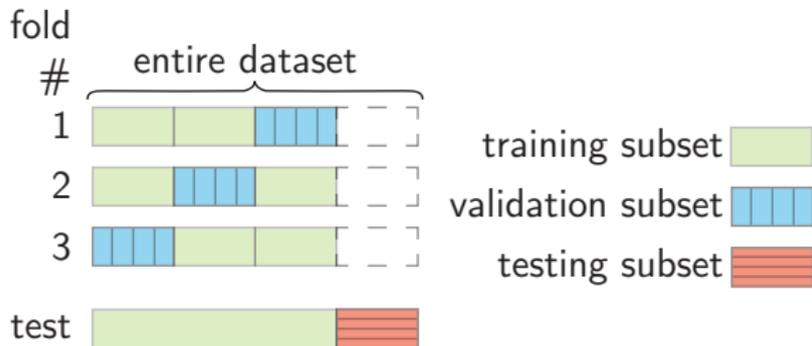
Linear regression is sensitive to over-fitting

Cross-validation – methodology

Over-fitting is prevented by employing cross-validation

Cross-validation: Do not employ the same data to train (i.e. estimate optimal parameters) and to evaluate the loss function.

N -fold cross-validation: Split the **shuffled** data into N subsets that are employed one at a time as an independent validation set



Algorithm 1: N -fold cross-validation

```

1  define  $\mathbf{x}, \mathbf{y}, g_m(\mathbf{x})$ 
2  CV_idx = randperm(D)
3  nCV = round(D/N)
4  for  $m = 1, 2, \dots, M$  do
5      for  $i = 1, 2, \dots, N$  do
6          val_idx = CV_idx((i - 1) × nCV + 1 : min(X, i × nCV))
7          x_train =  $\mathbf{x}$ , x_train(val_idx,:) = []
8          y_train =  $\mathbf{y}$ , y_train(val_idx,:) = []
9          x_val =  $\mathbf{x}$ (val_idx,:)
10         y_val =  $\mathbf{y}$ (val_idx,:)
11         Estimate  $\mathbf{b}_i^* = (\mathbf{X}_{train}^T \mathbf{X}_{train})^{-1} \mathbf{X}_{train}^T \mathbf{y}_{train}$ 
12         Compute  $J_i(\mathbf{b}^*, \mathbf{x}_{val}, \mathbf{y}_{val})$ 
13      $J_{val}^m = \sum_{i=1}^N J_i$ 
14 Select  $m^* : J_{val}^{m^*} = \min J_{val}^m$ 
15 Estimate  $\mathbf{b}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ 

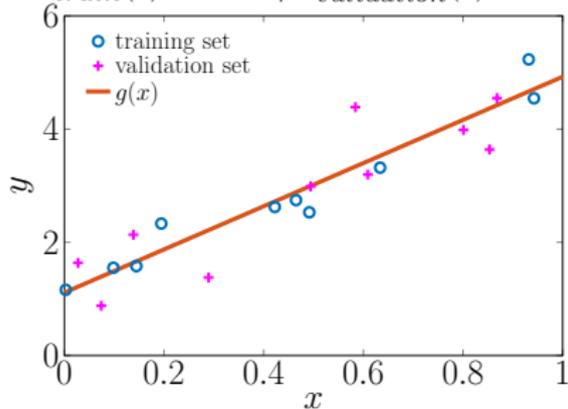
```

Example #3 – Cross-validation

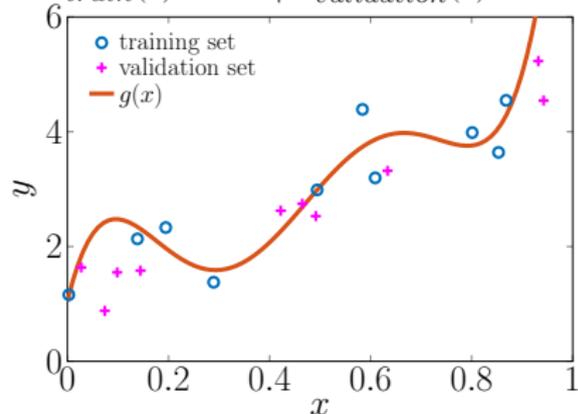
$$g(x) = b_0 + b_1x$$

$$g(x) = b_0 + b_1x + b_2x^2 + b_3x^3 + b_4x^4 + b_5x^5$$

$$J_{train}(b) = 0.43 \mid J_{validation}(b) = 1.56$$



$$J_{train}(b) = 0.7 \mid J_{validation}(b) = 5.42$$



$$\left. \begin{aligned} \sum J_{train} &= 1.97 \\ \sum J_{validation} &= 2.02 \end{aligned} \right\} \triangle$$

$$\left. \begin{aligned} \sum J_{train} &= 0.81 \\ \sum J_{validation} &= 20.89 \end{aligned} \right\} \triangle \times$$

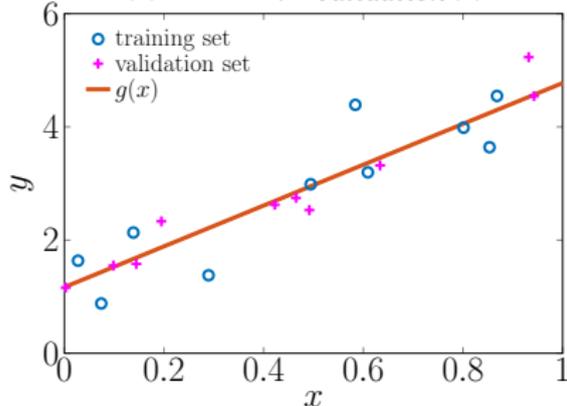
[CIV_ML/Linear_regression_1D.py]

Example #3 – Cross-validation

$$g(x) = b_0 + b_1x$$

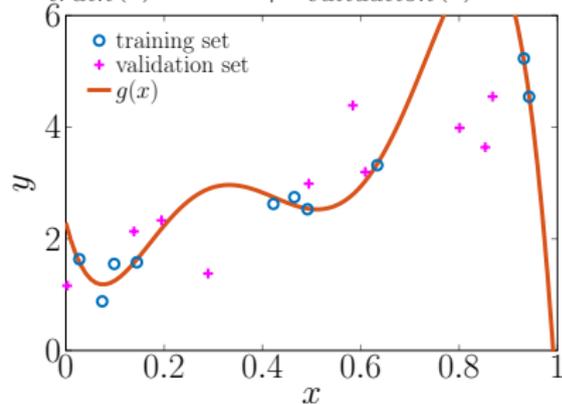
$$g(x) = b_0 + b_1x + b_2x^2 + b_3x^3 + b_4x^4 + b_5x^5$$

$$J_{train}(b) = 1.53 \mid J_{validation}(b) = 0.46$$



$$\left. \begin{array}{l} \sum J_{train} = 1.97 \\ \sum J_{validation} = 2.02 \end{array} \right\} \triangle$$

$$J_{train}(b) = 0.11 \mid J_{validation}(b) = 15.47$$



$$\left. \begin{array}{l} \sum J_{train} = 0.81 \\ \sum J_{validation} = 20.89 \end{array} \right\} \triangle \times$$

[CIV_ML/Linear_regression_1D.py]

N -fold cross-validation

The greater is the number of folds N , the more accurate are results.

If $N = D$ it is called **leave-one-out cross-validation (LOOCV)**

In practice LOOCV too computationally demanding when D is large because it requires as many calibration loops

Instead in practice, it is common to employ **5 to 20-fold cross-validation**

Regression in $> 1D$

Data

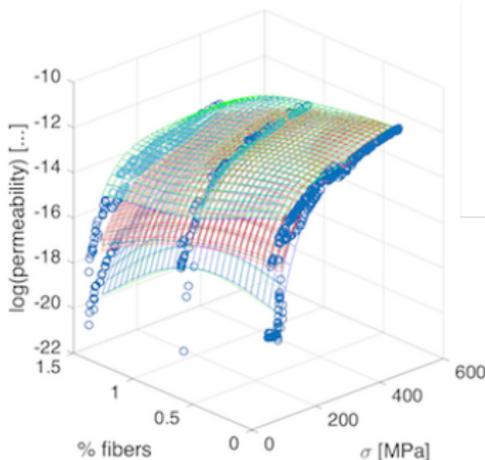
$$\mathcal{D} = \{(\mathbf{x}_i, y_i), \forall i = 1 : D\}$$

$$\mathbf{x}_i = [x_1, x_2]^T \in \mathbb{R}^2 : \begin{cases} \text{Covariates} \\ \text{attributes} \\ \text{regressors} \end{cases}$$

$$y_i \in \mathbb{R}^1 : \text{Observation}$$

Model

$$g(\mathbf{x}) \equiv \text{fct}(\mathbf{x}) \in \mathbb{R}^1$$



General case: $\mathbf{x}_i = [x_1, x_2, \dots, x_x]^T \in \mathbb{R}^x, y_i \in \mathbb{R}^1, g(\mathbf{x}_i) \in \mathbb{R}^1$

Mathematical formulation – general case >1D

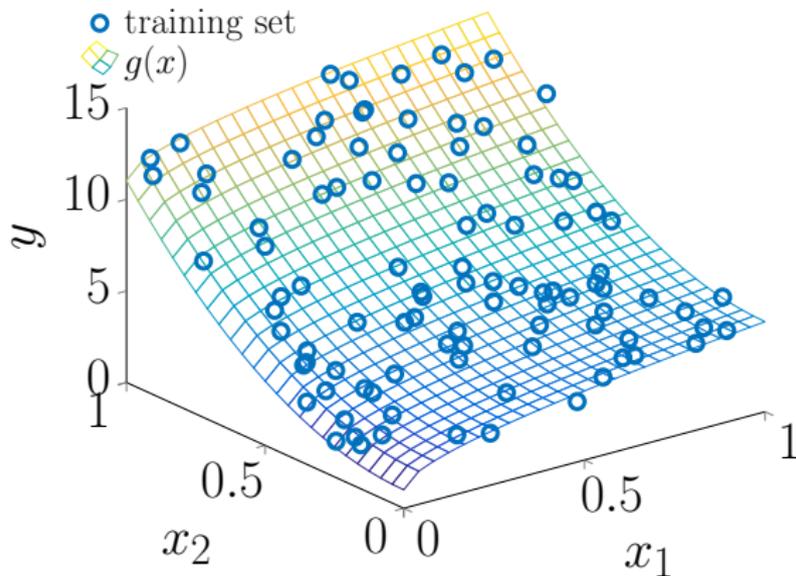
$$g(\mathbf{x}) = b_0 + b_1\phi_1(x_1) + b_2\phi_2(x_2) + \dots = \mathbf{X}\mathbf{b}$$

For the entire dataset

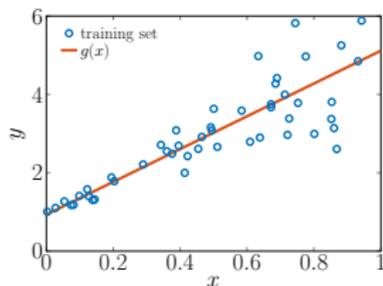
$$\mathbf{b} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_B \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_D \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & \phi_1(x_{1,1}) & \phi_2(x_{1,2}) & \cdots \\ 1 & \phi_1(x_{2,1}) & \phi_2(x_{2,2}) & \cdots \\ \vdots & \vdots & \vdots & \ddots \\ 1 & \phi_1(x_{D,1}) & \phi_2(x_{D,2}) & \cdots \end{bmatrix}$$

Example #4 

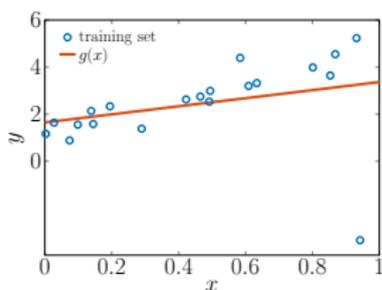
$$g(\mathbf{x}) = b_0 + b_1 x_1^{1/2} + b_2 x_2^2$$



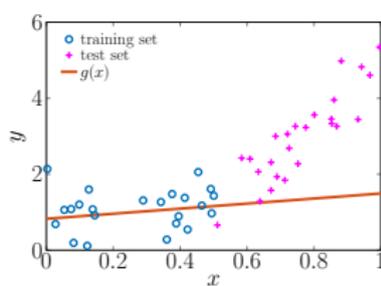
Limitations



Heteroscedasticity



Outliers



Extrapolation

Limitations

- ▶ Assumes **homoscedastic** errors
- ▶ Sensitive to **outliers**
- ▶ Does not differentiate between **interpolating and extrapolating**
- ▶ **⚠** The performance depends on the capacity to **hand-pick the correct transformation functions** (Difficult when $D > 1$)

Summary

In short: Linear regression is simple and quick, yet it is **outperformed by modern techniques** such as **Gaussian Process Regression**

Note: Cross-validation is not limited to linear regression and can be employed with any regression technique

Example – T° distribution along a pipeline



Given that we know the pipeline temperature to be $y = 8^\circ\text{C}$ at $x = 52\text{ km}$.

What is the temperature at $x = 52.1\text{ km}$?

What is the temperature at $x = 70\text{ km}$?



We will take advantage of the **conditional dependency between a system response y and covariate values x**

 Section Outline

Gaussian Process Regression

- 3.1 Introduction
 - 3.2 Definition
 - 3.3 Correlation functions
 - 3.4 Updating a GP using exact observations
 - 3.5 Updating a GP using noise-contaminated observations
 - 3.6 Multiple covariates: $\mathbf{x} = [x_1, x_2, \dots, x_X]^T$
 - 3.7 Parameter estimation
 - 3.8 Evaluating GPR's predictive capacity
 - 3.9 Code
-

Definition

Given a **system response** so that

$$\underbrace{g_i}_{\text{observation}} = \overbrace{g(\mathbf{x}_i)}^{\text{reality}}, \quad \underbrace{\mathbf{x}_i = [x_1, x_2, \dots, x_N]^T}_{\text{covariates}}$$

Data: $\mathcal{D} = \{(\mathbf{x}_i, g_i), \forall i = 1 : D\}$

$$\mathbf{g} = [g_1, g_2, \dots, g_D]^T, \quad \mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_D]_{X \times D}$$

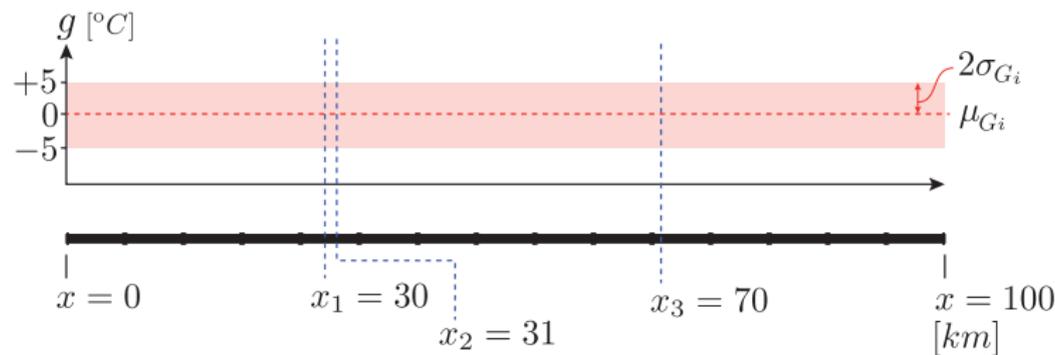
Gaussian process: $g(\mathbf{x}) : \mathbf{G} \sim \mathcal{N}(g(\mathbf{x}); \boldsymbol{\mu}_{\mathbf{G}}, \boldsymbol{\Sigma}_{\mathbf{G}})$

Set of discrete Gaussian random variables for which the pairwise correlation between G_i and G_j is a function of the distance between attributes

$$[\boldsymbol{\Sigma}_{\mathbf{G}}]_{ij} = \rho(x_i, x_j) \sigma_{G_i} \cdot \sigma_{G_j}, \quad \rho(x_i, x_j) = \text{fct}(|x_i - x_j|)$$

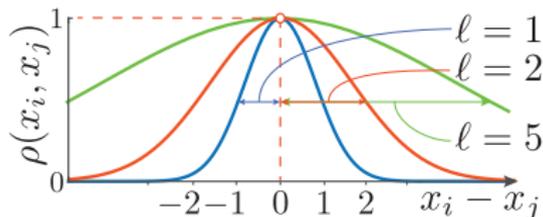
Example - T° distribution along a pipeline



Example - T° distribution along a pipeline

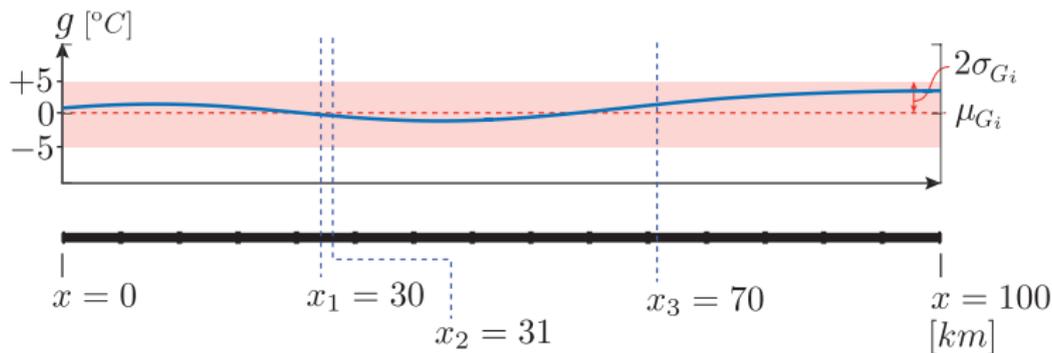
Prior knowledge: $\mu_{G_i} = \mu_G = 0^\circ\text{C}$, $\sigma_{G_i} = \sigma_G = 2.5^\circ\text{C}$

Gaussian correlation function: $\rho(x_i, x_j) = \exp\left(-\frac{1}{2}\frac{(x_i - x_j)^2}{\ell^2}\right)$



$$\text{for } \ell = 25 \text{ km} \begin{cases} \rho(x_1, x_2) = 0.999 \\ \rho(x_1, x_3) = 0.278 \\ \rho(x_2, x_3) = 0.296 \end{cases}$$

Example - T° distribution along a pipeline [Python]

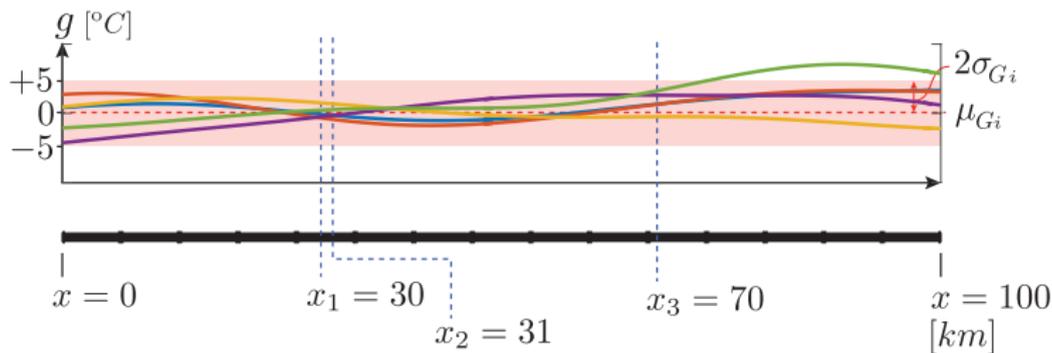


Pour $\mathbf{x} = [0, 1, \dots, 100]^\top$

$$\overbrace{\mathbf{G} \sim \mathcal{N}(g(\mathbf{x}); \boldsymbol{\mu}_{\mathbf{G}}, \boldsymbol{\Sigma}_{\mathbf{G}}), [\boldsymbol{\Sigma}_{\mathbf{G}}]_{ij} = \rho(x_i, x_j)\sigma_{G_i}^2}^{\text{Prior knowledge}}$$

$\mathbf{g}_i(\mathbf{x})$: realizations of our prior knowledge \mathbf{G}

Example - T° distribution along a pipeline [🐍]



Pour $\mathbf{x} = [0, 1, \dots, 100]^\top$

$$\overbrace{\mathbf{G} \sim \mathcal{N}(g(\mathbf{x}); \boldsymbol{\mu}_{\mathbf{G}}, \boldsymbol{\Sigma}_{\mathbf{G}}), [\boldsymbol{\Sigma}_{\mathbf{G}}]_{ij} = \rho(x_i, x_j)\sigma_{G_i}^2}^{\text{Prior knowledge}}$$

$\mathbf{g}_i(\mathbf{x})$: realizations of our prior knowledge \mathbf{G}

Updating a GP using exact observations

Given $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{g}_i), i = 1, \dots, D\}$ a set of D observations and \mathbf{x}_* a set of \mathbf{x}_* covariates for which we want to predict

$$f(\mathbf{g}_* | \mathbf{x}_*, \mathcal{D})$$

Reminder: Gaussian conditionals are also gaussian

$$\underbrace{\left\{ \begin{array}{c} \mathbf{G} \\ \mathbf{G}_* \end{array} \right\}, \boldsymbol{\mu} = \left\{ \begin{array}{c} \boldsymbol{\mu}_{\mathbf{G}} \\ \boldsymbol{\mu}_* \end{array} \right\}, \boldsymbol{\Sigma} = \left[\begin{array}{c|c} \boldsymbol{\Sigma}_{\mathbf{G}} & \boldsymbol{\Sigma}_{\mathbf{G}_*} \\ \hline \boldsymbol{\Sigma}_{\mathbf{G}_*}^{\top} & \boldsymbol{\Sigma}_* \end{array} \right]}_{\text{Prior knowledge}}$$

Prior knowledge

$$\begin{aligned} [\boldsymbol{\Sigma}_{\mathbf{G}}]_{ij} &= \rho(x_i, x_j) \sigma_G^2 \\ [\boldsymbol{\Sigma}_{\mathbf{G}_*}]_{ij} &= \rho(x_i, x_{*j}) \sigma_G^2 \\ [\boldsymbol{\Sigma}_*]_{ij} &= \rho(x_{*i}, x_{*j}) \sigma_G^2 \end{aligned}$$

$$f_{\mathbf{G}_* | \mathbf{x}_*, \mathcal{D}}(\mathbf{g}_* | \mathbf{x}_*, \mathcal{D}) = \mathcal{N}(\mathbf{g}_*; \boldsymbol{\mu}_* | \mathcal{D}, \boldsymbol{\Sigma}_* | \mathcal{D})$$

$$\boldsymbol{\mu}_* | \mathcal{D} \equiv \boldsymbol{\mu}_* | \mathcal{D} = \boldsymbol{\mu}_{\mathbf{G}_*} + \boldsymbol{\Sigma}_{\mathbf{G}_*}^{\top} \boldsymbol{\Sigma}_{\mathbf{G}}^{-1} (\mathbf{g} - \boldsymbol{\mu}_{\mathbf{G}})$$

$$\boldsymbol{\Sigma}_* | \mathcal{D} \equiv \boldsymbol{\Sigma}_* | \mathcal{D} = \boldsymbol{\Sigma}_* - \boldsymbol{\Sigma}_{\mathbf{G}_*}^{\top} \boldsymbol{\Sigma}_{\mathbf{G}}^{-1} \boldsymbol{\Sigma}_{\mathbf{G}_*}$$

Posterior knowledge

Example – Multivariate gaussian conditional

Soit $\underbrace{\mu_G = 0^\circ C, \sigma_G = 2.5^\circ C, \rho(x_1, x_2) = 0.8}_{\text{prior knowledge}}, \underbrace{g_2 = -2^\circ C}_{\text{observation}}$

$$f_{G_1 G_2}(g_1, g_2) = \mathcal{N}(\mu_G, \Sigma_G) \left\{ \begin{array}{l} \mu_G = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\ \Sigma_G = \begin{bmatrix} 2.5^2 & 0.8 \cdot 2.5^2 \\ 0.8 \cdot 2.5^2 & 2.5^2 \end{bmatrix} \end{array} \right.$$

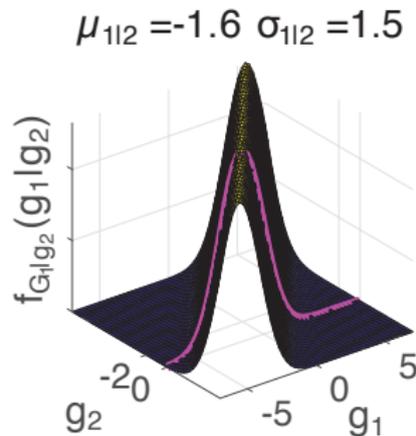
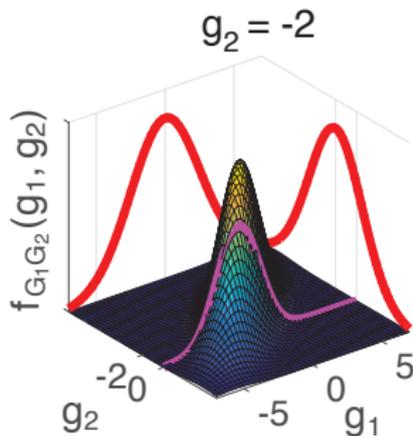
$$f_{G_1|g_2}(g_1|g_2) = \frac{f_{G_1 G_2}(g_1, g_2)}{f_{G_2}(g_2)} = \mathcal{N}(\mu_{1|2}, \sigma_{1|2}^2) \left\{ \begin{array}{l} \mu_{1|2} = \mu_1 + \rho \sigma_1 \frac{g_2 - \mu_2}{\sigma_2} \\ \sigma_{1|2} = \sigma_1 \sqrt{1 - \rho^2} \end{array} \right.$$

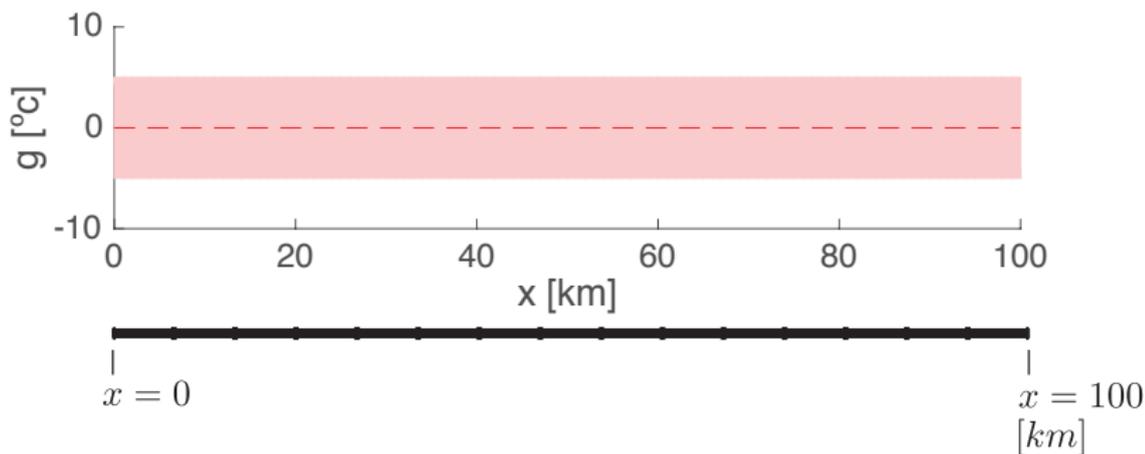
$$\mu_{1|2} = 0 + 0.8 \times 2.5 \frac{-2 - 0}{2.5} = -1.6^\circ C$$

$$\sigma_{1|2} = 2.5 \sqrt{1 - 0.8^2} = 1.5^\circ C$$

Example – Multivariate gaussian conditional

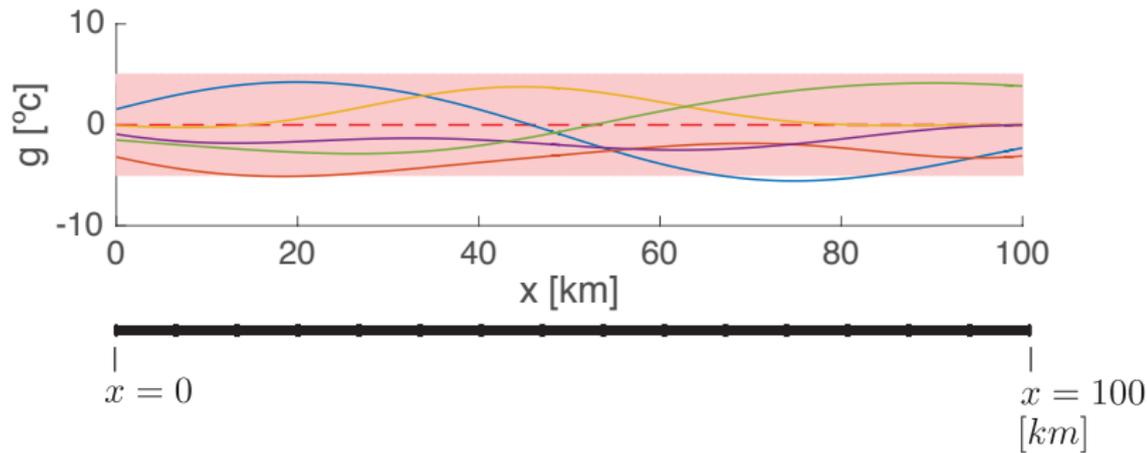
$$f_{G_1|g_2}(g_1|g_2) = \frac{f_{G_1 G_2}(g_1, g_2)}{f_{G_2}(g_2)} = \mathcal{N}(\mu_{1|2}, \sigma_{1|2}^2) \begin{cases} \mu_{1|2} = \mu_1 + \rho\sigma_1 \frac{y_2 - \mu_2}{\sigma_2} \\ \sigma_{1|2} = \sigma_1 \sqrt{1 - \rho^2} \end{cases}$$



Example - T° distribution along a pipeline [🐍]

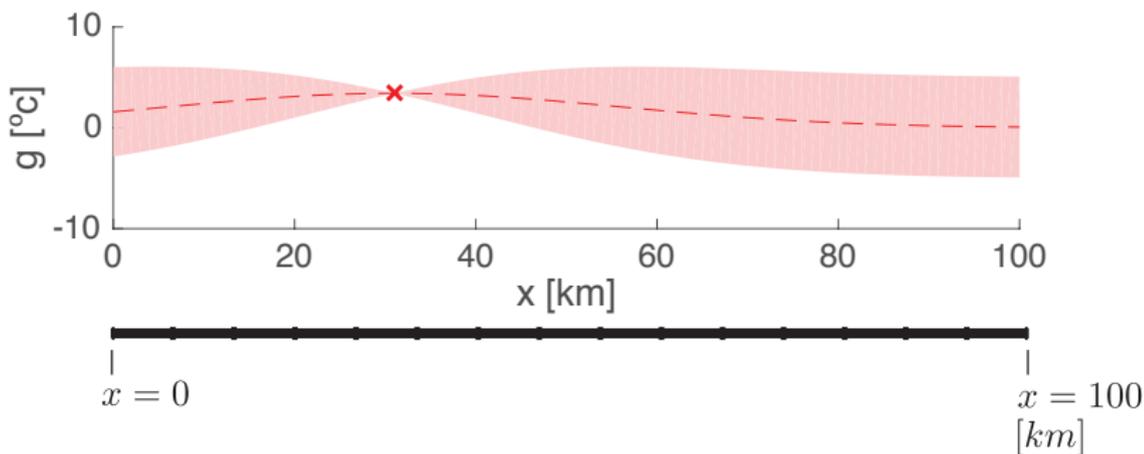
Observations: $\mathcal{D} = \left\{ \left(\underbrace{[\]}_{\mathbf{x}}, \underbrace{[\]}_{\mathbf{g}} \right) \right\}$

Example - T° distribution along a pipeline

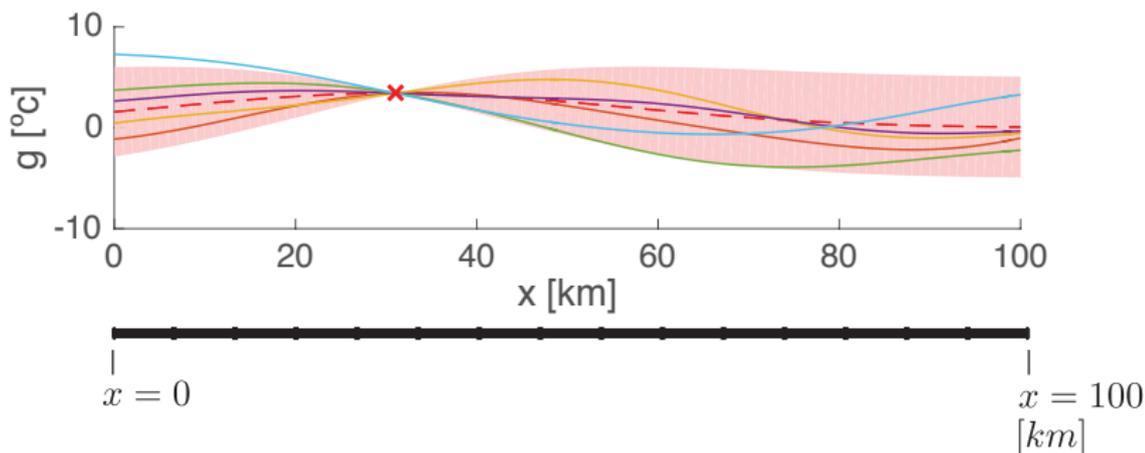


Observations: $\mathcal{D} = \left\{ \left(\underbrace{[[]}_{\mathbf{x}}, \underbrace{[[]}_{\mathbf{g}} \right) \right\}$

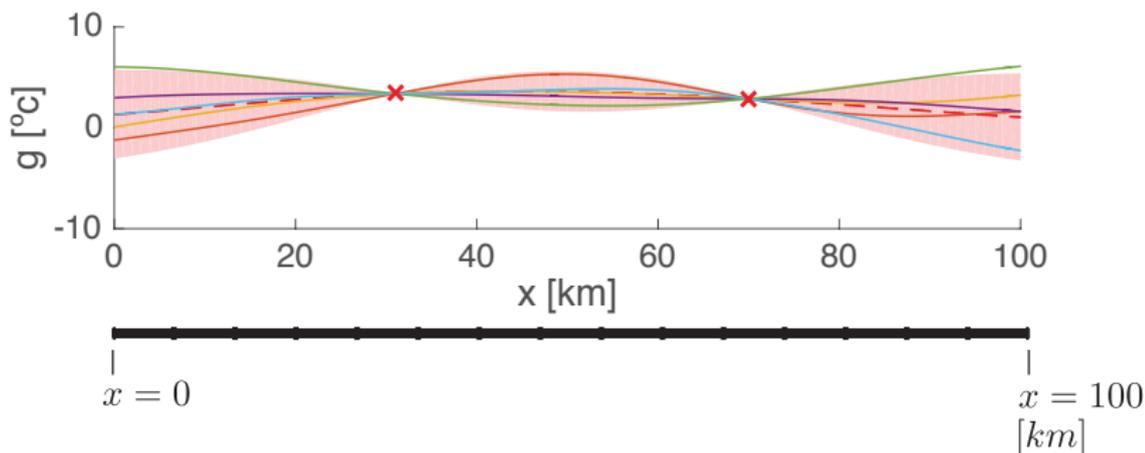
Example - T° distribution along a pipeline [🐍]



Observations: $\mathcal{D} = \left\{ \left(\underbrace{[31]}_x, \underbrace{[4.13]}_g \right) \right\}$

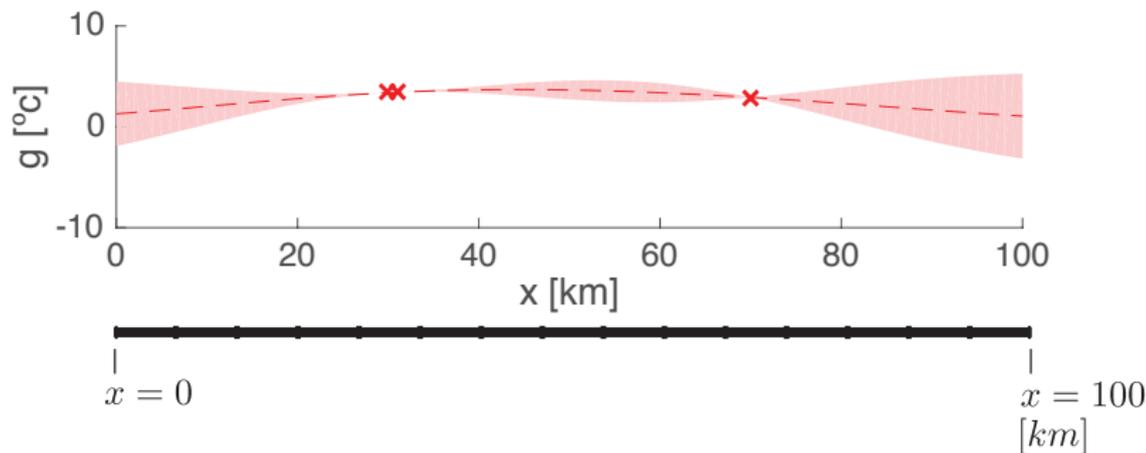
Example - T° distribution along a pipeline [🐍]

Observations: $\mathcal{D} = \left\{ \left(\underbrace{[31]}_x, \underbrace{[4.13]}_g \right) \right\}$

Example - T° distribution along a pipeline [🐍]

Observations: $\mathcal{D} = \left\{ \left(\underbrace{[31, 70]}_x, \underbrace{[4.13, 3.62]}_g \right) \right\}$

Example - T° distribution along a pipeline [🐍]



Observations: $\mathcal{D} = \left\{ \left(\underbrace{[31, 70, 30]}_x, \underbrace{[4.13, 3.62, 3.71]}_y \right) \right\}$

Noise-contaminated observations

Given a **system response** contaminated by **measurement noise** so that

$$\underbrace{Y_i}_{\text{observation}} = \underbrace{g(\mathbf{x}_i)}_{\text{reality}} + \underbrace{V}_{\text{measurement error}}, \quad V \sim \mathcal{N}(v; 0, \sigma_V^2)$$

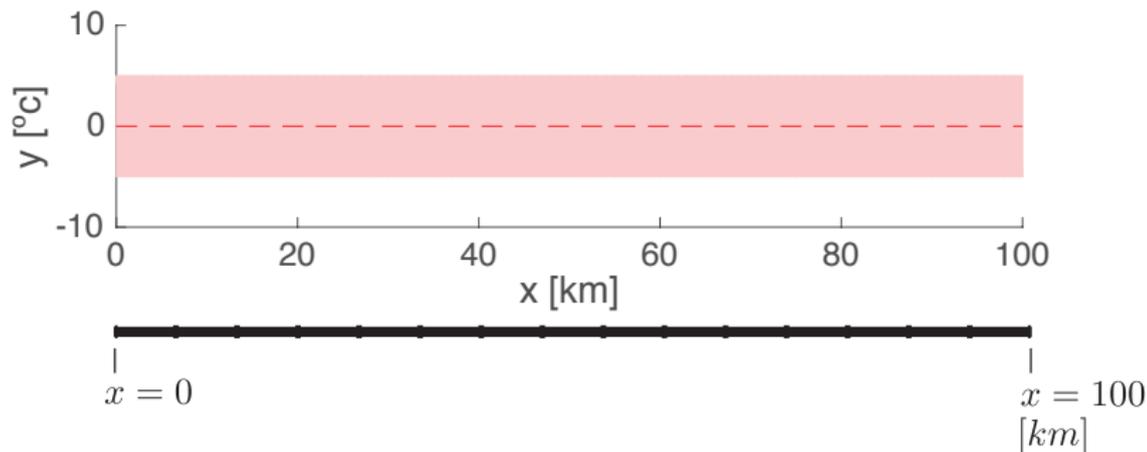
$$[\boldsymbol{\Sigma}_Y]_{ij} = \rho(\mathbf{x}_i, \mathbf{x}_j) \sigma_G^2 + \sigma_V^2 \delta_{ij}, \quad \delta_{ij} = 1 \text{ if } i = j$$

$$\underbrace{\left\{ \begin{array}{c} \mathbf{Y} \\ \mathbf{G}_* \end{array} \right\}, \quad \boldsymbol{\mu} = \left\{ \begin{array}{c} \boldsymbol{\mu}_Y \\ \boldsymbol{\mu}_{\mathbf{G}_*} \end{array} \right\}, \quad \boldsymbol{\Sigma} = \left[\begin{array}{c|c} \boldsymbol{\Sigma}_Y & \boldsymbol{\Sigma}_{Y*} \\ \hline \boldsymbol{\Sigma}_{Y*}^T & \boldsymbol{\Sigma}_* \end{array} \right]}_{\text{prior knowledge}}$$

$$f_{\mathbf{G}_*|\mathbf{x}_*, D}(\mathbf{g}_*|\mathbf{x}_*, D) = \mathcal{N}(\mathbf{g}_*; \boldsymbol{\mu}_{*|D}, \boldsymbol{\Sigma}_{*|D})$$

$$\underbrace{\boldsymbol{\mu}_{*|D} = \boldsymbol{\mu}_{\mathbf{G}_*} + \boldsymbol{\Sigma}_{Y*}^T \boldsymbol{\Sigma}_Y^{-1} (\mathbf{y} - \boldsymbol{\mu}_Y), \quad \boldsymbol{\Sigma}_{*|D} = \boldsymbol{\Sigma}_* - \boldsymbol{\Sigma}_{Y*}^T \boldsymbol{\Sigma}_Y^{-1} \boldsymbol{\Sigma}_{Y*}}_{\text{posterior knowledge}}$$

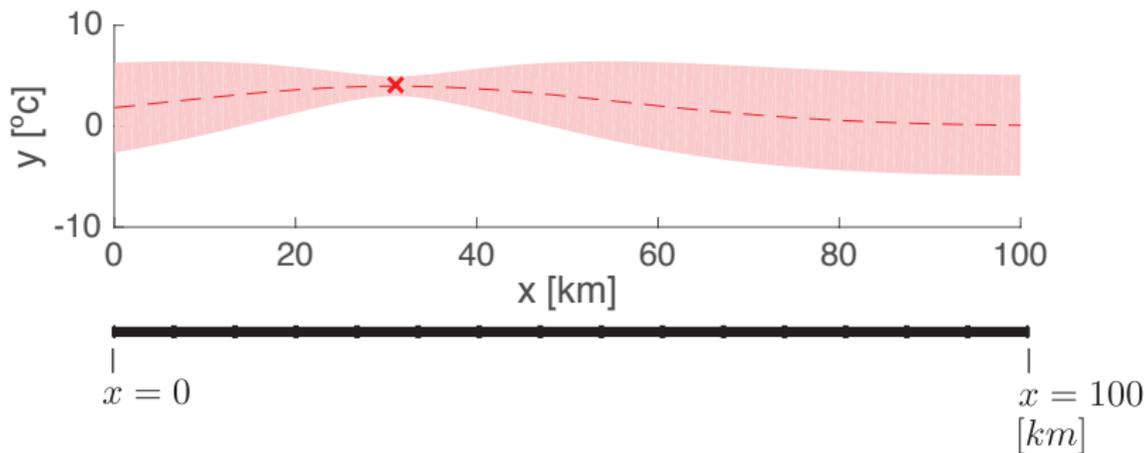
Example - T° distribution ($V \neq 0$)



Measurement error: $\sigma_V = 0.5$

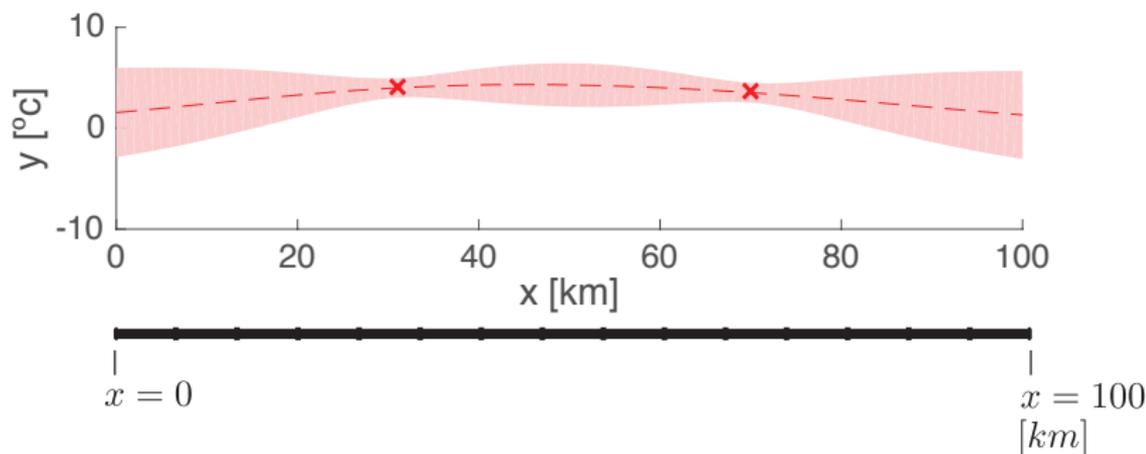
Observations: $\mathcal{D} = \{(\underbrace{[\]}_x, \underbrace{[\]}_y)\}$

Example - T° distribution ($V \neq 0$) [Python]



Measurement error: $\sigma_V = 0.5$

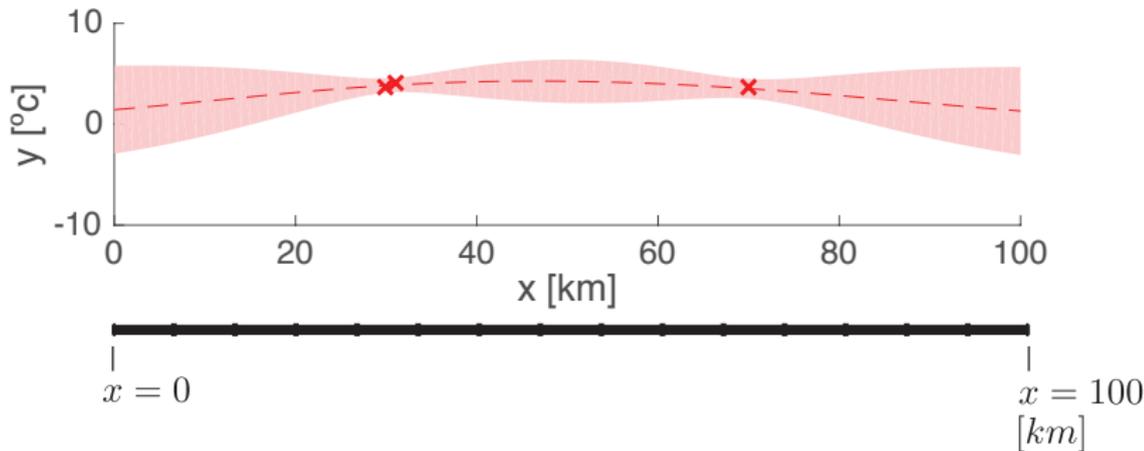
Observations: $\mathcal{D} = \{(\underbrace{[31]}_x, \underbrace{[3.42]}_y)\}$

Example - T° distribution ($V \neq 0$) [🐍]

Measurement error: $\sigma_V = 0.5$

Observations: $\mathcal{D} = \left\{ \left(\underbrace{[31, 70]}_x, \underbrace{[3.42, 2.92]}_y \right) \right\}$

Example - T° distribution ($V \neq 0$)



Measurement error: $\sigma_V = 0.5$

Observations: $\mathcal{D} = \{(\underbrace{[31, 70, 30]}_x, \underbrace{[3.42, 2.92, 3.38]}_y)\}$



Multiple covariates: $\mathbf{x} = [x_1, x_2, \dots, x_X]_i^\top$

Multiple covariates: $\mathbf{x} = [x_1, x_2, \dots, x_X]_i^\top$

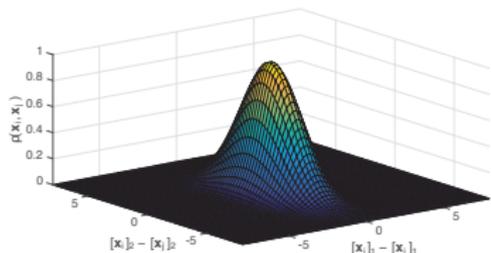
$$\underbrace{y_i}_{\text{observation}} = \overbrace{g(\mathbf{x}_i)}^{\text{reality}} + v_i, \quad \underbrace{\mathbf{x}_i}_{\text{covariates}} = [x_1, x_2, \dots, x_X]_i^\top, \quad \text{for } i = 1, \dots, D$$

$$\mathbf{y} = [y_1, y_2, \dots, y_D]^\top, \quad \mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_D]_{X \times D}, \quad \underbrace{\boldsymbol{\ell} = [\ell_1, \ell_2, \dots, \ell_X]_{X \times 1}^\top}_{\text{one } \ell_k \text{ per covariate } \mathbf{x}_k}$$

$$\begin{aligned} \rho(\mathbf{x}_i, \mathbf{x}_j) &= \exp\left(-\frac{1}{2} \sum_{k=1}^X \frac{([\mathbf{x}_i]_k - [\mathbf{x}_j]_k)^2}{\ell_k^2}\right) \\ &= \exp\left(-\frac{1}{2} (\mathbf{x}_i - \mathbf{x}_j)^\top \text{diag}(\ell^2)^{-1} (\mathbf{x}_i - \mathbf{x}_j)\right) \end{aligned}$$

Multiple covariates: $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$ Multivariate correlation function 

$$\begin{aligned}\rho(\mathbf{x}_i, \mathbf{x}_j) &= \exp\left(-\frac{1}{2} \sum_{k=1}^X \frac{([\mathbf{x}_i]_k - [\mathbf{x}_j]_k)^2}{\ell_k^2}\right) \\ &= \exp\left(-\frac{1}{2}(\mathbf{x}_i - \mathbf{x}_j)^T \text{diag}(\ell)^{-2}(\mathbf{x}_i - \mathbf{x}_j)\right)\end{aligned}$$

Raw datasi $\ell_k \gg \sigma_{X_k}$: $\text{Imp}(X_k \rightarrow y) \approx 0$ si $\ell_k < \sigma_{X_k}$: $\text{Imp}(X_k \rightarrow y) \gg 0$ **Standardized data** $(x_k - \mu_{X_k})/\sigma_{X_k}$ si $\ell_k \gg 1$: $\text{Imp}(X_k \rightarrow y) \approx 0$ si $\ell_k < 1$: $\text{Imp}(X_k \rightarrow y) \gg 0$ [\[CIV_ML/M14_multivariate_corr_fct.py\]](#)

Parameter estimation σ_G & ℓ

Up to now, we made the hypothesis that we knew σ_G and $\ell = [\ell_1, \ell_2, \dots, \ell_x]^T$.

In practice, we need to learn **hyperparameters** $\theta = [\sigma_G, \ell]^T$ using training data $\mathcal{D} = \{\mathcal{D}_x, \mathcal{D}_y\} = \{(\mathbf{x}_i, y_i), i = 1, \dots, D\}$.

Bayes rule:

$$\underbrace{f(\theta | \mathcal{D}_x, \mathcal{D}_y)}_{\text{posterior}} = \frac{\overbrace{f(\mathcal{D}_y | \mathcal{D}_x, \theta)}^{\text{likelihood}} \cdot \overbrace{f(\theta)}^{\text{prior}}}{\underbrace{f(\mathcal{D}_y)}_{\text{cte.}}}$$

⚠ In practice, evaluating $f(\theta | \mathcal{D}_x, \mathcal{D}_y)$ is computationally expensive...

MLE approximation

$$\underbrace{f(\boldsymbol{\theta} | \mathcal{D}_x, \mathcal{D}_y)}_{\text{posterior}} = \frac{\overbrace{f(\mathcal{D}_y | \mathcal{D}_x, \boldsymbol{\theta})}^{\text{likelihood}} \cdot \overbrace{f(\boldsymbol{\theta})}^{\text{prior}}}{\underbrace{f(\mathcal{D}_y)}_{\text{cte.}}}$$

If $f(\boldsymbol{\theta}) = \text{cte.}$, $\forall \boldsymbol{\theta} \rightarrow f(\boldsymbol{\theta} | \mathcal{D}_x, \mathcal{D}_y) \propto f(\mathcal{D}_y | \mathcal{D}_x, \boldsymbol{\theta})$

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} f(\mathcal{D}_y | \mathcal{D}_x, \boldsymbol{\theta})$$

$$f(\mathcal{D}_y | \mathcal{D}_x, \boldsymbol{\theta}) = \mathcal{N}(\mathcal{D}_y; \boldsymbol{\mu}_Y, \boldsymbol{\Sigma}_Y)$$

MLE approximation (cont.)

When the size of $\Sigma_{\mathbf{Y}}$ increases, $f(\mathcal{D}_y|\mathcal{D}_x, \boldsymbol{\theta}) \rightarrow 0$ causing **zero underflow** \triangle .

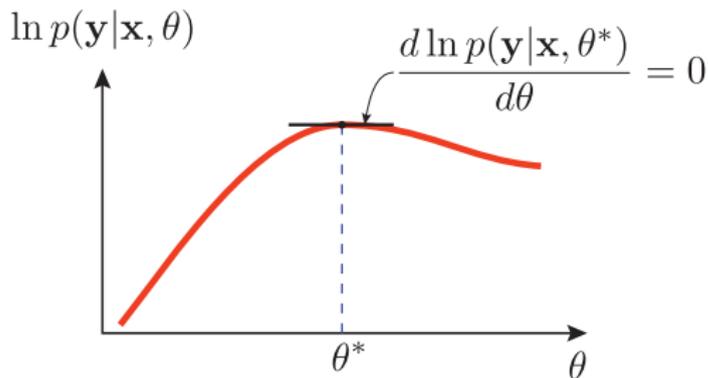
Solution:

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \overbrace{f(\mathcal{D}_y|\mathcal{D}_x, \boldsymbol{\theta})}^{\text{likelihood}} \equiv \arg \max_{\boldsymbol{\theta}} \overbrace{\log f(\mathcal{D}_y|\mathcal{D}_x, \boldsymbol{\theta})}^{\text{log-likelihood}}$$

$f(\mathcal{D}_y|\mathcal{D}_x, \boldsymbol{\theta})$'s maximum corresponds to $\log f(\mathcal{D}_y|\mathcal{D}_x, \boldsymbol{\theta})$'s maximum and there are no more zero underflow \triangle .

$$\begin{aligned} \log f(\mathcal{D}_y|\mathcal{D}_x, \boldsymbol{\theta}) &= \log \mathcal{N}(\mathcal{D}_y; \boldsymbol{\mu}_{\mathbf{Y}}, \Sigma_{\mathbf{Y}}) \\ &= -\frac{1}{2}(\mathcal{D}_y - \boldsymbol{\mu}_{\mathbf{Y}})^\top \Sigma_{\mathbf{Y}}^{-1}(\mathcal{D}_y - \boldsymbol{\mu}_{\mathbf{Y}}) - \frac{1}{2} \log |\Sigma_{\mathbf{Y}}| - \frac{D}{2} \log 2\pi \end{aligned}$$

Maximization



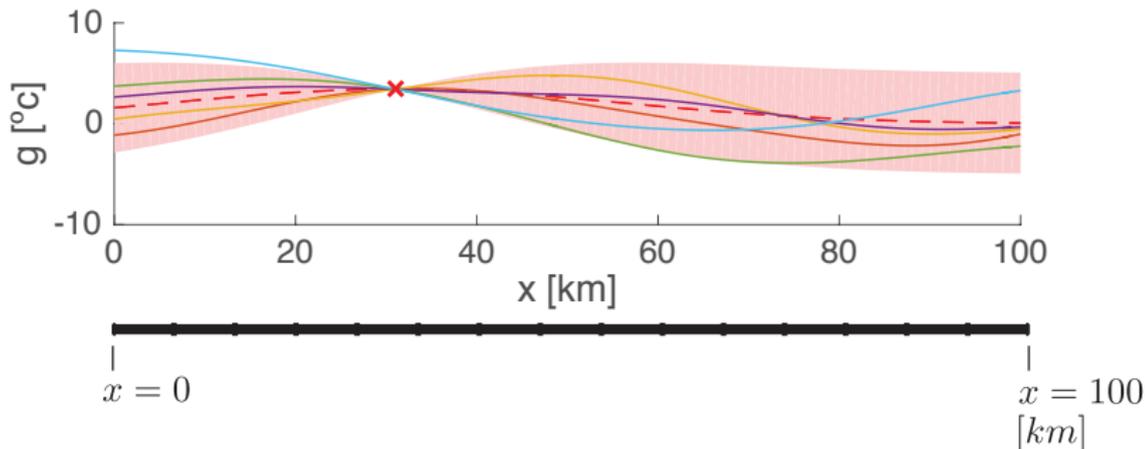
How to find θ^* ?

$$\frac{\partial}{\partial \theta_j} \log f(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{2} \mathbf{y}^\top \boldsymbol{\Sigma}_{\mathbf{Y}}^{-1} \frac{\partial \boldsymbol{\Sigma}_{\mathbf{Y}}}{\partial \theta_j} \boldsymbol{\Sigma}_{\mathbf{Y}}^{-1} \mathbf{y} - \frac{1}{2} \text{tr}(\boldsymbol{\Sigma}_{\mathbf{Y}}^{-1} \frac{\partial \boldsymbol{\Sigma}_{\mathbf{Y}}}{\partial \theta_j}) = 0$$

We employ iterative optimization techniques to find θ^*
(such as **Newton-Raphson**)

Efficient algorithms are already implemented in **GPy/GPML**.

Interactive example - T° distribution along a pipeline



Evaluating GPR's predictive capacity

We should employ **N-fold cross-validation** to estimate the GPR predictive capacity

1. set-up the GPR model
2. use CV to estimate hyper-parameters $\{\ell, \sigma_G, \sigma_V\}$ the current training set and then predict only data points that are not used in the training set
3. repeat 2. for each for the N cross-validation fold
4. plot the histogram of normalized prediction errors: $\frac{\mu_{\mathbf{Y}|\mathcal{D},i} - y_i}{\sigma_{\mathbf{Y}|\mathcal{D},i}}$
5. plot a scatter of predicted and observed values

GPy library - Option #1 Raw data

```

1  import GPy
2  import numpy as np
3
4  # Data
5  x_obs, y_obs, x_pred #Covariates, observed values, pred. covariates
6  nb_dim = x_obs.shape[1] #Number of covariates
7
8  # GPy
9  l, var_G, var_V = 10, 2, 0.1 #Lengthscale, Process/observation err. variance
10 kernel = GPy.kern.RBF(input_dim = nb_dim, variance = var_G, lengthscale = 1)
11 lik = GPy.likelihoods.Gaussian()
12 model = GPy.models.GPRegression(x_obs, y_obs, kernel = kernel ,ARD = True)
13 model.Gaussian_noise.variance = var_V
14 model.optimize()
15 print(model)
16
17 # Estimated parameter values
18 print(kernel.lengthscale)
19 print(kernel.variance)
20 print(model.log_likelihood())
21
22 # Predictions
23 m_G, var_G = model.predict(x_pred, full_cov = False) #Diagonal covariance
24
25 # Plotting
26 fig = model.plot()
27 GPy.plotting.show(fig)

```

GPy library - Option #2 Standardized data

```

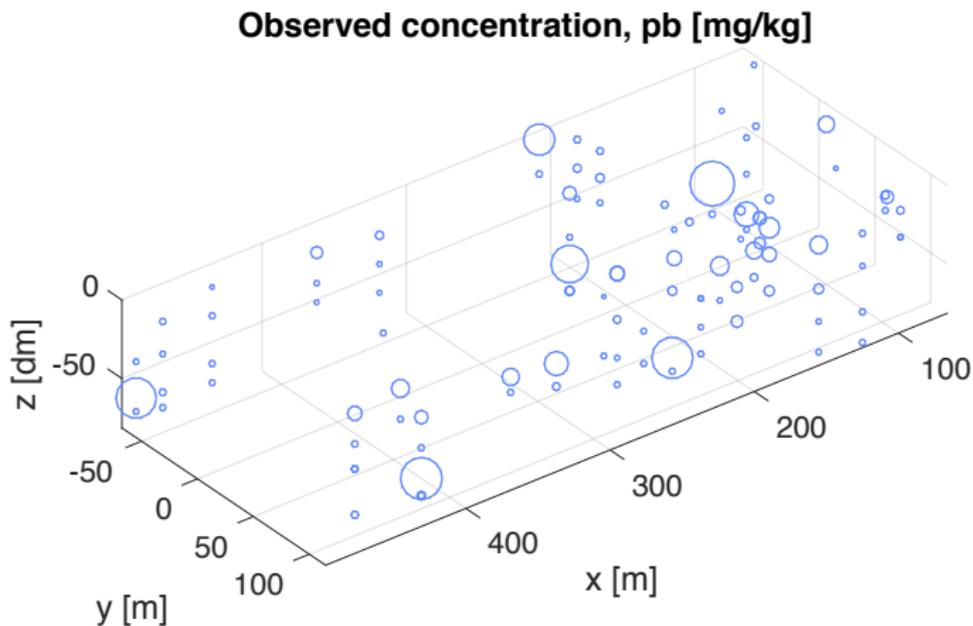
1 import GPy
2 import numpy as np
3
4 # Data
5 x_obs, y_obs, x_pred #Covariates, observed values, pred. covariates
6 nb_dim = x_obs.shape[1] #Number of covariates
7
8 # Standardization
9 mx = np.mean([x_obs],axis=1)
10 sx = np.std([x_obs],axis=1)
11 my = np.mean([y_obs])
12 sy = np.std([y_obs])
13 x_std = (x_obs-mx)/sx
14 y_std = (y_obs-my)/sy
15 x_pred_std = (x_pred-mx)/sx
16
17 # GPy
18 kernel = GPy.kern.RBF(input_dim = nb_dim, ARD = True)
19 model = GPy.models.GPRegression(x_std, y_std, kernel = kernel)
20 model.optimize()
21
22 # Standardized Predictions
23 m_G_std, var_G_std = model.predict(x_pred_std, full_cov = False) #Diag. cov.
24
25 # Unstandardize predictions
26 m_G = m_G_std*sy+my
27 var_G = var_G_std*sy**2

```

Section Outline

Examples

- 4.1 Soil contamination characterization
- 4.2 UCI - Concrete Slump

Soil contamination characterization 

[CIV_ML/Soil_contamination_data.fig]

Problem set-up

Covariates and observations:

$$\underbrace{\mathbf{l}_i^{\{Y\}} = [x, y, z]_i^\top}_{\text{geodesic coordinates}}$$

$$\mathcal{D} = \{\mathcal{D}_x, \mathcal{D}_y\} = \{(\mathbf{l}_i^{\{Y\}}, \mathbf{y}_i), i = 1 : 116\}$$

Observation model:

true contaminant $[\]$ in log space

$$\log \underbrace{y_i}_{\text{observation}} = \underbrace{g(\mathbf{l}_i^{\{Y\}})}_{\text{meas.}} + \underbrace{v_i}_{\text{error in log space}}, \quad v : V \sim \mathcal{N}(0, \sigma_V^2)$$

Problem set-up (cont.)

Prior:

$$\boldsymbol{\mu} = \begin{Bmatrix} \boldsymbol{\mu}_Y \\ \boldsymbol{\mu}_* \end{Bmatrix} = \begin{Bmatrix} \mathbf{0} \\ \mathbf{0} \end{Bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_Y & \boldsymbol{\Sigma}_{Y^*} \\ \boldsymbol{\Sigma}_{Y^*}^\top & \boldsymbol{\Sigma}_* \end{bmatrix}$$

$$[\boldsymbol{\Sigma}_Y]_{ij} = \rho(\mathbf{l}_i^{\{Y\}}, \mathbf{l}_j^{\{Y\}}) \sigma_G^2 + \sigma_V^2 \delta_{ij}, \quad \delta_{ij} = 1 \text{ if } i = j, \text{ else } \delta_{ij} = 0$$

$$[\boldsymbol{\Sigma}_*]_{kl} = \rho(\mathbf{l}_k^{\{G\}}, \mathbf{l}_l^{\{G\}}) \sigma_G^2$$

$$[\boldsymbol{\Sigma}_{Y^*}]_{ik} = \rho(\mathbf{l}_i^{\{Y\}}, \mathbf{l}_k^{\{G\}}) \sigma_G^2$$

$$\rho(\mathbf{l}_i, \mathbf{l}_j) = \exp\left(-\frac{1}{2}(\mathbf{l}_i - \mathbf{l}_j)^\top \text{diag}(\ell^2)^{-1}(\mathbf{l}_i - \mathbf{l}_j)\right)$$

Posterior knowledge:

$$\boldsymbol{\mu}_{*|\mathcal{D}} = \boldsymbol{\mu}_* + \boldsymbol{\Sigma}_{Y^*}^\top \boldsymbol{\Sigma}_Y^{-1}(\mathbf{y} - \boldsymbol{\mu}_Y)$$

$$\boldsymbol{\Sigma}_{*|\mathcal{D}} = \boldsymbol{\Sigma}_* - \boldsymbol{\Sigma}_{Y^*}^\top \boldsymbol{\Sigma}_Y^{-1} \boldsymbol{\Sigma}_{Y^*}$$

Problem set-up (cont.)

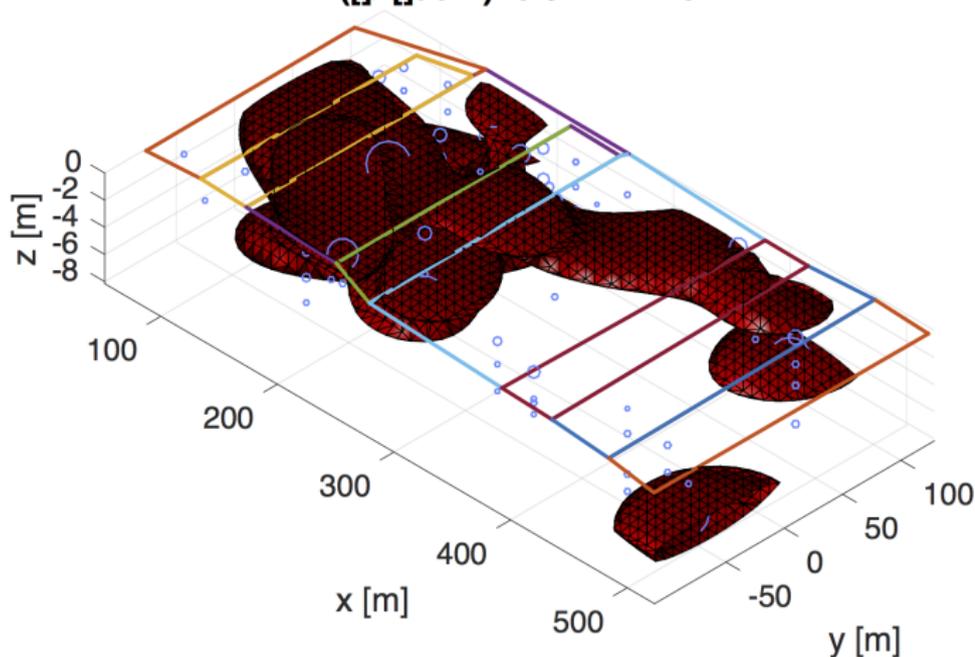
Hyper-parameters estimation:

$$\boldsymbol{\theta} = \{l_x, l_y, l_z, \sigma_G, \sigma_V\}$$

$$\begin{aligned} \boldsymbol{\theta}^* &= \arg \max_{\boldsymbol{\theta}} \overbrace{\log p(\mathcal{D}_y | \mathcal{D}_x, \boldsymbol{\theta})}^{\text{log-likelihood}} \\ &= \left\{ \begin{array}{l} l_x = 56.4 \text{ m} \\ l_y = 53.7 \text{ m} \\ l_z = 0.64 \text{ m} \\ \sigma_G = 2.86 \\ \sigma_V = 0.18 \end{array} \right\} \end{aligned}$$

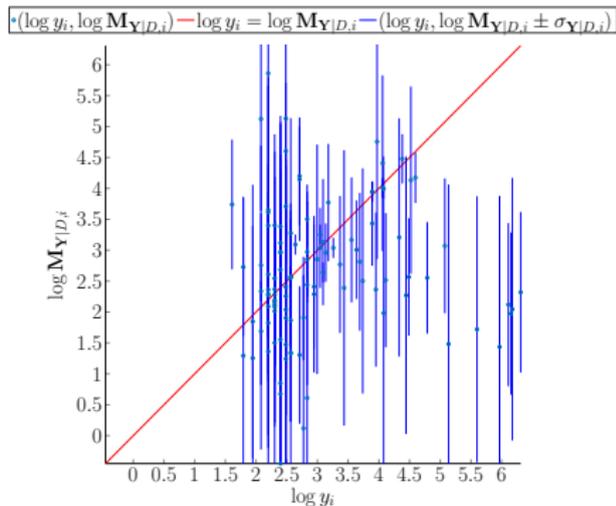
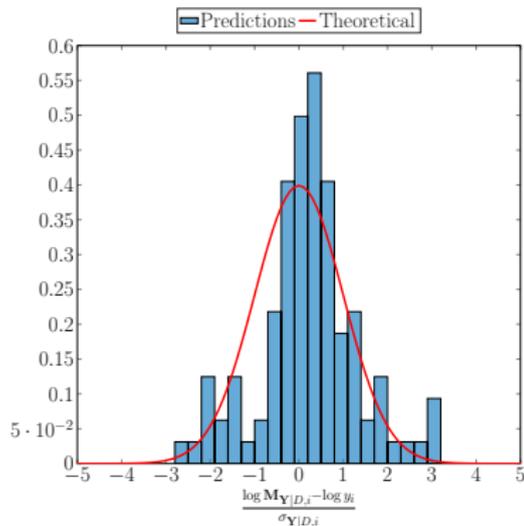
Results – contaminant concentration

$\Pr([\cdot] > [\cdot]_{adm}) = 0.5 \mid 1D \mid 10m$



[CIV_ML/Soil_contamination_concentration.fig//Soil_contaimnation/GPR_SC_main.m]

Results – Prediction quality



⚠ Overall, the predictive capacity is poor → need more data

UCI - Concrete Slump dataset



[About](#) [Citation Policy](#) [Donate a Data Set](#)
[Contact](#)

Repository Web [CURL](#)

[View ALL Data Sets](#)

Concrete Slump Test Data Set

Download: [Data Folder](#) [Data Set Description](#)

Abstract: Concrete is a highly complex material. The slump flow of concrete is not only determined by the water content, but that is also influenced by other concrete ingredients.



Data Set Characteristics:	Multivariate	Number of Instances:	103	Area:	Computer
Attribute Characteristics:	Real	Number of Attributes:	10	Date Donated	2009-04-30
Associated Tasks:	Regression	Missing Values?	N/A	Number of Web Hits:	132011

Source:

Donor: I-Cheng Yeh

Email: icyeh@i.cs.cmu.edu.tw

Institution: Department of Information Management, Chung-Hua University (Republic of China)

Other contact information: Department of Information Management, Chung-Hua University, Hsin Chu, Taiwan 30067, R.O.C.

<https://archive.ics.uci.edu/ml/datasets/Concrete+Slump+Test>

1 output [MPa]

- ▶ 28-day compressive strength

7 covariates [kg/m³]

- ▶ Cement
- ▶ Slag
- ▶ Fly ash
- ▶ Water
- ▶ Superplasticizer
- ▶ Coarse aggregate
- ▶ Fine aggregate

Problem set-up

Covariates and observations:

$$\mathbf{x}_i = [x_1, x_2, \dots, x_7]^T_i$$

$$\tilde{\mathcal{D}} = \{\tilde{\mathcal{D}}_x, \tilde{\mathcal{D}}_y\} = \left\{ \left(\frac{x_i - \mu_X}{\sigma_X}, \frac{y_i - \mu_Y}{\sigma_Y} \right), i = 1 : 103 \right\}$$

Observation model:

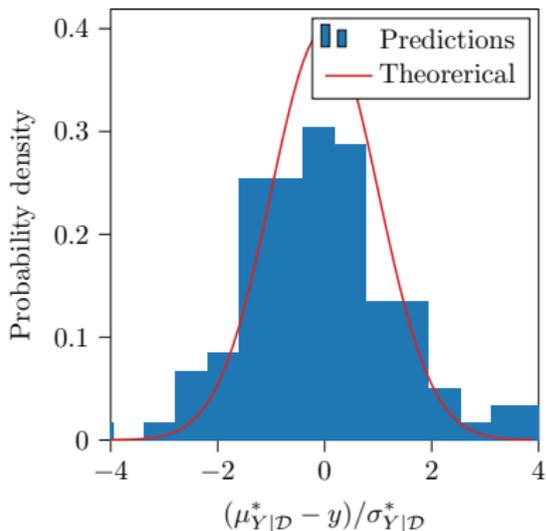
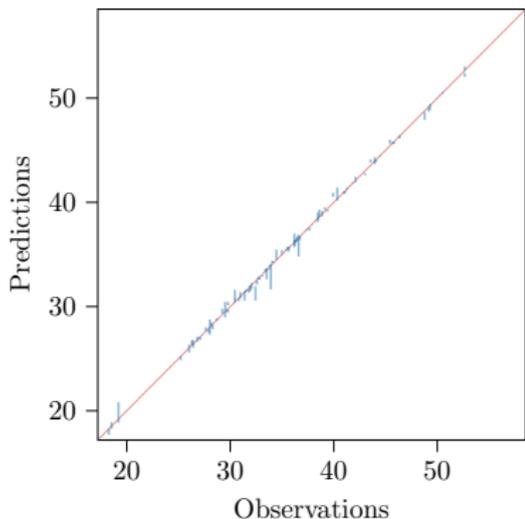
$$\underbrace{\tilde{y}_i}_{\text{observation}} = \underbrace{g(\tilde{\mathbf{x}}_i)}_{\text{true strength}} + \underbrace{v_i}_{\text{meas. error}}, v : V \sim \mathcal{N}(0, \sigma_V^2)$$

Hyper-parameters estimation

$$\theta^* = \arg \max_{\theta} \overbrace{\log p(\mathcal{D}_y | \mathcal{D}_x, \theta)}^{\text{log-likelihood}}$$

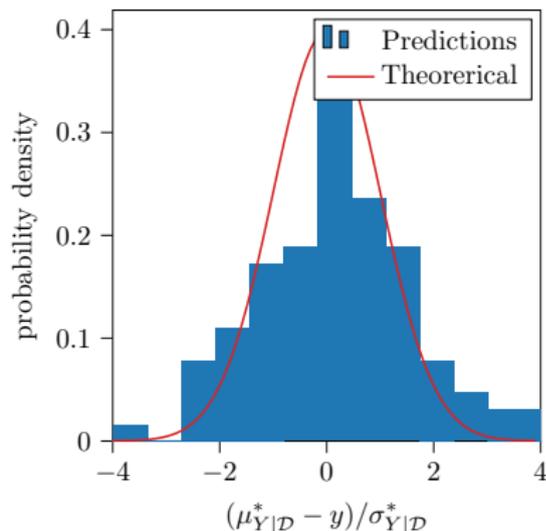
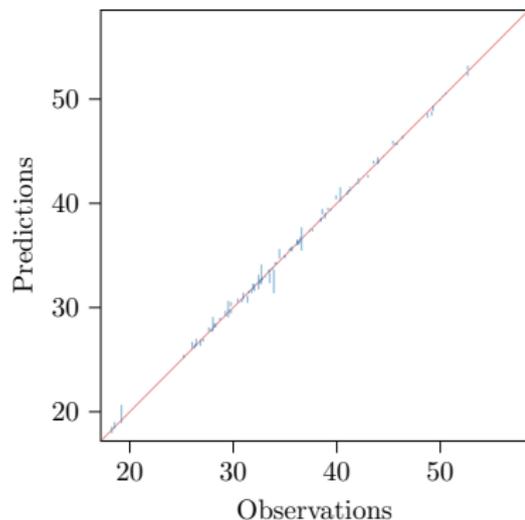
$$= \left\{ \begin{array}{l} \ell_1 = 5.4 \\ \ell_2 = 7.2 \\ \ell_3 = 4.8 \\ \ell_4 = 2.4 \\ \ell_5 = 18.5 \quad \triangle \\ \ell_6 = 28644 \quad \triangle \\ \ell_7 = 3.9 \\ \sigma_G = 2.7 \\ \sigma_V \approx 0 \end{array} \right.$$

Results 20-folds CV – $\mathbf{x} = [x_1, x_2, x_3, x_4, x_5, x_6, x_7]^T$



Mean square error (MSE)	0.118
Mean absolute error (MAE)	0.212
Mean error (ME)	0.016
Std error (RMSE)	0.343

Results 20-folds CV – $\mathbf{x} = [x_1, x_2, x_3, x_4, x_5, \dots, x_7]^T$



Mean square error (MSE)	0.079 < 0.118
Mean absolute error (MAE)	0.181
Mean error (ME)	0.029
Std error (RMSE)	0.279

[CIV_ML/GPR_example_Concrete.py]



Section Outline

Adv. topics

- 5.1 Prior knowledge – mean values
 - 5.2 Covariance functions
 - 5.3 Homo/Heteroscedasticity
 - 5.4 Other methods
-

Prior knowledge – mean values

So far, we have modelled our prior mean as either **0** or a **cte.**

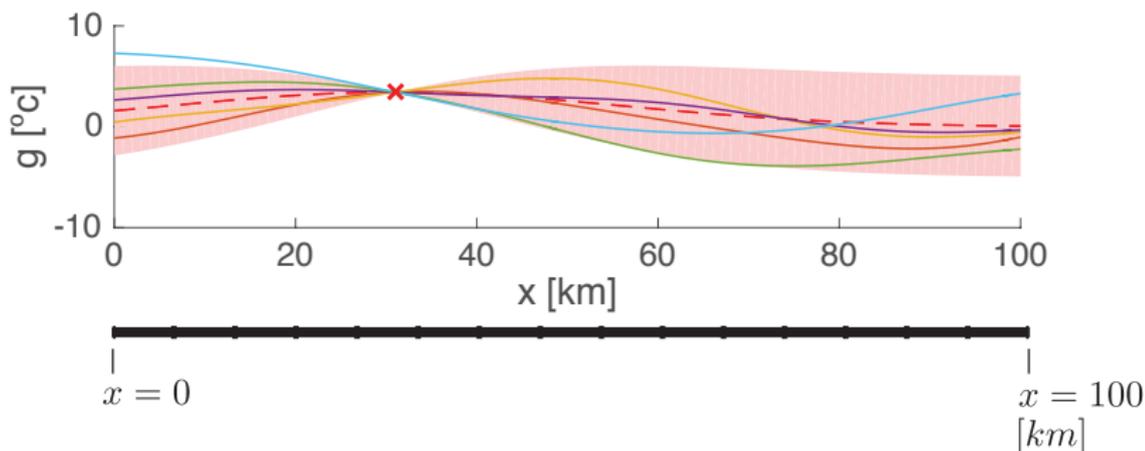
When predicting in the neighbourhood of observed data the prior knowledge has a negligible impact. **However, the impact of the prior increase as the distance from observed data increases**

Solution: **your prior knowledge can be parameterized by a function depending on covariates.**

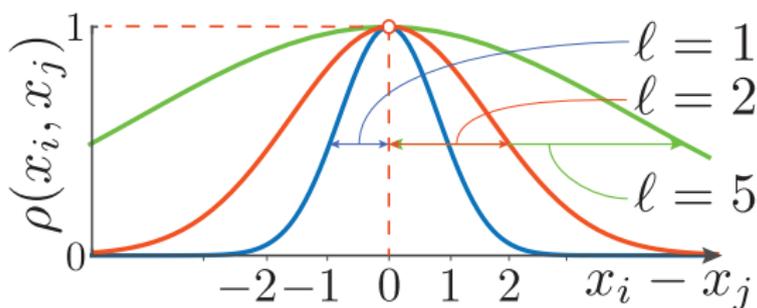
$$\text{e.g.: } m(x) = ax + b$$

Where $\{a, b\}$ are **hyper-parameters** that needs to be estimated using data.

Interactive example - T° distribution along a pipeline []



Square-exponential covariance function



$$\begin{aligned} \rho(\mathbf{x}_i, \mathbf{x}_j) &= \exp\left(-\frac{1}{2}(\mathbf{x}_i - \mathbf{x}_j)^\top \text{diag}(\ell)^{-2}(\mathbf{x}_i - \mathbf{x}_j)\right) \\ &= \exp\left(-\sum_{k=1}^x \frac{([x_i]_k - [x_j]_k)^2}{[\ell]_k}\right) \end{aligned}$$

Covariance function and cross-validation

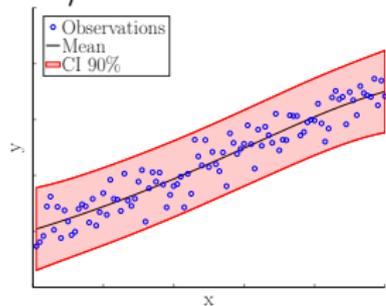
How to choose between two **covariance function structures**?

Cross-Validation (easy way)

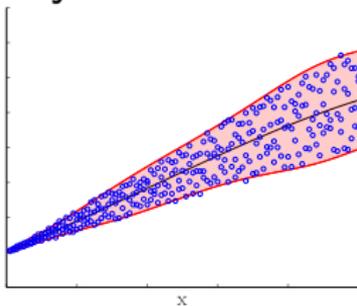
or

Bayesian Model selection (hard way)

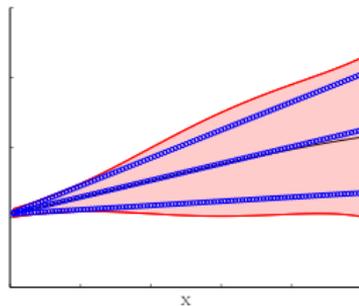
Homo/Heteroscedasticity



(a) Homoscedastic and independent



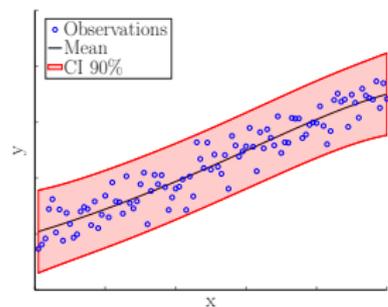
(b) Heteroscedastic and independent



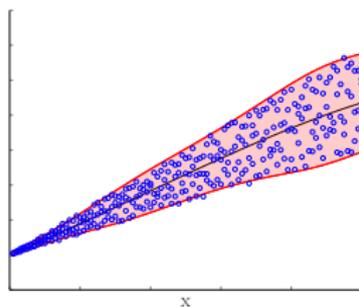
(c) Heteroscedastic and dependent

- a) Standard GPR, **GPpy/GPML**
- b) Tolvanen et al. (2014). *Expectation propagation for non-stationary heteroscedastic gaussian process regression*. IEEE-MLSP, **GPpy/GPStuff**
- c) Tabor et al. (2017). *Probabilistic modeling of heteroscedastic laboratory experiments using gaussian process regression*. Journal of Engineering Mechanics

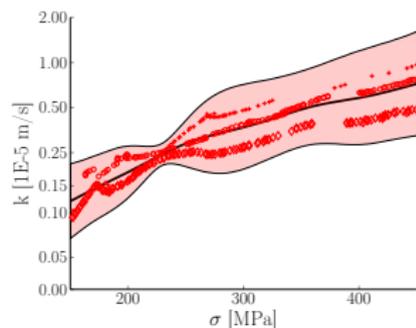
Homo/Heteroscedasticity



(a) Homoscedastic and independent



(b) Heteroscedastic and independent



(c) Heteroscedastic and dependent

- a) Standard GPR, **GPY/GPML**
- b) Tolvanen et al. (2014). *Expectation propagation for non-stationary heteroscedastic gaussian process regression*. IEEE-MLSP, **GPY/GPStuff**
- c) Tabor et al. (2017). *Probabilistic modeling of heteroscedastic laboratory experiments using gaussian process regression*. Journal of Engineering Mechanics

Regression methods

Linear regression and Gaussian Process regression are not the only approaches available

- ▶ Neural-networks (Deep-Learning)
- ▶ Support vector machine (SVM)
- ▶ Decision tree (Random-forest)
- ▶ ...++++

Most of current research carries on Neural Networks



Section Outline

Neural Networks

- 6.1 Introduction
 - 6.2 Linear Regression
 - 6.3 Activation Functions
 - 6.4 Multi-Layers Perceptron
 - 6.5 Parameter Estimation
 - 6.6 Examples
 - 6.7 Code
-

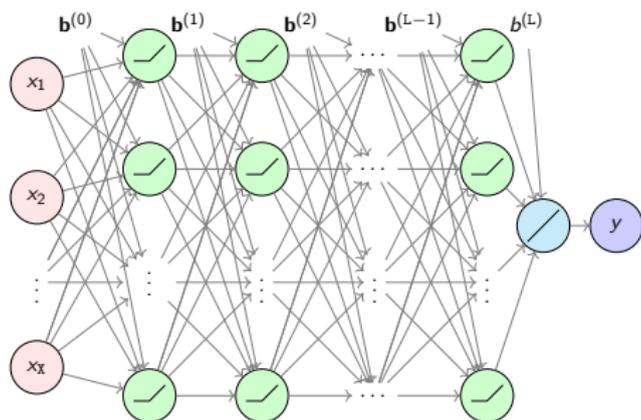
Introduction to the Neural Networks

Data

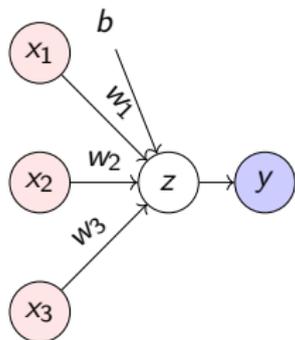
$$\mathcal{D} = \{(\mathbf{x}_i, y_i), \forall i = 1 : D\}$$

$$\mathbf{x}_i \in \mathbb{R}^X : \begin{cases} \text{Covariate} \\ \text{attribute} \\ \text{regressor} \end{cases}$$

$$y_i \in \mathbb{R} : \text{Observation}$$



Linear Regression → Neural Network



Covariates

$$\mathbf{x} = [x_1, x_2, x_3]^T$$

Hidden state variable

$$z = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

Observation equation

$$y = z + v, \quad v : V \sim \mathcal{N}(v; 0, \sigma_V^2)$$

Model parameters

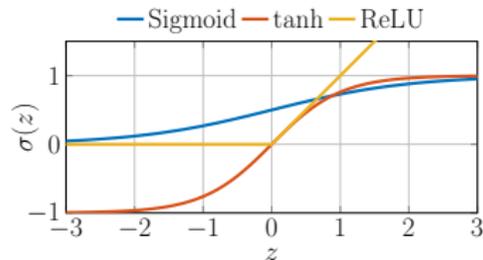
$$\theta = \{w_1, w_2, w_3, b\}$$

Linear Regression → Limitation



No matter how many layers you include → linear relationship between x and y

Activation Functions



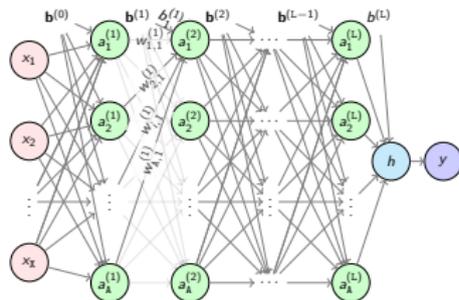
$$\sigma^s(z) = \frac{1}{1 + \exp(-z)} \quad (\text{Sigmoid})$$

$$\sigma^t(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad (\text{tanh})$$

$$\sigma^r(z) = \max(0, z) \quad (\text{ReLU})$$

Activation unit: $a_i^{(j)} = \sigma(z_i^{(j)})$

Multi-layer perceptron (MLP)



Activation unit in layer $j + 1$

Hidden variable in layer $j + 1$

Weights

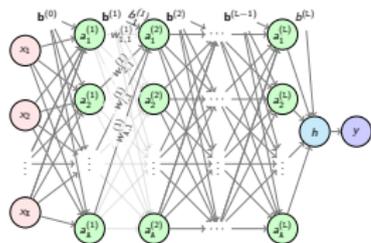
Bias

$$a_i^{(j+1)} = \sigma \left(z_i^{(j+1)} \right) = \sigma \left(w_{1,i}^{(j)} a_1^{(j)} + w_{2,i}^{(j)} a_2^{(j)} + \dots + w_{A,i}^{(j)} a_A^{(j)} + b_i^{(j)} \right)$$

Activation function

Activation units from layer j

Multi-layer perceptron (MLP) – Matrix Notation



$$\begin{aligned} \mathbf{z}^{(j+1)} &= \mathbf{W}^{(j)} \mathbf{a}^{(j)} + \mathbf{b}^{(j)} \\ &= \mathbf{W}^{(j)} \sigma(\mathbf{z}^{(j)}) + \mathbf{b}^{(j)} \end{aligned}$$

$$h = z^{(L)} = w_1^{(L)} a_1^{(L)} + w_2^{(L)} a_2^{(L)} + \dots + w_A^{(L)} a_A^{(L)} + b^{(L)}$$

$$\underbrace{\begin{bmatrix} z_1^{(j+1)} \\ z_2^{(j+1)} \\ \vdots \\ z_A^{(j+1)} \end{bmatrix}}_{\mathbf{z}^{(j+1)}} \quad A \times 1 = \underbrace{\begin{bmatrix} w_{1,1}^{(j)} & w_{1,2}^{(j)} & \dots & w_{1,A}^{(j)} \\ w_{2,1}^{(j)} & w_{2,2}^{(j)} & \dots & w_{2,A}^{(j)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{A,1}^{(j)} & w_{A,2}^{(j)} & \dots & w_{A,A}^{(j)} \end{bmatrix}}_{\mathbf{W}^{(j)}} \quad A \times A \times \underbrace{\sigma \left(\underbrace{\begin{bmatrix} z_1^{(j)} \\ z_2^{(j)} \\ \vdots \\ z_A^{(j)} \end{bmatrix}}_{\mathbf{a}^{(j)}} \right)}_{A \times 1} + \underbrace{\begin{bmatrix} b_1^{(j)} \\ b_2^{(j)} \\ \vdots \\ b_A^{(j)} \end{bmatrix}}_{\mathbf{b}^{(j)}} \quad A \times 1$$

Multi-layer perceptron (MLP) – Parameters

For fully connected network with L layers, each having A activation units,

$$\begin{aligned} \boldsymbol{\theta}_w &\in \mathbb{R}^{A \times A \times (L-1)} && \text{weight parameters, } w_{i,j}^l \\ \boldsymbol{\theta}_b &\in \mathbb{R}^{A \times L} && \text{bias parameters, } b_j^l \end{aligned}$$

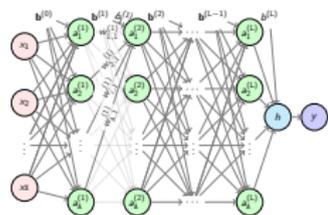
For $L = 2$, $A = 1000$

$$\boldsymbol{\theta} = \{\boldsymbol{\theta}_w, \boldsymbol{\theta}_b\} \in \mathbb{R}^{\approx 10^6}$$

Joint Log-likelihood – $\theta = \{\theta_w, \theta_b\}$

$$\mathcal{D} = \{\mathcal{D}_x, \mathcal{D}_y\} = \{(\mathbf{x}_i, y_i), \forall i = 1 : D\}$$

$$f(\mathcal{D}_y | \mathcal{D}_x, \theta) = \prod_{i=1}^D \mathcal{N}(y_i; h(\mathbf{x}_i; \theta), \sigma_V^2)$$



$$\ln f(\mathcal{D}_y | \mathcal{D}_x, \theta) = \sum_{i=1}^D \ln(\mathcal{N}(y_i; h(\mathbf{x}_i; \theta), \sigma_V^2))$$

$$\propto -\frac{1}{D} \sum_{i=1}^D \frac{1}{2} \|h(\mathbf{x}_i; \theta) - y_i\|^2$$

$$= -\frac{1}{D} \sum_{i=1}^D J(\mathbf{x}_i, y_i; \theta)$$

$$= -J(\mathcal{D}; \theta)$$

$$\theta^* = \arg \min_{\theta} J(\mathcal{D}; \theta)$$

Parameter estimation → Gradient descent

 λ : **Learning rate**

$$w_{i,j}^{(l)} \leftarrow w_{i,j}^{(l)} - \lambda \nabla_{w_{i,j}^{(l)}} J(\mathcal{D}; \theta)$$

 $\nabla_{w_{i,j}^{(l)}} J(\mathcal{D}; \theta)$: gradient w.r.t $w_{i,j}^{(l)}$

$$b_j^{(l)} \leftarrow b_j^{(l)} - \lambda \nabla_{b_j^{(l)}} J(\mathcal{D}; \theta)$$

 $\nabla_{b_j^{(l)}} J(\mathcal{D}; \theta)$: gradient w.r.t $b_j^{(l)}$ Gradients are estimated using **backpropagation** (Chain rule)

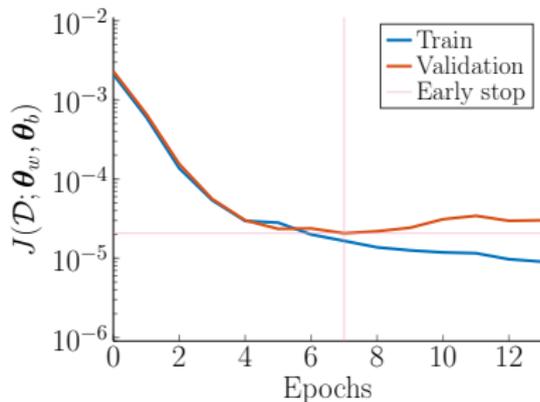
Estimating θ^* is done using **(stochastic) gradient descent** (ADAM, ADAGRAD, etc...) over multiple **epochs**

- ▶ **SGD**: Compute ∇ using small batches of data
- ▶ **Epoch**: Seeing the entire dataset once
- ▶ **Early stop**: Stop when the loss on the validation set increases

Parameter estimation \rightarrow Gradient descent

$$w_{i,j}^{(l)} \leftarrow w_{i,j}^{(l)} - \lambda \nabla_{w_{i,j}^{(l)}} J(\mathcal{D}; \theta)$$

$$b_j^{(l)} \leftarrow b_j^{(l)} - \lambda \nabla_{b_j^{(l)}} J(\mathcal{D}; \theta)$$



Estimating θ^* is done using **(stochastic) gradient descent** (ADAM, ADAGRAD, etc...) over multiple **epochs**

- ▶ **SGD**: Compute ∇ using small batches of data
- ▶ **Epoch**: Seeing the entire dataset once
- ▶ **Early stop**: Stop when the loss on the validation set increases

Parameter estimation & data standardization

In order to perform well, NN need to work on **standardized data**

$$\begin{aligned}\tilde{\mathbf{x}}_i &= \frac{\mathbf{x}_i - \mu_X}{\sigma_X} \\ \tilde{\mathbf{y}}_i &= \frac{\mathbf{y}_i - \mu_Y}{\sigma_Y}\end{aligned}$$

Then you need to destandardize the NN output $g(\tilde{\mathbf{x}}; \theta) \odot \sigma_Y + \mu_Y$

Why? Because otherwise, you would need to initialize

$$\theta = \underbrace{\text{fct}(\mu_X, \sigma_X, \mu_Y, \sigma_Y)}_{\Delta_X}$$

UCI - Concrete Slump dataset

Concrete Slump Test Data Set

Download: [Data Folder](#) [Data Set Description](#)

Abstract: Concrete is a highly complex material. The slump flow of concrete is not only determined by the water content, but that is also influenced by other concrete ingredients.



Data Set Characteristics:	Multivariate	Number of Instances:	103	Area:	Computer
Attribute Characteristics:	Real	Number of Attributes:	10	Date Donated	2009-04-30
Associated Tasks:	Regression	Missing Values?	N/A	Number of Web Hits:	132011

Source:

Donor: I-Cheng Yeh
 Email: icyeh@chu.edu.tw
 Institution: Department of Information Management, Chung-Hua University (Republic of China)
 Other contact information: Department of Information Management, Chung-Hua University, Hsin Chu, Taiwan 30067, R.O.C.

<https://archive.ics.uci.edu/ml/datasets/Concrete+Slump+Test>

1 output [MPa]

- ▶ 28-day compressive strength

7 covariates [kg/m³]

- ▶ Cement
- ▶ Slag
- ▶ Fly ash
- ▶ Water
- ▶ Superplasticizer
- ▶ Coarse aggregate
- ▶ Fine aggregate

Problem set-up

Covariates and observations:

$$\mathbf{x}_i = [x_1, x_2, \dots, x_7]^T$$

$$\tilde{\mathcal{D}} = \{\tilde{\mathcal{D}}_x, \tilde{\mathcal{D}}_y\} = \left\{ \left(\frac{x_i - \mu_X}{\sigma_X}, \frac{y_i - \mu_Y}{\sigma_Y} \right), i = 1 : 103 \right\}$$

Observation model:

$$\underbrace{\tilde{y}_i}_{\text{observation}} = g(\underbrace{\tilde{\mathbf{x}}_i}_{\text{true strength}})$$

Neural network setup

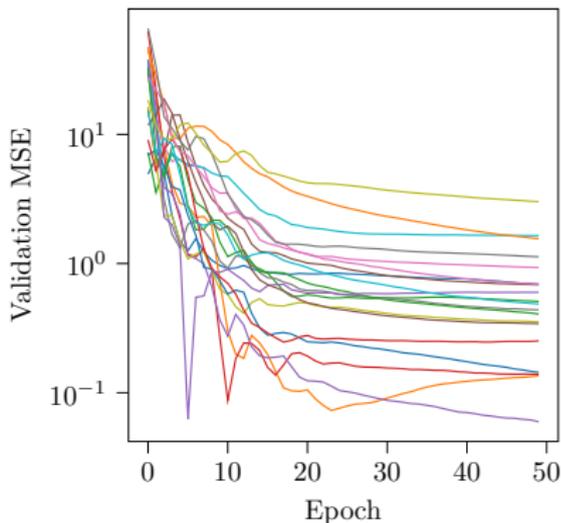
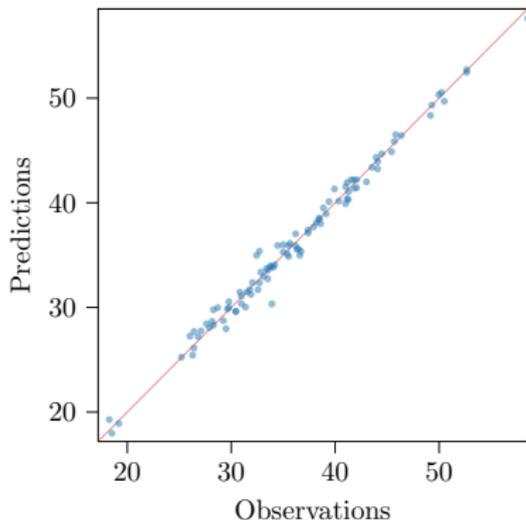
Loss function

$$\theta^* = \arg \min_{\theta} \overbrace{\frac{1}{D} \sum_{i=1}^D (\tilde{y}_i - g(\tilde{\mathbf{x}}_i; \theta))^2}^{\text{Mean square error}}$$

Network configuration

- ▶ $L = 1$ (nb. hidden layer)
- ▶ $A = 200$ (nb. hidden units/layer)
- ▶ $E = 50$ (nb. epochs)
- ▶ $B = 20$ (Batch size)
- ▶ $R = 5$ (nb. of repetitions)
- ▶ learning rate = 0.005

Results 20-folds CV \times 5 repetitions



Mean square error (MSE)	$0.696 \pm 0.14 \gg 0.118$ (GPR) 
Mean absolute error (MAE)	0.572 ± 0.03
Mean error (ME)	-0.003 ± 0.04
Std error (RMSE)	0.830 ± 0.08

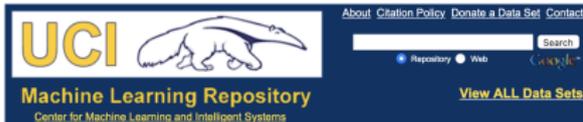
UCI - Concrete Strength dataset: **GPR v.s. NN**

1 output [MPa]

- ▶ 28-day concrete strength

8 covariates

- ▶ Cement [kg/m^3]
- ▶ Blast Furnace Slag [kg/m^3]
- ▶ Fly ash [kg/m^3]
- ▶ Water [kg/m^3]
- ▶ Superplasticizer [kg/m^3]
- ▶ Coarse aggregate [kg/m^3]
- ▶ Fine aggregate [kg/m^3]
- ▶ Age [days]



Concrete Compressive Strength Data Set

Download: [Data Folder](#), [Data Set Description](#)

Abstract: Concrete is the most important material in civil engineering. The concrete compressive strength is a highly nonlinear function of age and ingredients.



Data Set Characteristics:	Multivariate	Number of Instances:	1030	Area:	Physical
Attribute Characteristics:	Real	Number of Attributes:	9	Date Donated	2007-06-03
Associated Tasks:	Regression	Missing Values?	N/A	Number of Web Hits:	316898

Source:

Original Owner and Donor
 Prof. J-Cheng Yeh
 Department of Information Management
 Chung-Hua University,
 Heih Chu, Taiwan 33067, R.O.C.
 e-mail:cyeh@im.chu.edu.tw
 TEL:886-3-5186511

Date Donated: August 3, 2007

<https://archive.ics.uci.edu/ml/datasets/concrete+compressive+strength>

Problem set-up

Covariates and observations:

$$\mathbf{x}_i = [x_1, x_2, \dots, x_8]^T$$

$$\tilde{\mathcal{D}} = \{\tilde{\mathcal{D}}_x, \tilde{\mathcal{D}}_y\} = \left\{ \left(\frac{\mathbf{x}_i - \boldsymbol{\mu}_X}{\boldsymbol{\sigma}_X}, \frac{y_i - \mu_Y}{\sigma_Y} \right), i = 1 : 1030 \right\}$$

Observation models:

GPR

$$\underbrace{\tilde{y}_i}_{\text{observation}} = \underbrace{g(\tilde{\mathbf{x}}_i)}_{\text{true strength}} + \underbrace{v_i}_{\text{meas. error}}$$

NN

$$\underbrace{\tilde{y}_i}_{\text{observation}} = \underbrace{g(\tilde{\mathbf{x}}_i)}_{\text{true strength}}$$

Model setup

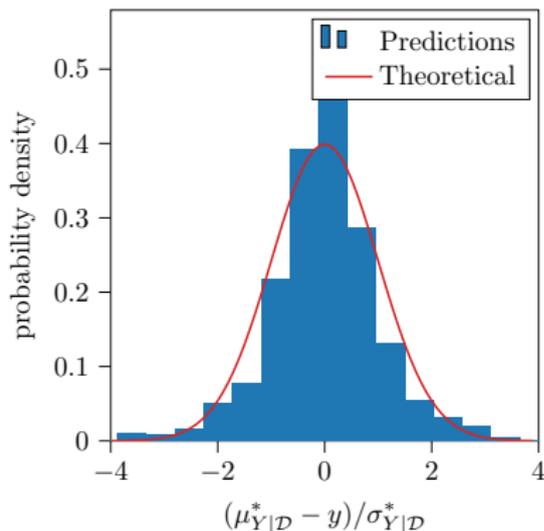
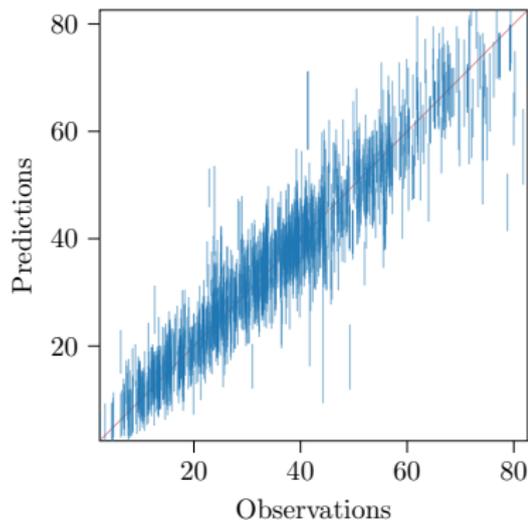
GPR

- ▶ Square exponential correlation fct. (RBF)
- ▶ $\emptyset \rightarrow$ default parameters

NN

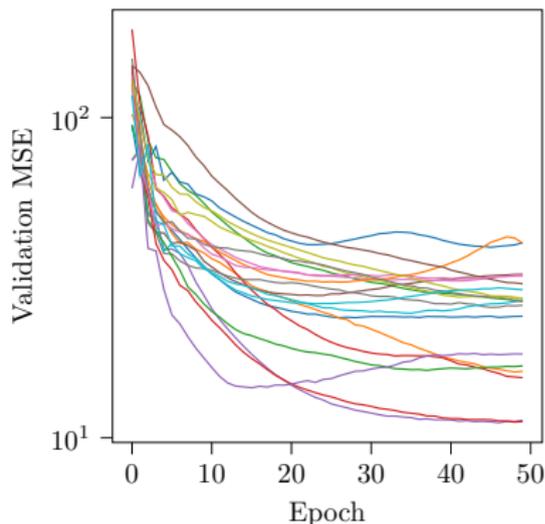
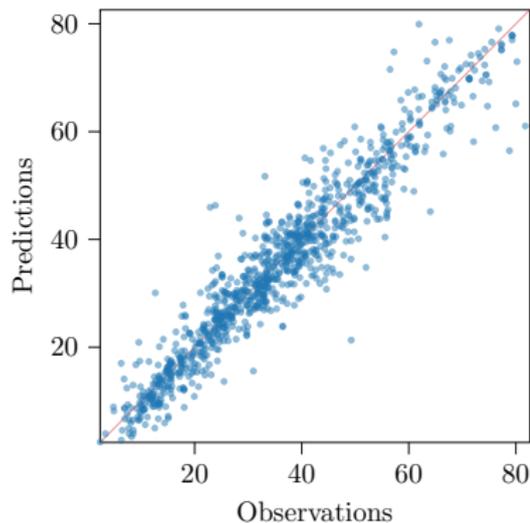
- ▶ Loss function: MSE
- ▶ $L = 1$ (nb. hidden layer)
- ▶ $A = 500$ (nb. hidden units)
- ▶ $E = 50$ (nb. epochs)
- ▶ $B = 100$ (Batch size)
- ▶ $R = 5$ (nb. of repetitions)
- ▶ learning rate = 0.005

GPR 20-folds CV [Python]



Mean square error (MSE)	23.9
Mean absolute error (MAE)	3.4
Mean error (ME)	-0.006
Std error (RMSE)	4.9 ± 0.08

NN 20-folds CV \times 5 repetitions



Mean square error (MSE)	$23.1 \pm 0.1 \approx 23.9$ (GPR) 
Mean absolute error (MAE)	3.5 ± 0.02
Mean error (ME)	-0.2 ± 0.03
Std error (RMSE)	4.9 ± 0.01

Defining a FNN with PyTorch

```
1 import torch
2 import numpy as np
3 import pandas as pd
4
5 #Load the data
6 data = pd.read_excel('Data.xls', 'Sheet1', header=None, skiprows=1)
7 data = pd.read_csv('Data.csv', header=None, skiprows=1)
8
9 data_array = np.array(data)
10 data_array = np.random.permutation(data_array) #Randomly shuffle the rows
11
12 input_size, output_size, nb_units = 10, 1, 25
13 x = data_array[:,0:input_size]
14 y = data_array[:,input_size:input_size+output_size]
15
16 #Define the network
17 model = torch.nn.Sequential(
18     torch.nn.Linear(input_size, nb_units),
19     torch.nn.ReLU(),
20     torch.nn.Linear(nb_units, output_size))
21 criterion = torch.nn.MSELoss() #Define loss function
22 optimizer = torch.optim.Adam(model.parameters(), lr=0.005) #lr: learning rate
```

Note: you need to reinitialize the model for every CV fold

Training a FNN with PyTorch

```

1  #CV index selection
2  idx = ... #Define your index for the train/validation CV fold
3  x_train = np.delete(x, idx, 0)
4  y_train = np.delete(y, idx, 0)
5  x_val = x[idx]
6  #Standardization
7  mx = np.mean([x_train],axis=1)
8  sx = np.std([x_train],axis=1)
9  my = np.mean([y_train])
10 sy = np.std([y_train])
11 x_train = (x_train-mx)/sx
12 y_train = (y_train-my)/sy
13 x_val = (x_val-mx)/sx
14 x_train = torch.tensor(x_train,dtype=torch.float32) #Transform data to tensor
15 y_train = torch.tensor(y_train,dtype=torch.float32) #Transform data to tensor
16 #Training
17 for e in range(nb_epochs): #Loop over epochs
18     nb_batch = ...
19     for b in range(nb_batch): #Loop over batches
20         idx_b = ... #Define index indicating which observation to use
21         y_pred = model(x_train[idx_b]) #Predict
22         loss = criterion(y_pred,y_train[idx_b].reshape(-1,1)) #Calculate loss
23         optimizer.zero_grad() #Backprop
24         loss.backward()
25         optimizer.step()

```

Testing a FNN with PyTorch

```
1 #Prediction
2 x_val = torch.tensor(x_val, dtype=torch.float32) #Transform data to tensor
3 y_pred_val = model(x_val)*sy+my #Predict & unnormalize for the validation set
4 y_pred_val = np.concatenate(y_pred_val.detach().numpy())#Tensor to array
```

Summary

Linear Regression: outperformed by modern methods

Gaussian Process Regression:

- ▶ State-of-the-art for small datasets ($\approx D \leq 10^3 - 10^4$)
- ▶ Allows interpolating and extrapolating (uncertainty estimates)
- ▶ Few hyperparameters

Neural Networks:

- ▶ State-of-the-art for large datasets
- ▶ Requires large datasets to outperform other methods
- ▶ Several hyperparameters

Cross-Validation:

- ▶ Prevents **over-fitting**
- ▶ The standard method for quantifying accuracy & comparing models