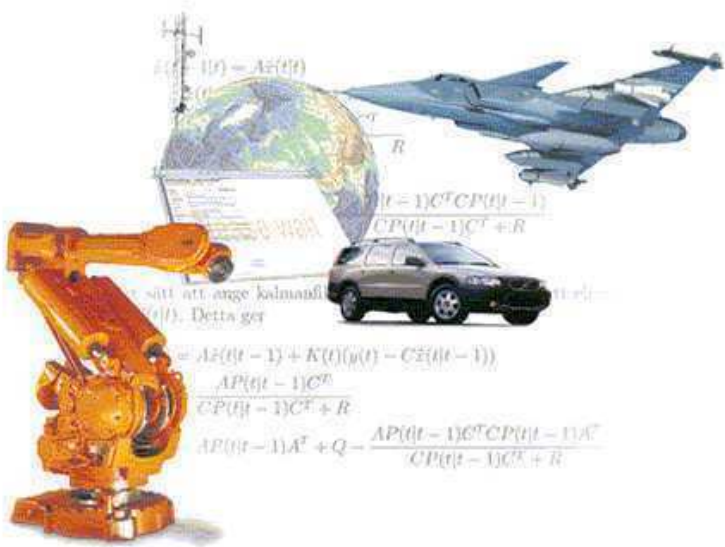


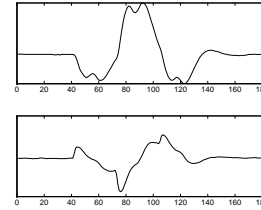
System Identification: *The Path from Data to Model*

Lennart Ljung

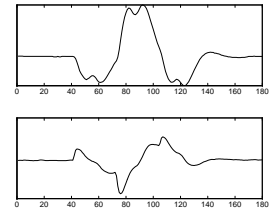
Division of Automatic Control
Linköping University
Sweden



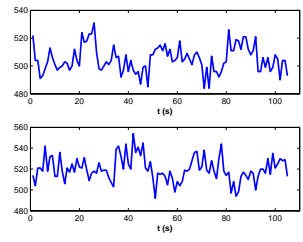
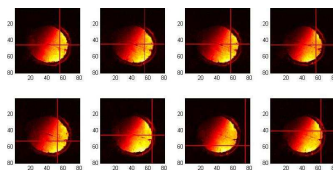
Aircraft Dynamics:



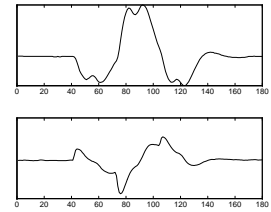
Aircraft Dynamics:



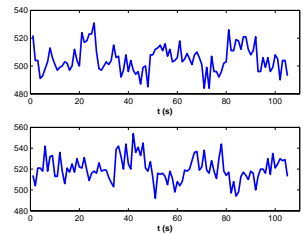
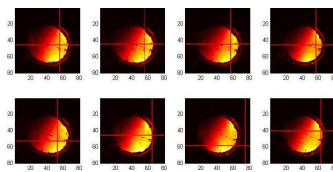
Brain Activity (fMRI):



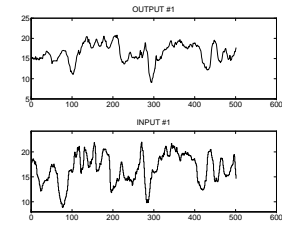
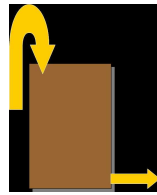
Aircraft Dynamics:



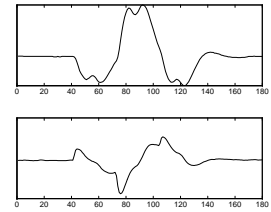
Brain Activity (fMRI):



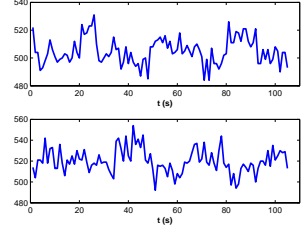
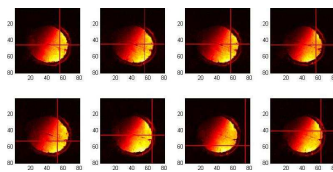
Pulp Buffer Vessel:



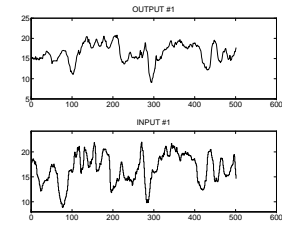
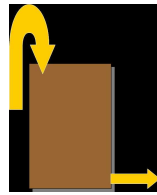
Aircraft Dynamics:



Brain Activity (fMRI):



Pulp Buffer Vessel:

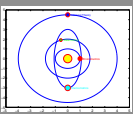


Industrial Engineering:



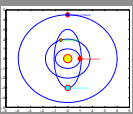
90 papers





Constructing (mathematical) models from data is a prime problem in many scientific fields and many application areas.

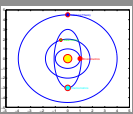




Constructing (mathematical) models from data is a prime problem in many scientific fields and many application areas.

Many communities and cultures around the area have grown, with their own nomenclatures and their own “social lives”.



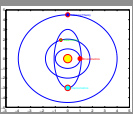


Constructing (mathematical) models from data is a prime problem in many scientific fields and many application areas.

Many communities and cultures around the area have grown, with their own nomenclatures and their own “social lives”.

This has created a very rich, and somewhat confusing, plethora of methods and approaches for the problem.





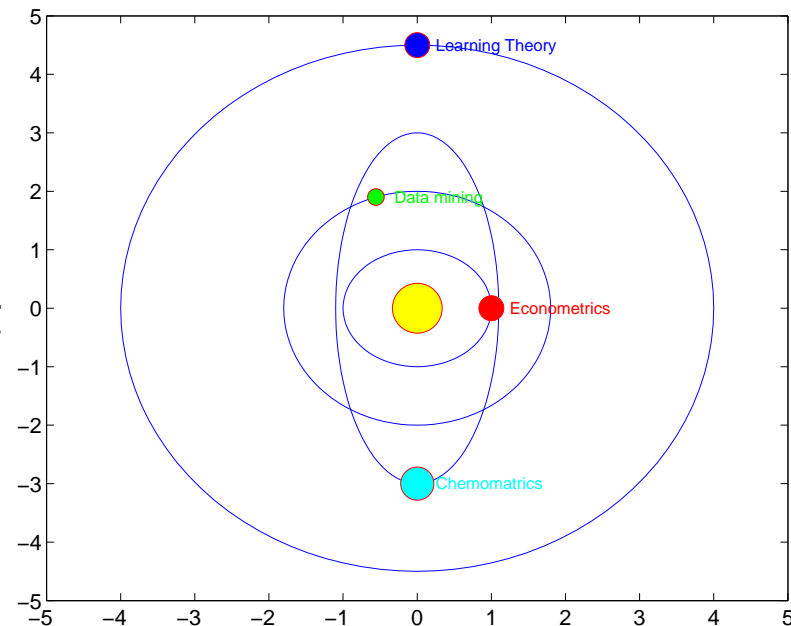
The Communities

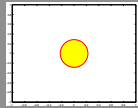
Constructing (mathematical) models from data is a prime problem in many scientific fields and many application areas.

Many communities and cultures around the area have grown, with their own nomenclatures and their own “social lives”.

This has created a very rich, and somewhat confusing, plethora of methods and approaches for the problem.

A picture: There is a core of central material, encircled by the different communities.

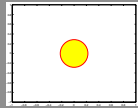




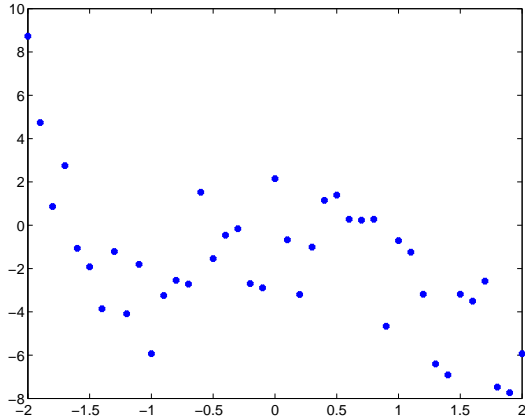
Central terms

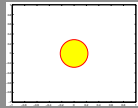
- Model m – Model Class \mathcal{M} – Complexity (Flexibility) \mathcal{C}
- Information \mathcal{I} – Data Z
- Estimation – Validation (Learning – Generalization)
- Model fit $\mathcal{F}(m, Z)$



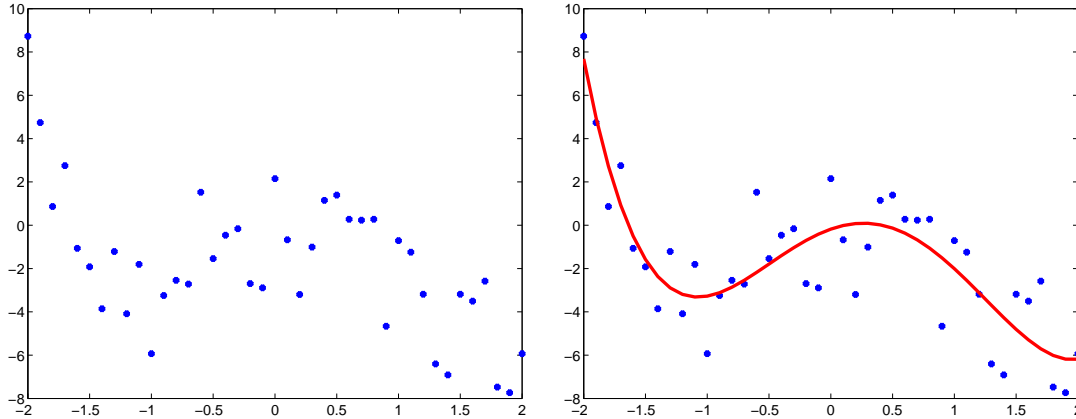


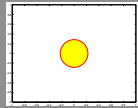
information in data



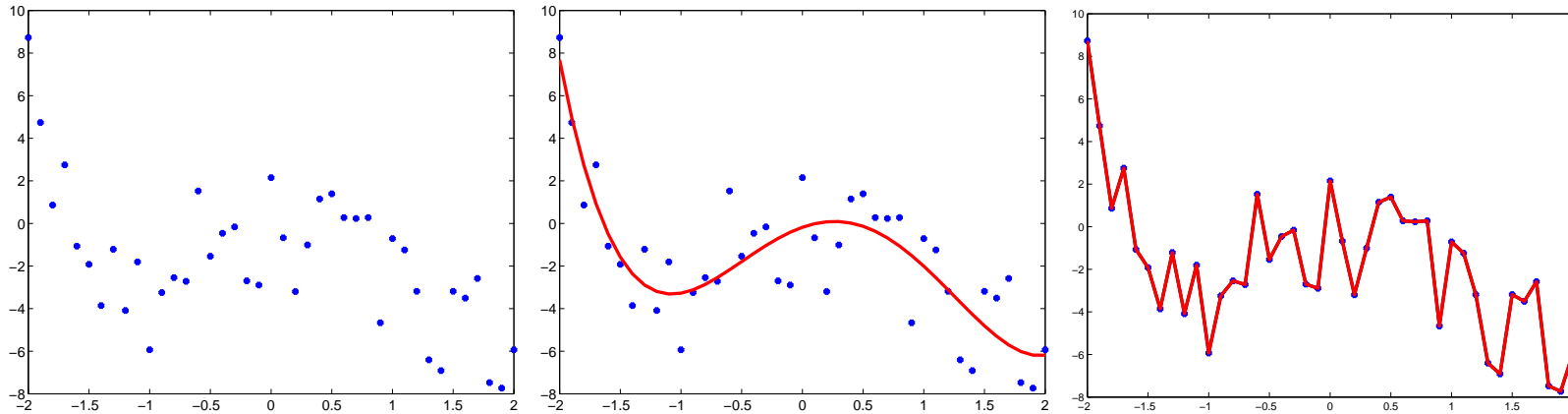


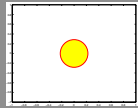
Squeeze out the relevant information in data



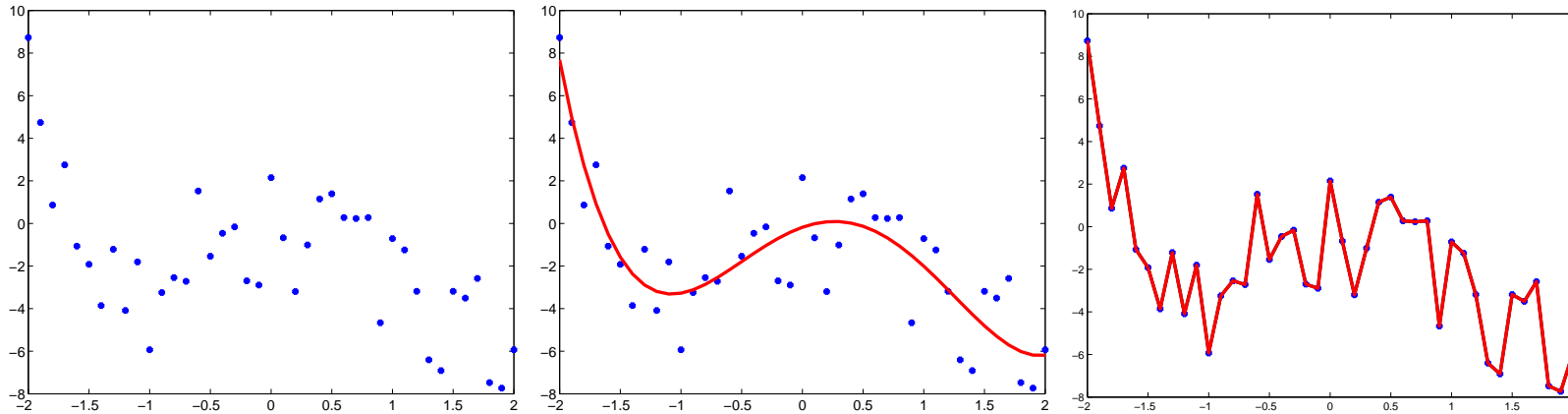


Squeeze out the relevant information in data. (BUT NOT MORE!)



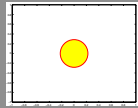


Squeeze out the relevant information in data. (BUT NOT MORE!)

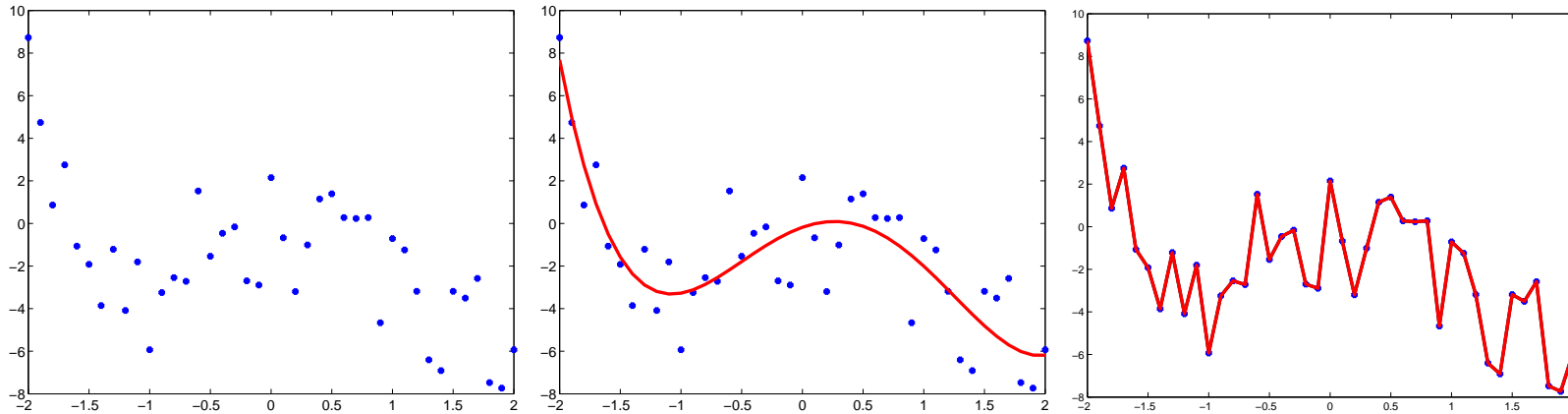


All data contain Information and Misinformation (“Signal and noise”).





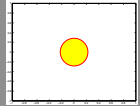
Squeeze out the relevant information in data. (BUT NOT MORE!)



All data contain Information and Misinformation (“Signal and noise”).

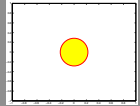
So need to meet the data with a prejudice!





Estimation Prejudices





Nature is Simple!

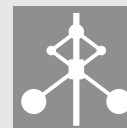


Nature is Simple! (Occam's razor, Lex Parsimoniae...)



Nature is Simple! (Occam's razor, Lex Parsimoniae...)

God is subtle, but He is not malicious (Einstein)



Nature is Simple! (Occam's razor, Lex Parsimoniae...)

God is subtle, but He is not malicious (Einstein)

So, conceptually:

$$\hat{m} = \arg \min_{m \in \mathcal{M}} (\text{Fit} + \text{Complexity Penalty})$$



Nature is Simple! (Occam's razor, Lex Parsimoniae...)

God is subtle, but He is not malicious (Einstein)

So, conceptually:

$$\hat{m} = \arg \min_{m \in \mathcal{M}} (\text{Fit} + \text{Complexity Penalty})$$

Examples:

- Search for a model in sets with a maximal Complexity
- The Akaike criterion
- Regularization



Fit to estimation data Z_e^N (N : Number of data points)

$$F(\hat{m}, Z_e^N) \quad (\text{"The empirical risk"})$$



Fit to estimation data Z_e^N (N : Number of data points)

$$F(\hat{m}, Z_e^N) \quad (\text{"The empirical risk"})$$

Now try your model on a fresh data set (Validation data Z_v):

$$EF(\hat{m}, Z_v) \approx F(\hat{m}, Z_e^N) + f(C(\mathcal{M}), N)$$



Fit to estimation data Z_e^N (N : Number of data points)

$$F(\hat{m}, Z_e^N) \quad (\text{"The empirical risk"})$$

Now try your model on a fresh data set (Validation data Z_v):

$$EF(\hat{m}, Z_v) \approx F(\hat{m}, Z_e^N) + f(\mathcal{C}(\mathcal{M}), N)$$

f increases with the complexity \mathcal{C} and decreases with N , so the more flexible the model set the worse fit to validation data.



Fit to estimation data Z_e^N (N : Number of data points)

$$F(\hat{m}, Z_e^N) \quad (\text{"The empirical risk"})$$

Now try your model on a fresh data set (Validation data Z_v):

$$EF(\hat{m}, Z_v) \approx F(\hat{m}, Z_e^N) + f(\mathcal{C}(\mathcal{M}), N)$$

f increases with the complexity \mathcal{C} and decreases with N , so the more flexible the model set the worse fit to validation data.

In words: If you have a model that describes the estimation data well, the fit will be (much) worse when you try it on validation data.



Fit to estimation data Z_e^N (N : Number of data points)

$$F(\hat{m}, Z_e^N) \quad (\text{"The empirical risk"})$$

Now try your model on a fresh data set (Validation data Z_v):

$$EF(\hat{m}, Z_v) \approx \mathcal{F}(\hat{m}, Z_e^N) + f(\mathcal{C}(\mathcal{M}), N)$$

f increases with the complexity \mathcal{C} and decreases with N , so the more flexible the model set the worse fit to validation data.

In words: If you have a model that describes the estimation data well, the fit will be (much) worse when you try it on validation data.

So don't be impressed by a good fit to estimation data in a flexible model set!



\mathcal{S} – True system \mathcal{M} – Model set \hat{m} – Estimate

m^* – Expected model $m^* = E\hat{m}$ Typically $m^* = \arg \min_{m \in \mathcal{M}} \|\mathcal{S} - m\|^2$

Then



\mathcal{S} – True system \mathcal{M} – Model set \hat{m} – Estimate

m^* – Expected model $m^* = E\hat{m}$ Typically $m^* = \arg \min_{m \in \mathcal{M}} \|\mathcal{S} - m\|^2$

Then

$$E\|\mathcal{S} - \hat{m}\|^2 = \|\mathcal{S} - m^*\|^2 + E\|\hat{m} - m^*\|^2$$

$$\text{MSE} = \text{B: BIAS} + \text{V: Variance}$$



\mathcal{S} – True system \mathcal{M} – Model set \hat{m} – Estimate

m^* – Expected model $m^* = E\hat{m}$ Typically $m^* = \arg \min_{m \in \mathcal{M}} \|\mathcal{S} - m\|^2$

Then

$$E\|\mathcal{S} - \hat{m}\|^2 = \|\mathcal{S} - m^*\|^2 + E\|\hat{m} - m^*\|^2$$

MSE = B: BIAS + V: Variance

Error: = Systematic+Random



\mathcal{S} – True system \mathcal{M} – Model set \hat{m} – Estimate

m^* – Expected model $m^* = E\hat{m}$ Typically $m^* = \arg \min_{m \in \mathcal{M}} \|\mathcal{S} - m\|^2$

Then

$$E\|\mathcal{S} - \hat{m}\|^2 = \|\mathcal{S} - m^*\|^2 + E\|\hat{m} - m^*\|^2$$

MSE = B: BIAS + V: Variance

Error: = Systematic+Random

As model complexity increases, Bias decreases and Variance increases



\mathcal{S} – True system \mathcal{M} – Model set \hat{m} – Estimate

m^* – Expected model $m^* = E\hat{m}$ Typically $m^* = \arg \min_{m \in \mathcal{M}} \|\mathcal{S} - m\|^2$

Then

$$E\|\mathcal{S} - \hat{m}\|^2 = \|\mathcal{S} - m^*\|^2 + E\|\hat{m} - m^*\|^2$$

MSE = B: BIAS + V: Variance

Error: = Systematic+Random

As model complexity increases, Bias decreases and Variance increases

This bias/variance trade-off is at the heart of estimation.

Note that the model complexity that minimizes the MSE typically has a non-zero systematic error.



Bottom line: There is a theoretic best accuracy that can be achieved in estimation, independent of methods and computational effort. This bound depends on prior knowledge and data quality.



Bottom line: There is a theoretic best accuracy that can be achieved in estimation, independent of methods and computational effort. This bound depends on prior knowledge and data quality.

Formalization: Observe Y . Let its probability density function (pdf) be $f_Y(x, \theta)$

The (Fisher) Information Matrix is

$$\mathcal{I} = E \ell'_Y (\ell'_Y)^T, \quad \ell'_Y = \frac{\partial}{\partial \theta} \log f_Y(x, \theta)$$



Bottom line: There is a theoretic best accuracy that can be achieved in estimation, independent of methods and computational effort. This bound depends on prior knowledge and data quality.

Formalization: Observe Y . Let its probability density function (pdf) be $f_Y(x, \theta)$

The **(Fisher) Information Matrix** is

$$\mathcal{I} = E \ell'_Y (\ell'_Y)^T, \quad \ell'_Y = \frac{\partial}{\partial \theta} \log f_Y(x, \theta)$$

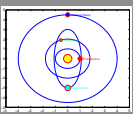
The **Cramér-Rao inequality** tells us that

$$\text{cov} \hat{\theta} \geq \mathcal{I}^{-1}$$

for any (unbiased) estimator $\hat{\theta}$ of the parameter.

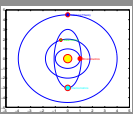
\mathcal{I} is thus a prime quantity for Experiment Design.





- **Statistics, The Mother Area**
 - Recent activities...
 - Bootstrap
 - Regularization to control complexity (LASSO, LARS,...)





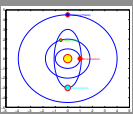
■ Statistics, The Mother Area

- Recent activities...
- Bootstrap
- Regularization to control complexity (LASSO, LARS,...)

■ Econometrics

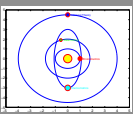
- Volatility Clustering (varying variance), GARCH
- Common roots for variations (co-integration)





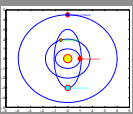
- **Statistics, The Mother Area**
 - Recent activities...
 - Bootstrap
 - Regularization to control complexity (LASSO, LARS,...)
- **Econometrics**
 - Volatility Clustering (varying variance), GARCH
 - Common roots for variations (co-integration)
- **Statistical Learning Theory**
 - Convex Formulations, SVM
 - VC-dimensions





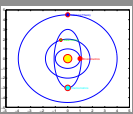
- **Statistics, The Mother Area**
 - Recent activities...
 - Bootstrap
 - Regularization to control complexity (LASSO, LARS,...)
- **Econometrics**
 - Volatility Clustering (varying variance), GARCH
 - Common roots for variations (co-integration)
- **Statistical Learning Theory**
 - Convex Formulations, SVM
 - VC-dimensions





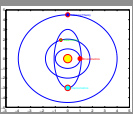
- **Chemometrics – Statistical Process Control**
 - High-dimensional Data Spaces (Many process variables)





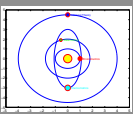
- **Chemometrics – Statistical Process Control**
 - High-dimensional Data Spaces (Many process variables)
- **Data Mining**
 - The Internet!





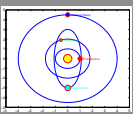
- **Chemometrics – Statistical Process Control**
 - High-dimensional Data Spaces (Many process variables)
- **Data Mining**
 - The Internet!
- **Machine Learning**





- **Chemometrics – Statistical Process Control**
 - High-dimensional Data Spaces (Many process variables)
- **Data Mining**
 - The Internet!
- **Machine Learning**
 - Grown out of artificial intelligence, more and more statistically oriented





- **Chemometrics – Statistical Process Control**
 - High-dimensional Data Spaces (Many process variables)
- **Data Mining**
 - The Internet!
- **Machine Learning**
 - Grown out of artificial intelligence, more and more statistically oriented
- **System Identification**
 - Dynamical systems

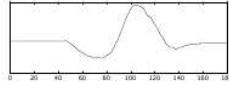
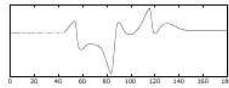


System

Input



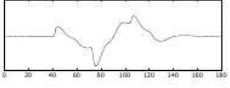
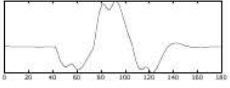
rudders
ailerons
thrust

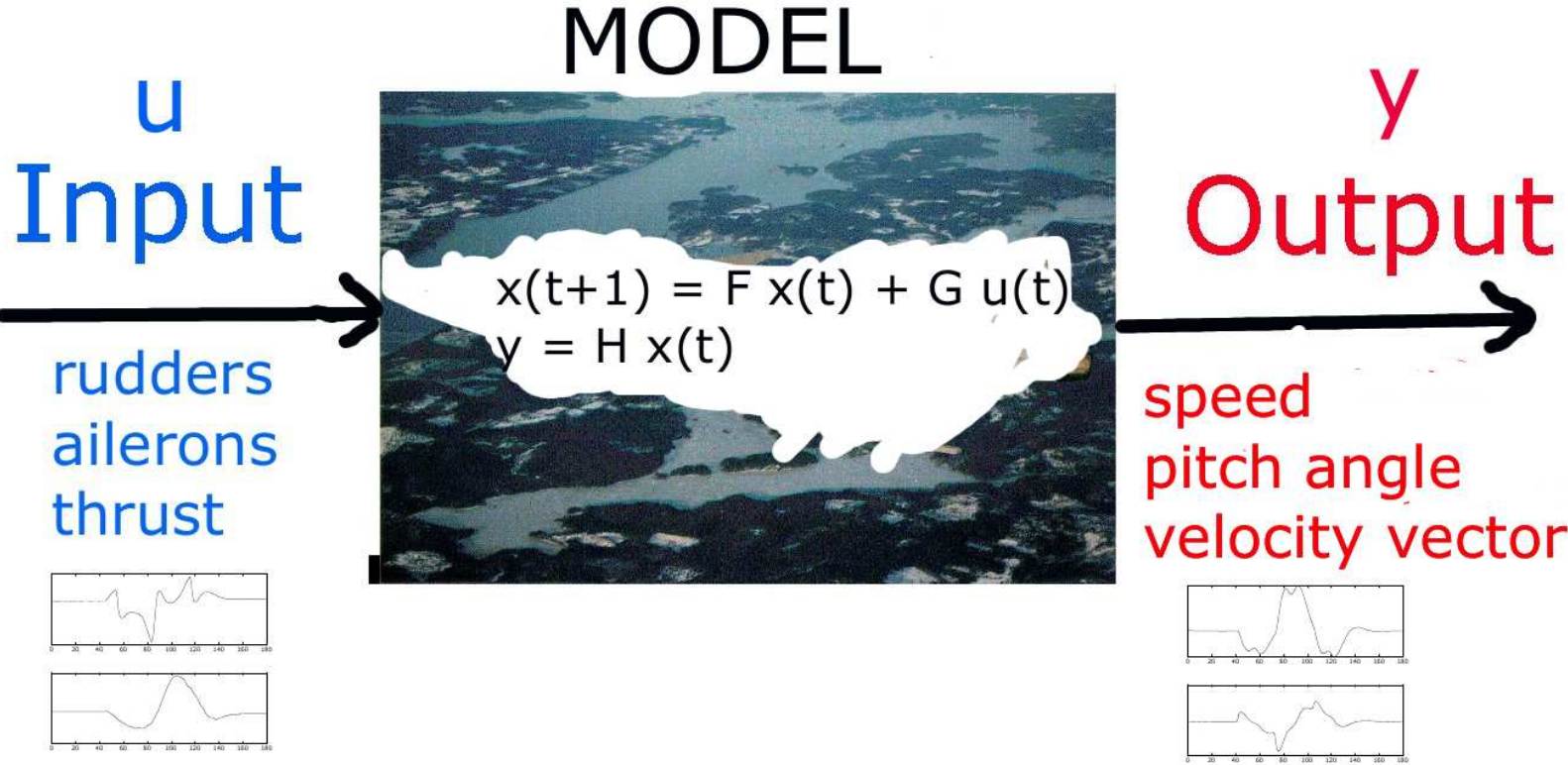


Output



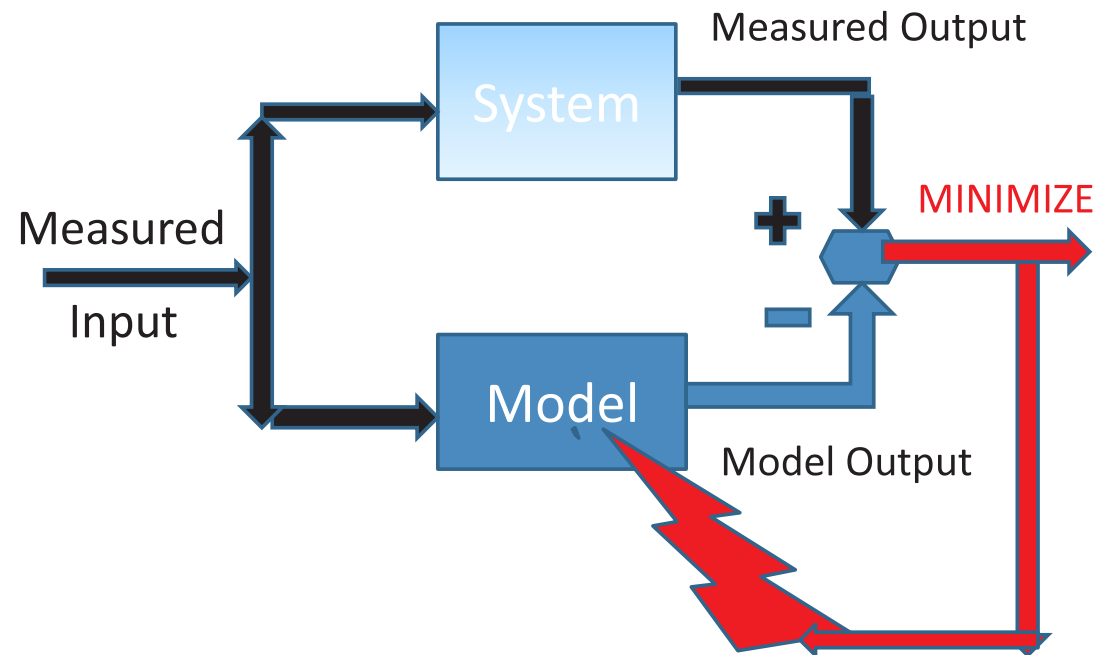
speed
pitch angle
velocity vector





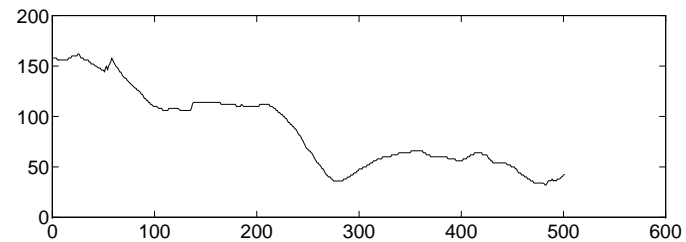
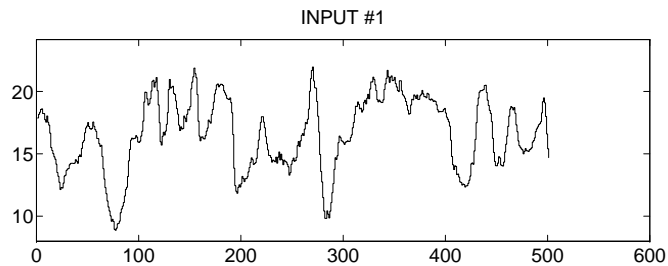
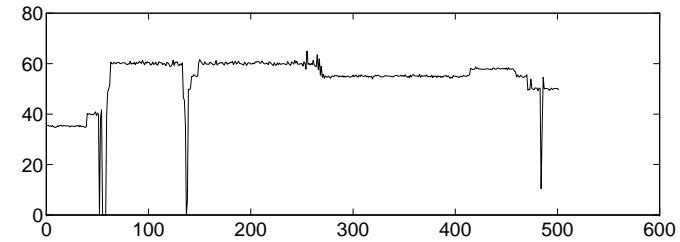
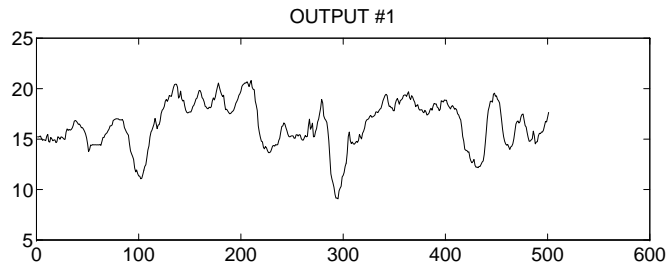
Let the model predict the next output and minimize the error in prediction:

Fitting Models to Data



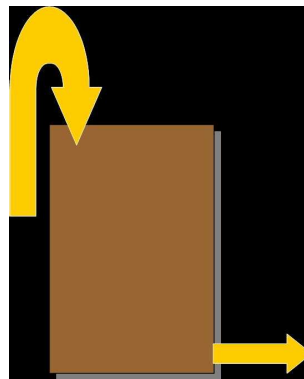
- Use physical insights/common sense together with standard flexible model classes:
- Color Coded: Black-box, Grey-box, White and Off-White Models.
- Let us look at a concrete example from Process Industry

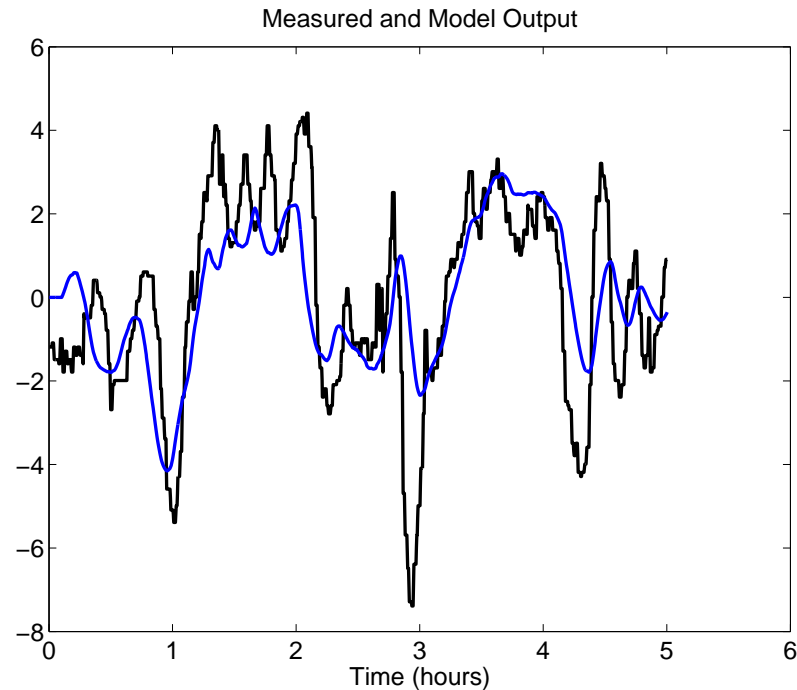
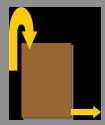




κ -number of outflow,
 κ -number of inflow,

flow
volume



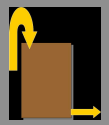


Black line: κ -number after the vessel, actual measurements.

Blue line: Simulated κ -number using the input only and a process model

estimated using the first 200 data points. $G(s) = \frac{0.818}{1+676s} e^{-480s}$

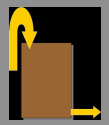




Now it's time to

Think:



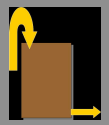


Think:

If no mixing in tank (“plug flow”) a particle that enters the top will exit T seconds later, where

$$T = \frac{\text{Tank Volume}}{\text{Flow}} : \left[\frac{m^3}{m^3/s} = s \right]$$





Think:

If no mixing in tank ("plug flow") a particle that enters the top will exit T seconds later, where

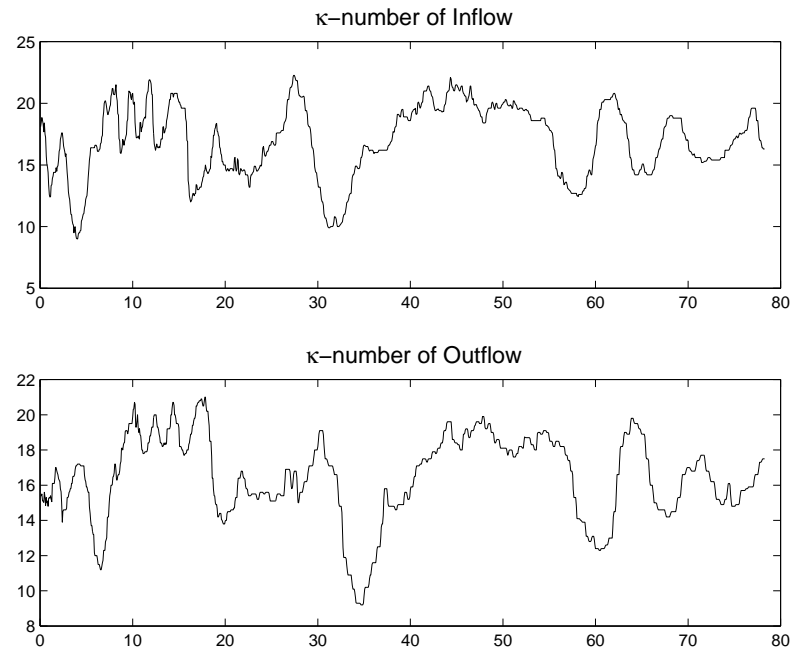
$$T = \frac{\text{Tank Volume}}{\text{Flow}} : \left[\frac{m^3}{m^3/s} = s \right]$$

But this "natural delay time" is time-varying: $T = T(t)$, since flow and volume changes

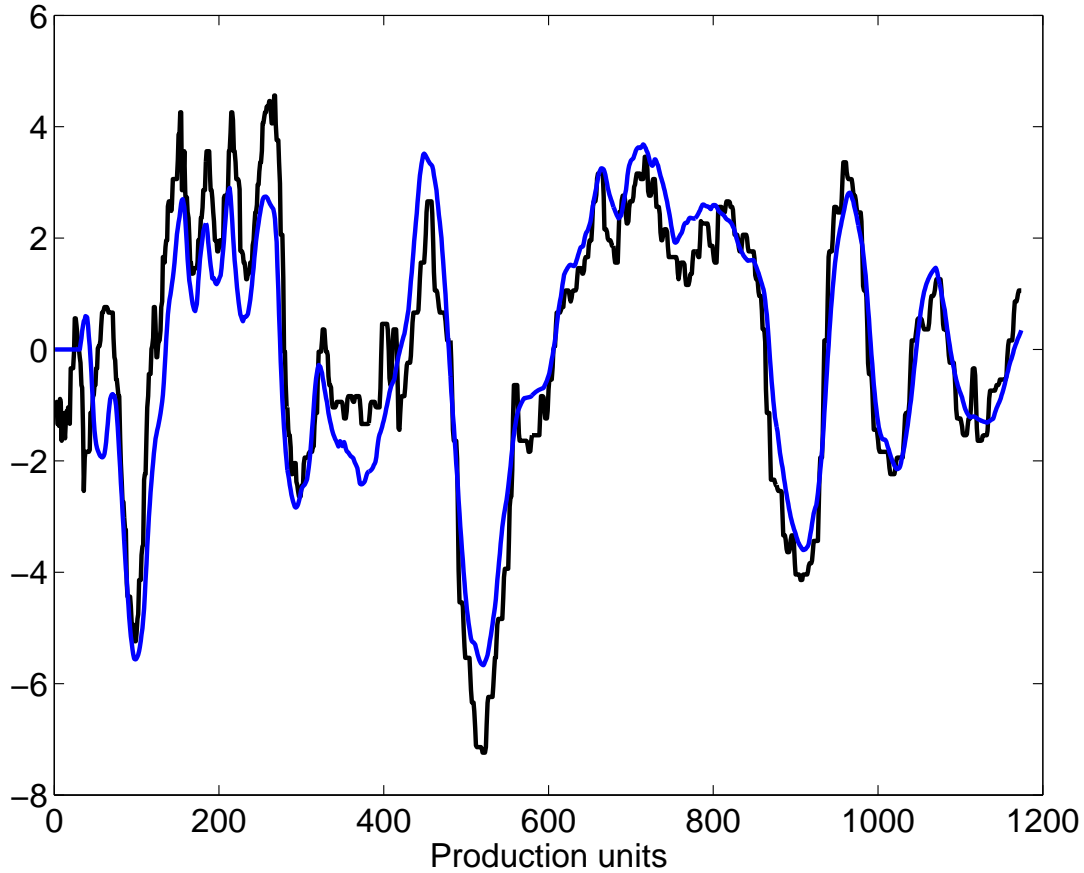


So: Resample Data!

```
z = [y,u]; pf = flow./level;  
t = 1:length(z)  
newt = interp1([cumsum(pf),t],[pf(1):sum(pf)]');  
newz = interp1([t,z], newt);
```

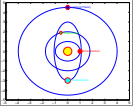


Measured and Model Output



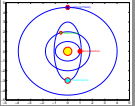
$$G(s) = \frac{0.8116}{1+110.28s} e^{-369.58s}$$





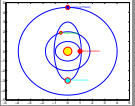
- Another satellite encircling the core.
- Deals with mathematical models of dynamic systems.





- Another satellite encircling the core.
- Deals with mathematical models of dynamic systems.
- **Typical themes:**
 - Useful model structures
 - Adapt and adopt the core's fundamentals: Prediction error methods
 - Experiment design (make \mathcal{I} large)
 - with intended model use in mind ("identification for control")





- Another satellite encircling the core.
- Deals with mathematical models of dynamic systems.
- **Typical themes:**
 - Useful model structures
 - Adapt and adopt the core's fundamentals: Prediction error methods
 - Experiment design (make \mathcal{I} large)
 - with intended model use in mind ("identification for control")

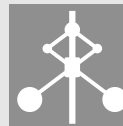




Machine Learning



Rather From Data to Action



Rather From Data to Action

Let us first watch a thought-provoking video (from Carl Rasmussen at Cambridge). It takes 1.5 minutes.

...



Rather From Data to Action

Let us first watch a thought-provoking video (from Carl Rasmussen at Cambridge). It takes 1.5 minutes.

...

What is going on?

At time t , let the position of the pendulum be $y(t)$ and the control action (the force on the cart) $u(t)$.

Somehow we need to understand the mapping $u(t) \Rightarrow y(t)$.



Rather From Data to Action

Let us first watch a thought-provoking video (from Carl Rasmussen at Cambridge). It takes 1.5 minutes.

...

What is going on?

At time t , let the position of the pendulum be $y(t)$ and the control action (the force on the cart) $u(t)$.

Somehow we need to understand the mapping $u(t) \Rightarrow y(t)$.

But the result of the action depends on where the pendulum is, so if

$x(t)$: pendulum's angle, angular velocity, cart position and cart velocity

we need to find a function $y(t) = f(x(t), u(t))$



Rather From Data to Action

Let us first watch a thought-provoking video (from Carl Rasmussen at Cambridge). It takes 1.5 minutes.

...

What is going on?

At time t , let the position of the pendulum be $y(t)$ and the control action (the force on the cart) $u(t)$.

Somehow we need to understand the mapping $u(t) \Rightarrow y(t)$.

But the result of the action depends on where the pendulum is, so if

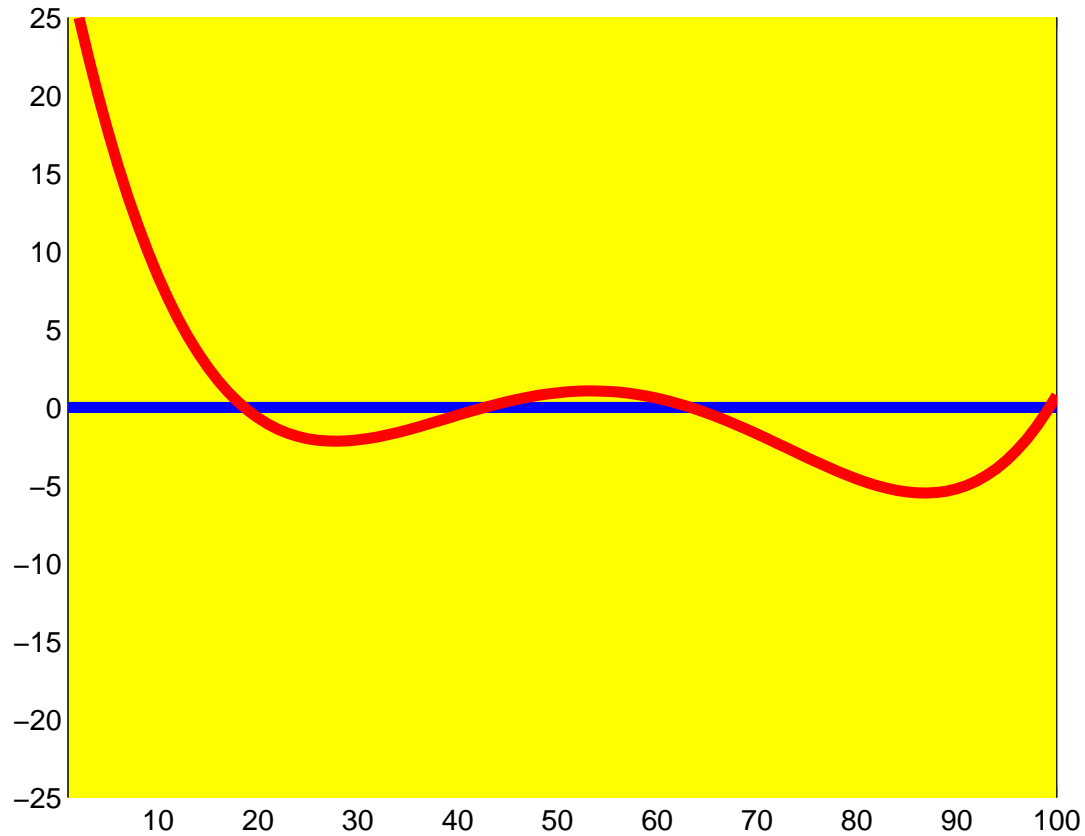
$x(t)$: pendulum's angle, angular velocity, cart position and cart velocity

we need to find a function $y(t) = f(x(t), u(t))$

How to estimate a function f ?

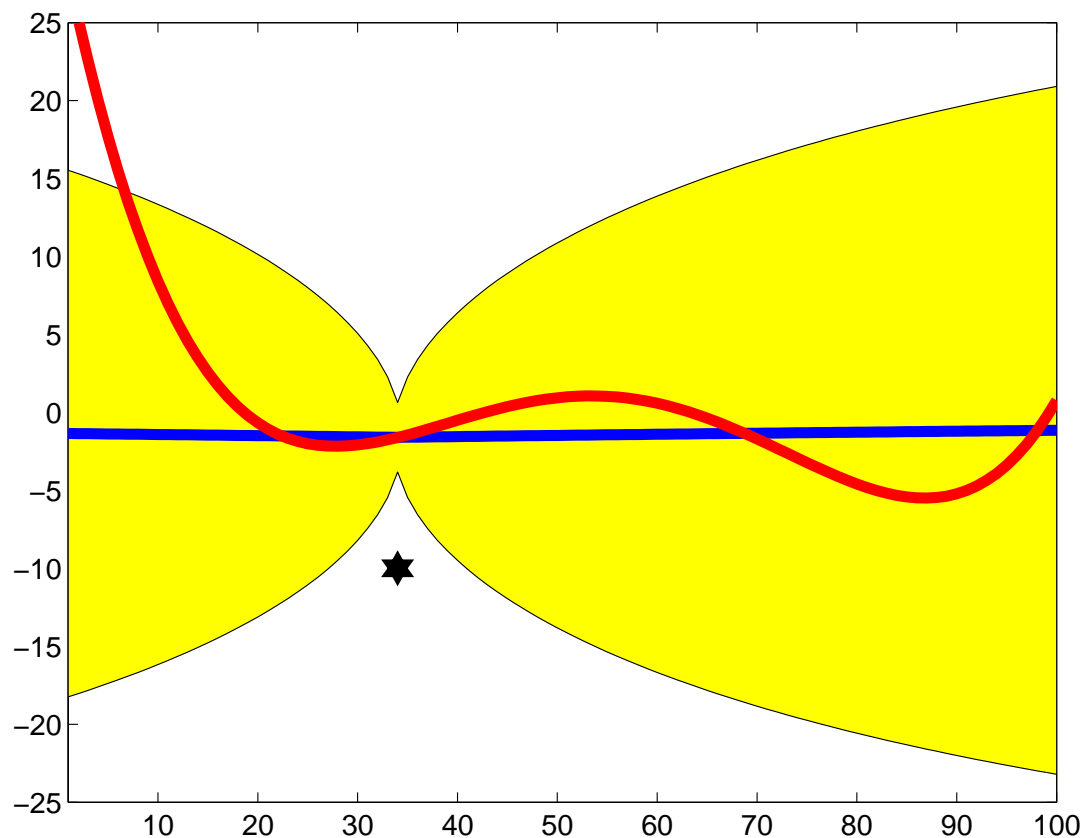


Function (curve) Estimation

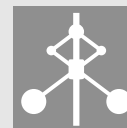


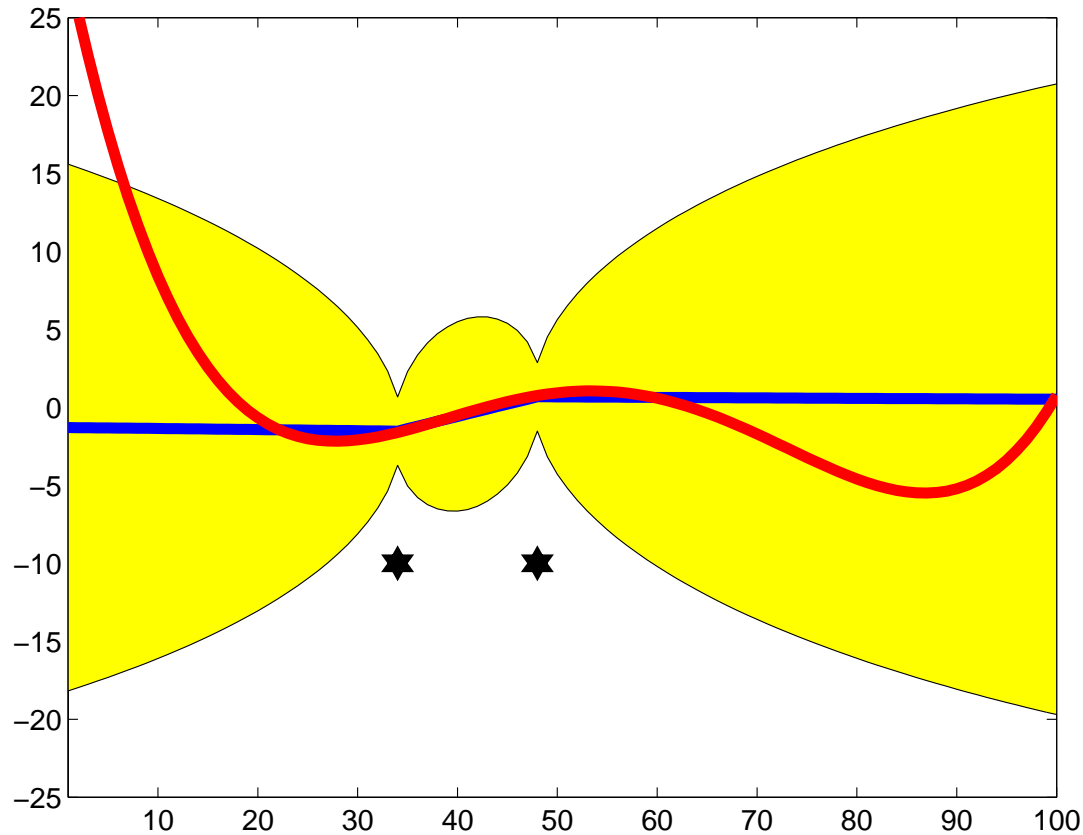
red: The unknown curve
blue: our current guess
yellow: our uncertainty





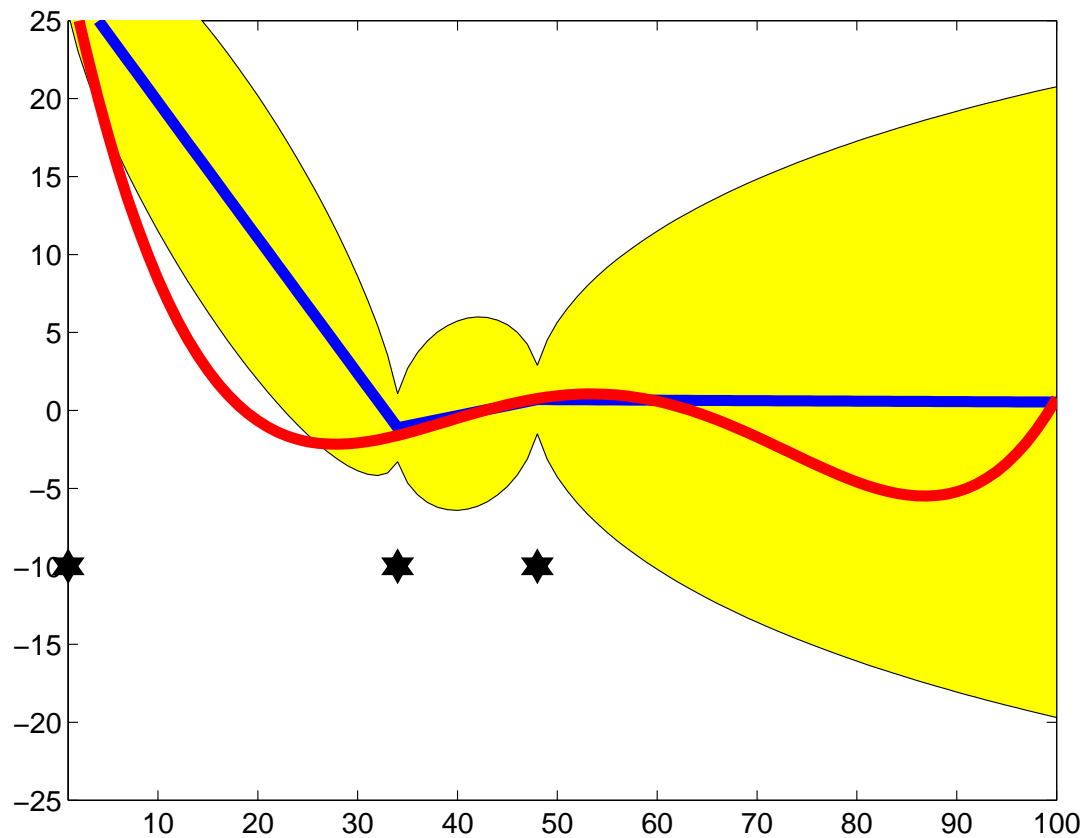
star: observation
blue: smooth curve
with certain correlation
yellow: our uncertainty
The info from the measurement propagates
due to the correlation



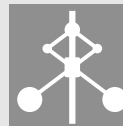


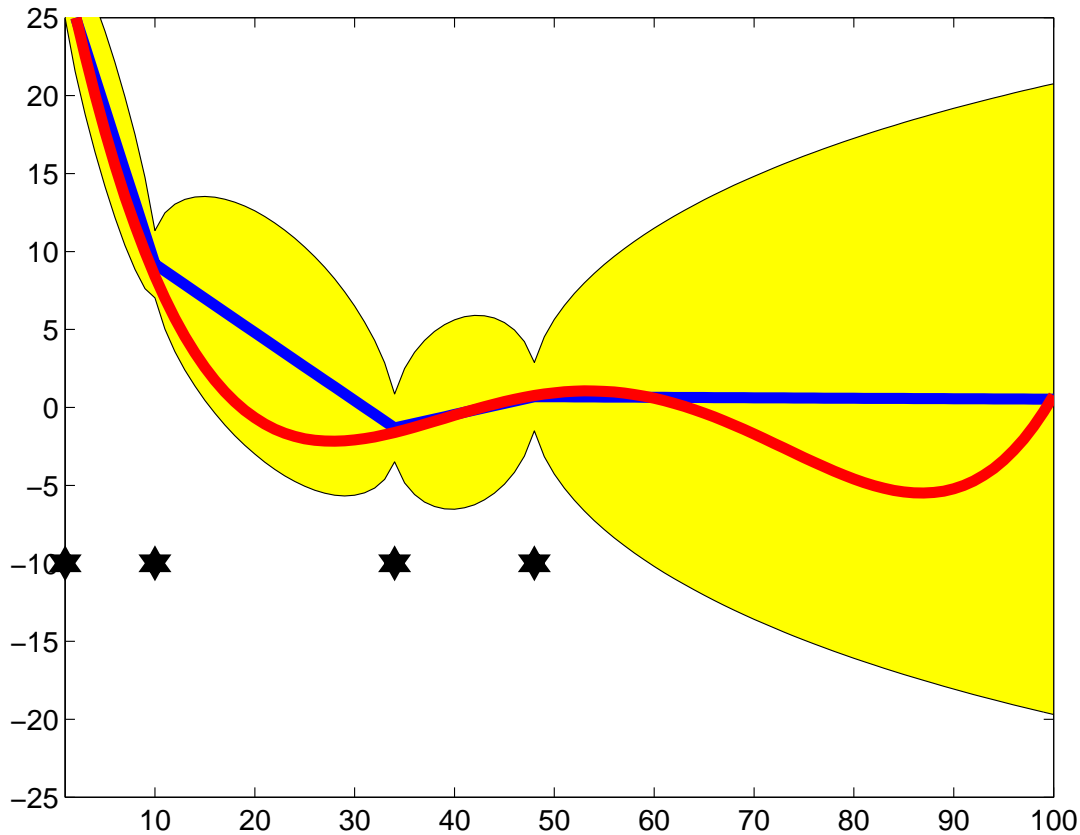
star: observation
blue: smooth curve
with certain correlation
yellow: our uncertainty
The info from the measurement propagates
due to the correlation





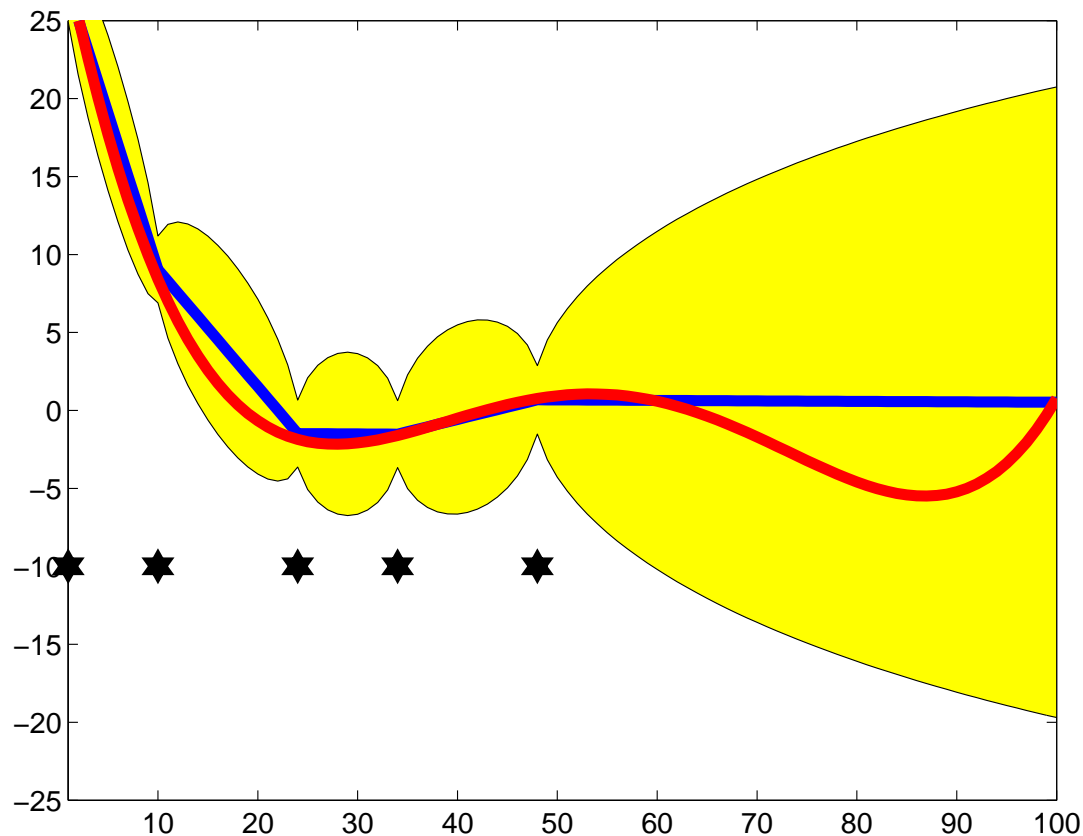
star: observation
blue: smooth curve
with certain correlation
yellow: our uncertainty
The info from the measurement propagates due to the correlation



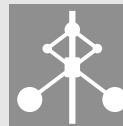


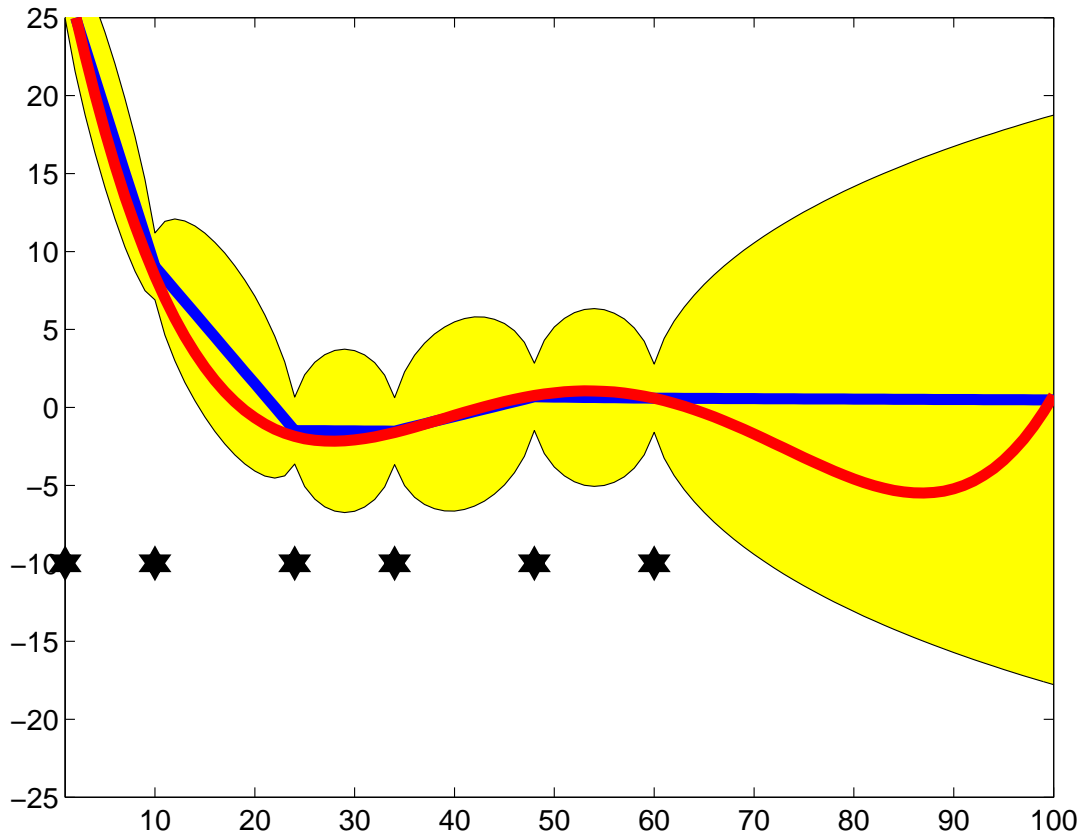
star: observation
blue: smooth curve
with certain correlation
yellow: our uncertainty
The info from the measurement propagates due to the correlation





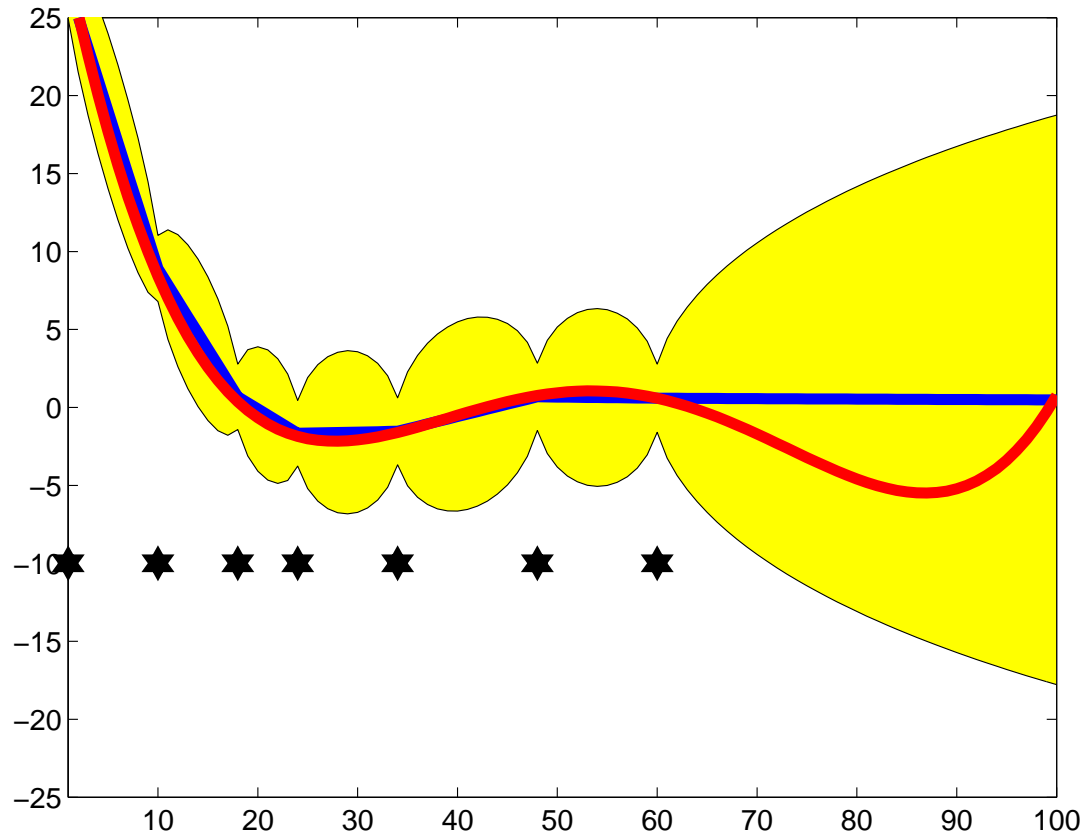
star: observation
blue: smooth curve
with certain correlation
yellow: our uncertainty
The info from the measurement propagates
due to the correlation





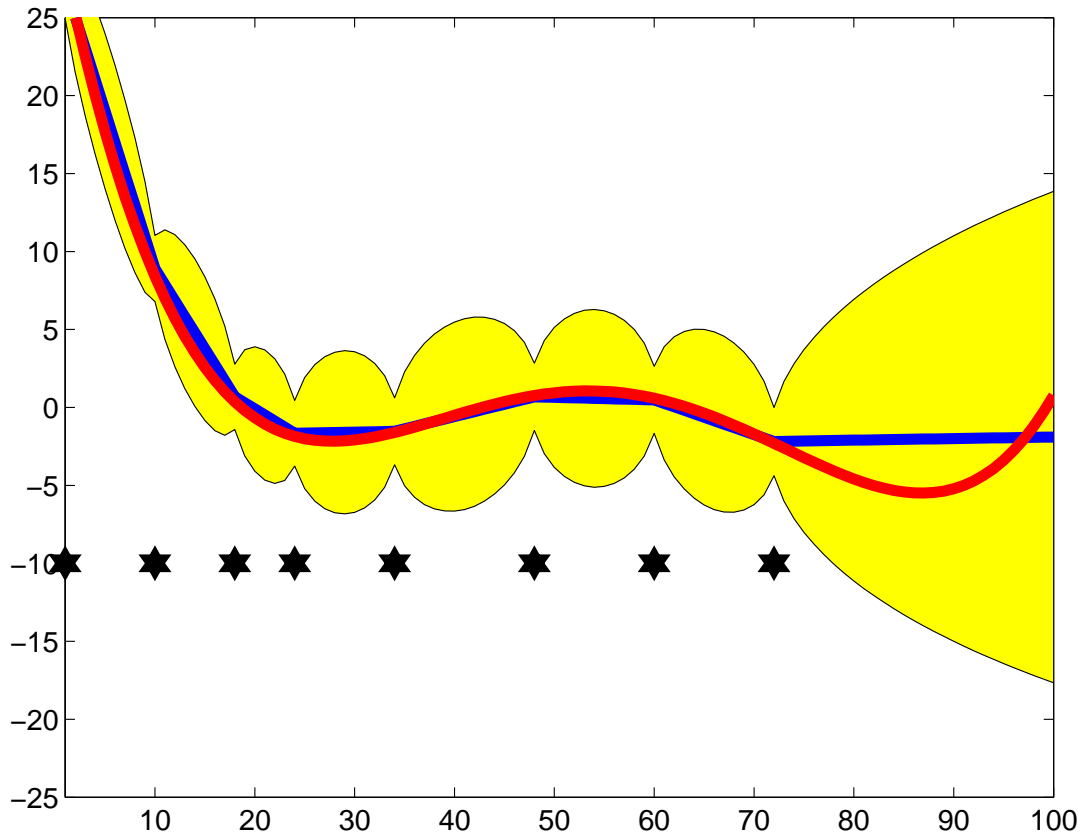
star: observation
blue: smooth curve
with certain correlation
yellow: our uncertainty
The info from the measurement propagates due to the correlation





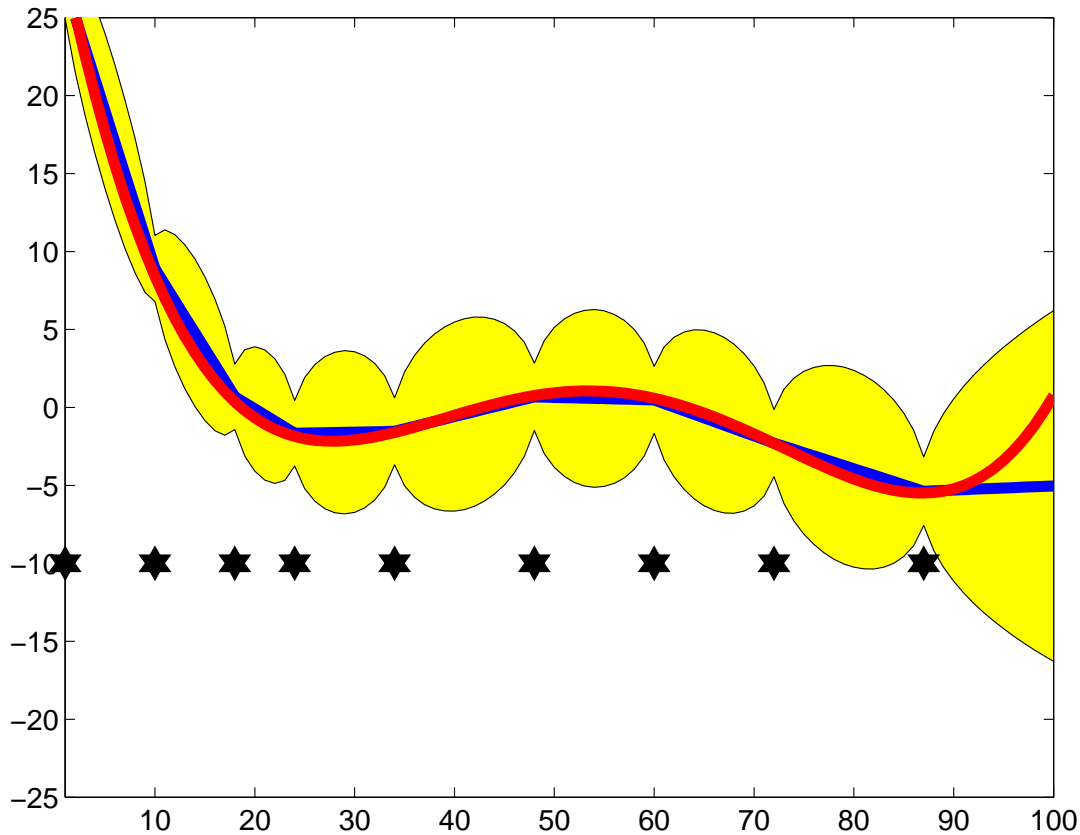
star: observation
blue: smooth curve
with certain correlation
yellow: our uncertainty
The info from the measurement propagates
due to the correlation





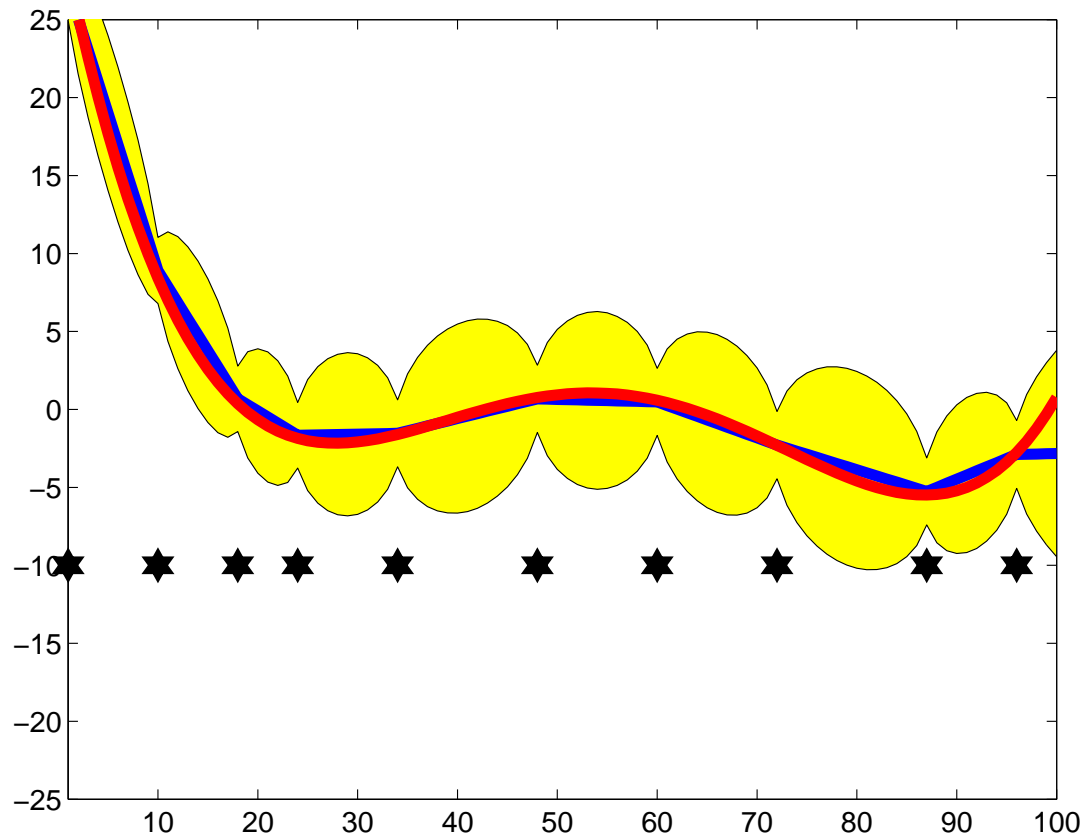
star: observation
blue: smooth curve
with certain correlation
yellow: our uncertainty
The info from the measurement propagates due to the correlation





star: observation
blue: smooth curve
with certain correlation
yellow: our uncertainty
The info from the measurement propagates due to the correlation





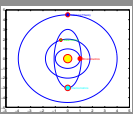
We have a good estimate of the curve with 10 observations (which are pretty dense compared to the 880 observations in the 5-dimensional space of the pendulum)



BTW, finding the **function f that predicts** the next output (pendulum position) from the current input (force on cart) and “state” (pendulum and cart movement) is the same problem as a **prediction error method**.

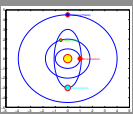
It is just that we have no parameterization of the prediction function, and don't want to introduce any physical knowledge, but just estimate it as a (smooth) curve.





- Very rich literature on building models from data (“The communities”)

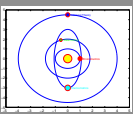




Conclusions

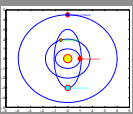
- Very rich literature on building models from data (“The communities”)
- Relatively few leading principles (“The core”)





- Very rich literature on building models from data (“The communities”)
- Relatively few leading principles (“The core”)
- Basic Principles are “classical”
 - Prediction error identification \Leftarrow K.F. Gauss (1777-1855)
 - Pendulum – Prediction function learning \Leftarrow T. Bayes (1702-1761)





- Very rich literature on building models from data (“The communities”)
- Relatively few leading principles (“The core”)
- Basic Principles are “classical”
 - Prediction error identification \Leftarrow K.F. Gauss (1777-1855)
 - Pendulum – Prediction function learning \Leftarrow T. Bayes (1702-1761)
- **Future:**
 - Can't beat the theoretical limits
 - Can use (massively) more data
 - More “data mining” in model building
 - More contacts between the Communities



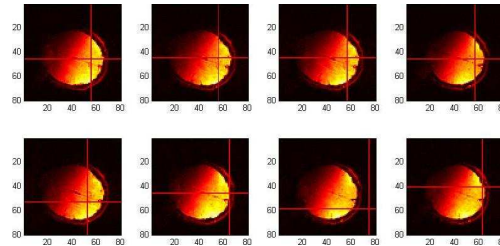
Find function for predicting next observation!



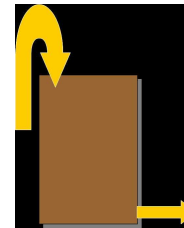
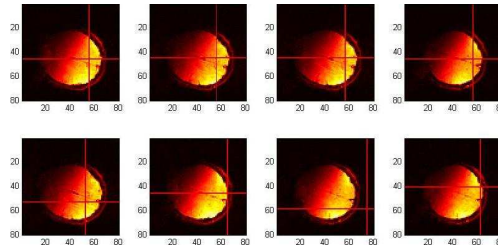
Find function for predicting next observation!



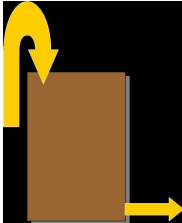
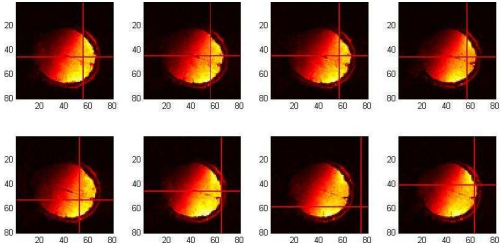
Find function for predicting next observation!



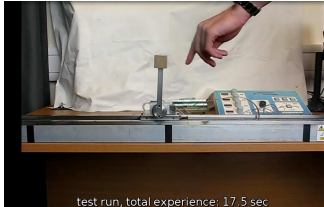
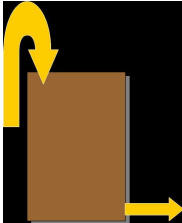
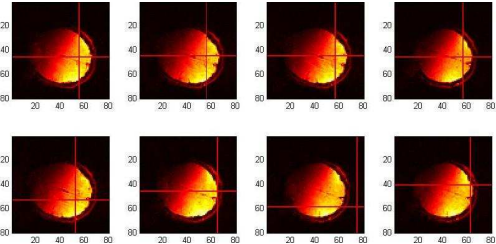
Find function for predicting next observation!



Find function for predicting next observation!



Find function for predicting next observation!



The IEEE International Conference on
**Industrial Engineering and
Engineering Management**
7 TO 10 DECEMBER 2010

