# An Introduction to
# Data Assimilation

Presenter: Yuxi Zhang

# Introduction

Assume:

A free-falling object with initial velocity of
$$v_0 = 1 \pm 0.3 \text{ m/s}$$

At $t = 1$ s, what is the speed of the object?

It is simple to calculate with the acceleration formula:
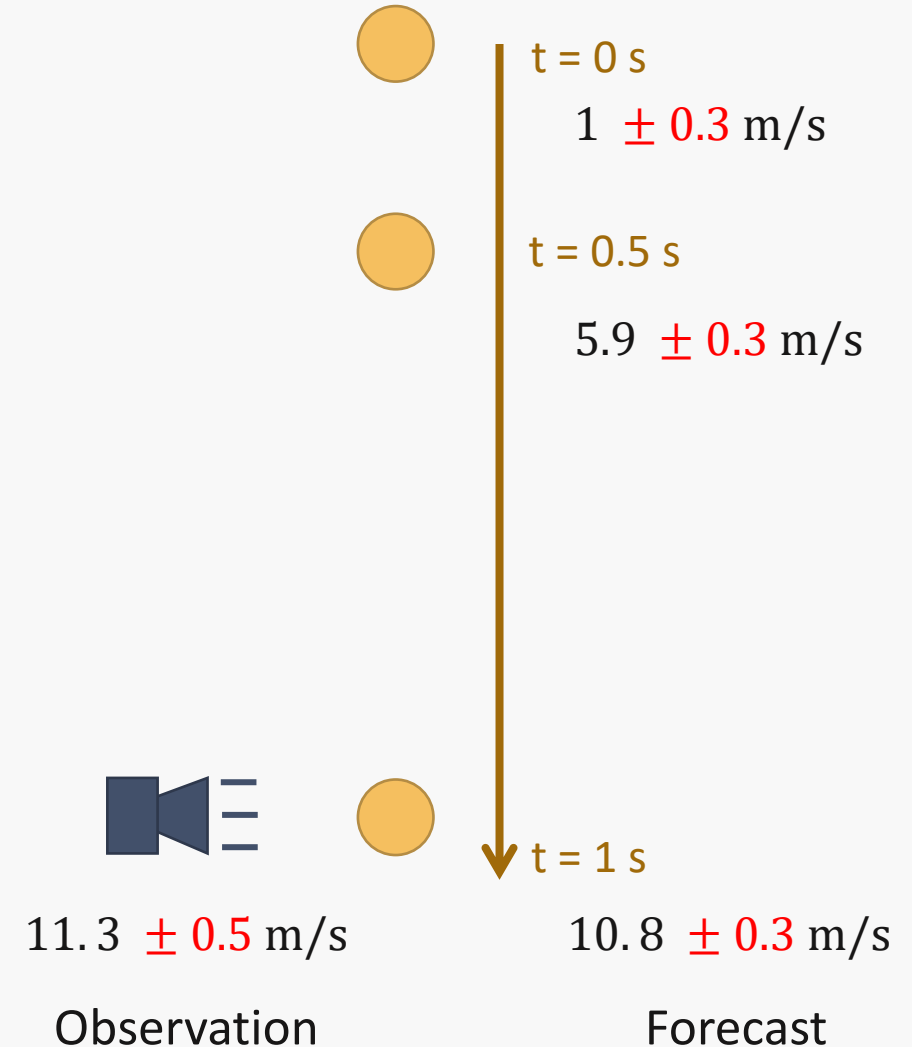$$v_{t=1} = v_0 + gt = 1 \,(\pm 0.3) + 9.8 \times 1 = 10.8 \,(\pm 0.3) \text{ m/s}.$$

Now assume we have an additional information:

At $t = 1$ s, a sensor measured the speed of the object as $11.3 \pm 0.5$ m/s.

In this case, how do you estimate the true velocity of the object at $t = 1$ s?

Which data do you trust?

**The purpose of data assimilation is to find the best estimate of the 'true' state.**

t = 0 s

$1 \pm 0.3$ m/s

t = 0.5 s

$5.9 \pm 0.3$ m/s

t = 1 s

$11.3 \pm 0.5$ m/s

Observation

$10.8 \pm 0.3$ m/s

Forecast

# Introduction

Forecast: $x_f = 10.8, \sigma_f = 0.3$

Observation: $x_o = 11.3, \sigma_o = 0.5$
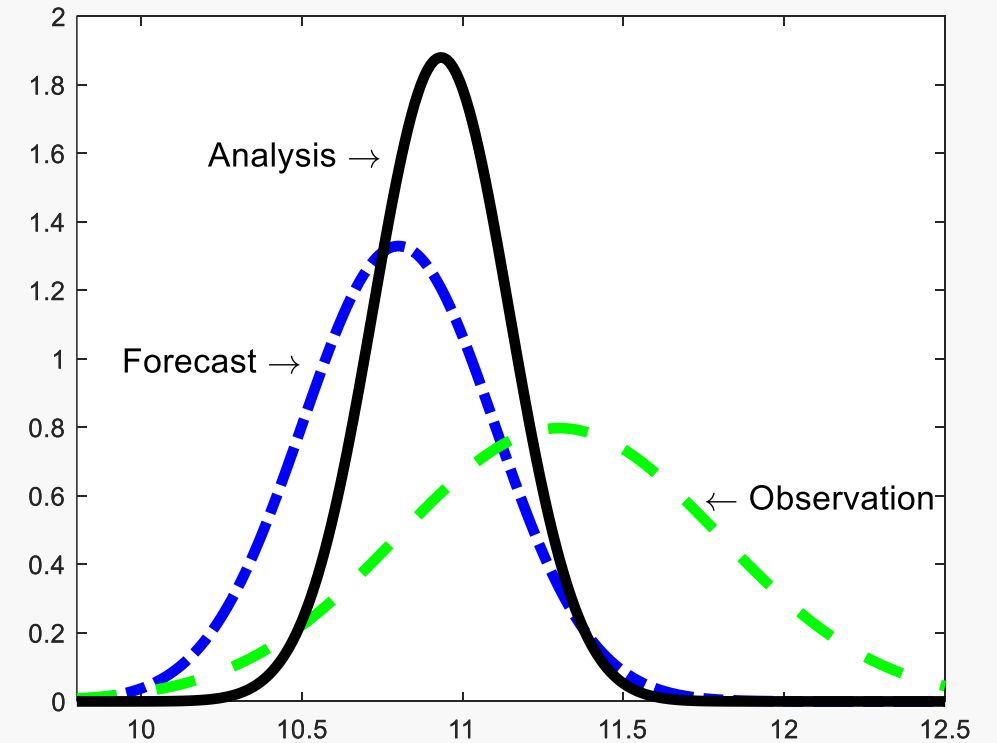
**Best Linear Unbiased Estimator (BLUE)**

Analysis:

$$x_a = x_f + \frac{\sigma_f^2}{\sigma_f^2 + \sigma_o^2}(x_o - x_f)$$

$$\sigma_a^{-2} = \sigma_f^{-2} + \sigma_o^{-2}$$

Thus,

$$x_a = 10.8 + \frac{0.3^2}{0.3^3 + 0.5^2}(11.3 - 10.8) = 10.93$$

$$\sigma_a = \frac{1}{1/0.3^2 + 1/0.5^2} = 0.21$$
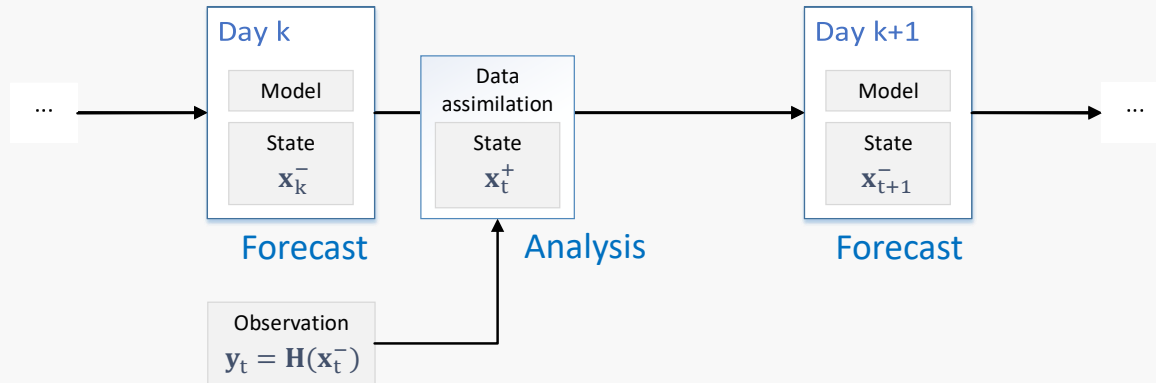


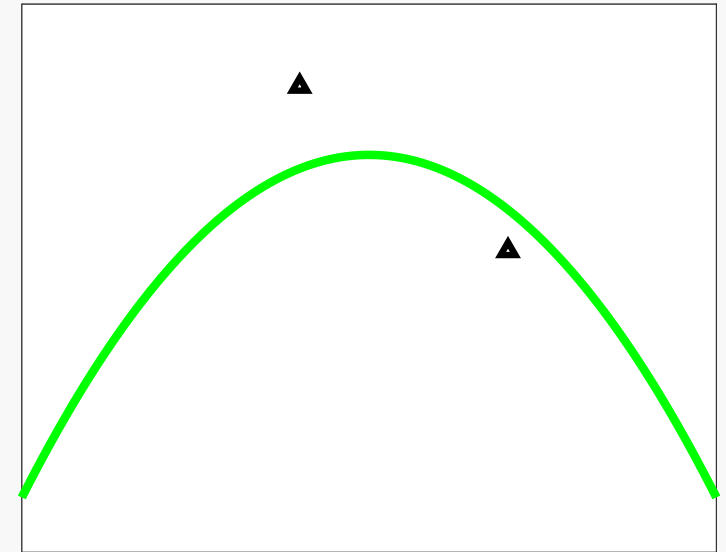$10.93 \pm 0.21$ m/s

Analysis

$11.3 \pm 0.5$ m/s          $10.8 \pm 0.3$ m/s

Observation          Forecast

# Concepts



Day k

| Model |
| State $\mathbf{x}_k^-$ |

Forecast

Data assimilation

| State $\mathbf{x}_t^+$ |

Analysis

Day k+1

| Model |
| State $\mathbf{x}_{t+1}^-$ |

Forecast
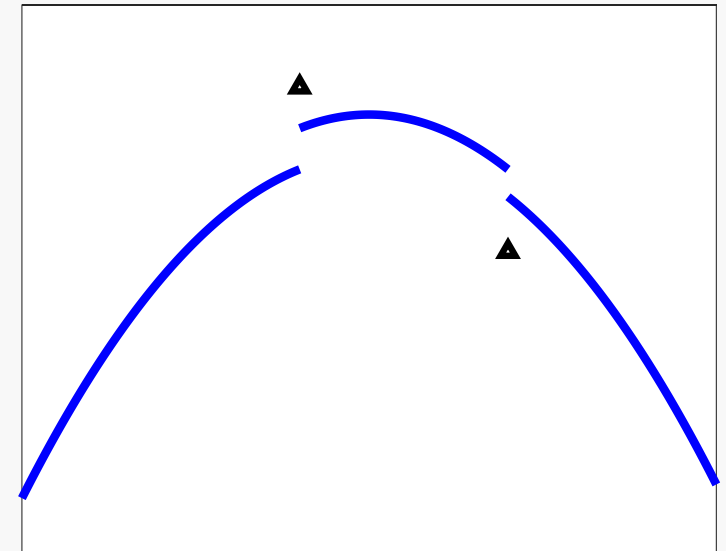
Observation
$$\mathbf{y}_t = \mathbf{H}(\mathbf{x}_t^-)$$

➢ State: $\mathbf{x}$

➢ Observation: $\mathbf{y}$

➢ State transfer function: $\mathbf{x}_t = M(\mathbf{x}_{t-1})$

➢ Error covariance: $\mathbf{P}$

➢ Prior/background/forecast: $\mathbf{x}^-, \mathbf{P}^-$

➢ Posterior/analysis: $\mathbf{x}^+, \mathbf{P}^+$



Open-loop



Closed-loop (data assimilation)

# Kalman filter

**Hypothesis**

Linear system with uncertainty:

$$\begin{aligned} \mathbf{x}_t &= M(\mathbf{x}_{t-1}) + \mathbf{w}_t \\ &= \mathbf{A}_t \mathbf{x}_{t-1} + \mathbf{B}_t \mathbf{f}_t + \mathbf{C}_t + \mathbf{w}_t \end{aligned}$$

Observations:

$$\begin{aligned} \mathbf{y}_t &= H(\mathbf{x}_t) + \mathbf{v}_t \\ &= \mathbf{H}_t \mathbf{x}_t + \mathbf{v}_t \end{aligned}$$

where $\mathbf{w}_t$ and $\mathbf{v}_t$ are independent uncertainties, following Gaussian distribution with a mean of zero and variance of $\mathbf{Q}$ and $\mathbf{R}$, respectively.

Mathematically,
$$\begin{aligned} \mathbf{w} &\sim \mathrm{N}(0, \mathbf{Q}) \\ \mathbf{v} &\sim \mathrm{N}(0, \mathbf{R}) \end{aligned}$$

$\mathbf{Q}$ and $\mathbf{R}$ are the variance of model and observational errors, respectively.

# Kalman filter

**Forecast**:

$$\mathbf{x}_t^- = M_t(\mathbf{x}_{t-1}^+)$$

The **prior** error:

$$
\begin{aligned}
\mathbf{e}_t^- &= \mathbf{x}_t^- - \mathbf{x}_t \\
&= M_t(\mathbf{x}_{t-1}^+) - M_t(\mathbf{x}_{t-1}) - \mathbf{w}_t \\
&= \mathbf{A}_t(\mathbf{x}_{t-1}^+ - \mathbf{x}_{t-1}) - \mathbf{w}_t \\
&= \mathbf{A}_t \mathbf{e}_{t-1}^+ - \mathbf{w}_t
\end{aligned}
$$

The **prior** error covariance

$$
\begin{aligned}
\mathbf{P}_t^- &= \mathrm{cov}(\mathbf{e}_t^-, \mathbf{e}_t^-) \\
&= \mathrm{E}\big(\mathbf{e}_t^- \mathbf{e}_t^{-\mathrm{T}}\big) \\
&= \mathbf{A}_t \mathbf{P}_{t-1}^+ \mathbf{A}_t^{\mathrm{T}} + \mathbf{Q}
\end{aligned}
$$

This equation links the prior error covariance to the posterior error covariance of the previous step.

**Analyse:**

$$
\begin{aligned}
\mathbf{y}_t &= \mathbf{H}_t \mathbf{x}_t + \mathbf{v}_t \\
\mathbf{x}_t^+ &= \mathbf{x}_t^- + \mathbf{K}_t(\mathbf{y}_t - \mathbf{H}_t \mathbf{x}_t^-)
\end{aligned}
$$

The **posterior** error:

$$
\begin{aligned}
\mathbf{e}_t^+ &= \mathbf{x}_t^+ - \mathbf{x}_t \\
&= \mathbf{x}_t^- + \mathbf{K}_t(\mathbf{H}_t \mathbf{x}_t + \mathbf{v}_t - \mathbf{H}_t \mathbf{x}_t^-) - \mathbf{x}_t \\
&= \mathbf{K}_t \mathbf{v}_t + (\mathbf{1} - \mathbf{K}_t \mathbf{H}_t)(\mathbf{x}_t^- - \mathbf{x}_t) \\
&= \mathbf{K}_t \mathbf{v}_t + (\mathbf{1} - \mathbf{K}_t \mathbf{H}_t) \mathbf{e}_t^-
\end{aligned}
$$

The **posterior** error covariance

$$
\begin{aligned}
\mathbf{P}_t^+ &= \mathrm{cov}(\mathbf{e}_t^+, \mathbf{e}_t^+) \\
&= \mathrm{E}\big(\mathbf{e}_t^+ \mathbf{e}_t^{+\mathrm{T}}\big) \\
&= (\mathbf{1} - \mathbf{K}_t \mathbf{H}_t) \mathbf{P}_t^- (\mathbf{1} - \mathbf{K}_t \mathbf{H}_t)^{\mathrm{T}} + \mathbf{K}_t \mathbf{R} \mathbf{K}_t^{\mathrm{T}}
\end{aligned}
$$

This equation links posterior error covariance to the prior error covariance of the same step.

# Kalman filter

In KF, the Kalman Gain ($\mathbf{K}_t$) is solved by minimising the squared difference between posterior estimates and the 'true' state, that is, the posterior error covariance ($\mathbf{P}_t^+$).

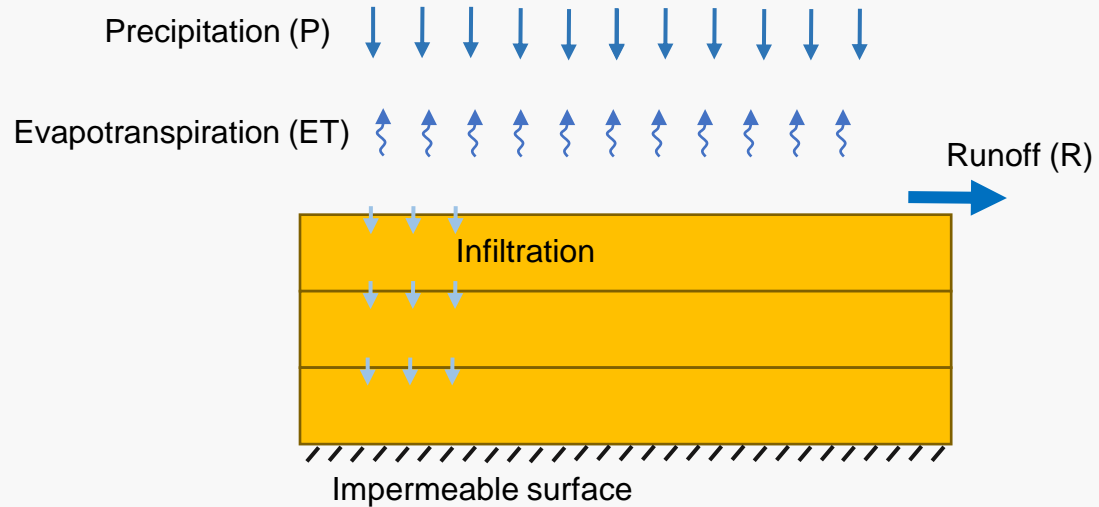By minimising $\mathbf{P}_t^+$, $\mathbf{K}_t$ is solved as

$$\mathbf{K}_t^* = \mathbf{P}_t^- \mathbf{H}_t^T \left( \mathbf{H}_t \mathbf{P}_t^- \mathbf{H}_t^T + \mathbf{R}_t \right)^{-1}$$

Once $\mathbf{K}_t$ is calculated, update the state and error covariance.

$$\mathbf{x}_t^+ = \mathbf{x}_t^- + \mathbf{K}_t (\mathbf{y}_t - \mathbf{H}_t \mathbf{x}_t^-)$$

$$\mathbf{P}_t^+ = (\mathbf{I} - \mathbf{K}_t^* \mathbf{H}_t) \, \mathbf{P}_t^-$$

# A simple water balance model

Precipitation (P)

Evapotranspiration (ET)

Runoff (R)

Infiltration

Impermeable surface

Assume depth = 1 m, the water balance model can be written as

Layer 1: $\theta(1)_{t+1} = \theta(1)_t + P_t - ET_t - R_t - I(1)_t$

Layer 2: $\theta(2)_{t+1} = \theta(2)_t + I(1)_t - I(2)_t$

Layer 3: $\theta(3)_{t+1} = \theta(3)_t + I(2)_t$

Rewrite in matrix

$$\underbrace{\begin{bmatrix} \theta(1) \\ \theta(2) \\ \theta(3) \end{bmatrix}_{t+1}}_{\mathbf{x_{t+1}}} = \underbrace{\begin{bmatrix} \theta(1) \\ \theta(2) \\ \theta(3) \end{bmatrix}_{t}}_{\mathbf{x_t}} + \begin{bmatrix} P_t - ET_t - R_t \\ 0 \\ 0 \end{bmatrix}_t + \begin{bmatrix} I(1) \\ I(1) - I(2) \\ I(2) \end{bmatrix}_t$$

State transfer function

$$\underbrace{\begin{bmatrix} \theta(1) \\ \theta(2) \\ \theta(3) \end{bmatrix}_{t+1}}_{\mathbf{x_{t+1}}} = \underbrace{\begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix}}_{A_t} \underbrace{\begin{bmatrix} \theta(1) \\ \theta(2) \\ \theta(3) \end{bmatrix}_{t}}_{\mathbf{x_t}} + \underbrace{\begin{bmatrix} 1 & -1 & -1 \\ & & \\ & & \end{bmatrix}}_{B_t} \underbrace{\begin{bmatrix} P_t \\ ET_t \\ R_t \end{bmatrix}_t}_{\text{weather}_t} + \underbrace{\begin{bmatrix} I(1) \\ I(1) - I(2) \\ I(2) \end{bmatrix}_t}_{C_t}$$

# Kalman filter

**Formula**

1. Define the model that runs one step forward.

$$\begin{bmatrix} \theta(1) \\ \theta(2) \\ \theta(3) \end{bmatrix}_{t+1} = \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} \begin{bmatrix} \theta(1) \\ \theta(2) \\ \theta(3) \end{bmatrix}_{t} + \begin{bmatrix} 1 & -1 & -1 \end{bmatrix} \begin{bmatrix} P_t \\ ET_t \\ R_t \end{bmatrix}_{t} + \begin{bmatrix} 1 & \\ 1 & -1 \\ & 1 \end{bmatrix} \begin{bmatrix} I(1) \\ I(2) \end{bmatrix}_{t}$$

$$\mathbf{x}_{t+1} = A_t \quad \mathbf{x}_t + B_t \; \text{weather}_t + C_t \; \text{parameter}_t$$

**MATLAB code**

HydroModel.m    ×    run_KF.m    ×    run_EnKF.m    ×    run_OL_deterministic.m    ×    run_OL_stochastic.m    ×    +

```matlab
1    function [state_out] = HydroModel(state_in, weather)
2
3        % Define matrices in model
4        A = [1 0 0; 0 1 0; 0 0 1];
5        B = [1 -1 -1; 0 0 0; 0 0 0] ;
6        C = [1 0; 1 -1; 0 1];
7        parameter = [20 10]';
8        state_out = A * state_in + B * weather * 0.001 + C * parameter  * 0.001;
9
10       % Constrain range
11       state_out(state_out<0.1) = 0.1;
12       state_out(state_out>0.9) = 0.9;
```

# Kalman filter

## Formula

2. Define observational operator.

$$\begin{bmatrix} y(1) \\ y(2) \end{bmatrix}_t = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}_t \begin{bmatrix} \theta(1) \\ \theta(2) \\ \theta(3) \end{bmatrix}_t$$

$$\mathbf{y}_t \quad = \quad \mathbf{H}_t \quad \cdot \quad \mathbf{x}_t$$

3. Set initial conditions

$$\begin{bmatrix} \theta(1) \\ \theta(2) \\ \theta(3) \end{bmatrix}_0 = \begin{bmatrix} 0.3 \\ 0.3 \\ 0.3 \end{bmatrix}$$

## MATLAB code

HydroModel.m   run_KF.m   run_EnKF.m   run_OL_deterministic.m   run_OL_sto

```matlab
1    %% Loading data
2    clear
3    close all
4    cd 'H:\My Drive\Slides\Introduction_DA\scripts'
5    weathers = load('data\ExampleData.csv');
6    observations = load('data\Obs.csv');
7
8    %% Model initial conditions
9    state= [0.3 0.3 0.3]';
10
11   %% KF matrices
12   % Define observational matrix H.
13   H = [1 0 0; 0, 1, 0];
14   % Set model and observational error variance by estimation.
15   Q = [0.1 0 0; 0 0.1 0; 0 0 0.1].^2;
16   R = [0.1 0; 0, 0.1].^2;
17   % Set  initial error covariance by guess.
18   % Note that P, Q, R are symmetric matrices.
19   P = [0.1 0 0; 0 0.1 0; 0 0 0.1].^2;
```

# Kalman filter

## Formula

4. Run the model one step forward.

If observation at the current timestep is available, run data assimilation.

$$\mathbf{K}_t = \mathbf{P}_t^- \mathbf{H}_t^T \left( \mathbf{H}_t \mathbf{P}_t^- \mathbf{H}_t^T + \mathbf{R}_t \right)^{-1}$$

$$\mathbf{P}_t^+ = (\mathbf{I} - \mathbf{K}_t^* \mathbf{H}_t)\,\mathbf{P}_t^-$$

When $\mathbf{K}_t$ is calculated, the posterior state can be calculated.
$$\mathbf{x}_t^+ = \mathbf{x}_t^- + \mathbf{K}_t(\mathbf{y}_t - \mathbf{H}_t \mathbf{x}_t^-)$$

Then the prior state will be replaced by the posterior state values.

## MATLAB code

HydroModel.m    run_KF.m    run_EnKF.m    run_OL_deterministic.m    run_OL_stochastic

```matlab
21      %% Run the model in a loop.
22      days = length(weathers);
23      % Create empty matrices to save data.
24      state_prior = nan(length(state),days);
25      state_posterior = nan(size(state_prior));
26      P_all = nan([size(P) days]);
27      tmp = size(H);
28      K_all = nan(length(state), tmp(1), days);
29
30      for i = 1 : days
31          % get daily data
32          weather = weathers (i, :)';
33          obs = observations (i, :)';
34          state = HydroModel(state, weather);
35          state_prior(:, i) = state;
36
37          % Check if observation is available.
38          % If yes, run data assimilation to update the state vector.
39          if ~isnan(obs)
40              K = P * H' * (H * P * H' + R ) ^ (-1);
41              P = (1 - K * H) * P;
42              state = state + K * (obs - H * state);
43              % Save or export the values of theta, K, P if necessary.
44              P_all(:,:,i) = P;
45              K_all(:,:,i) = K;
46              disp('State is updated.')
47          else
48              disp('Observation is not available!')
49          end
50          state_posterior(:, i) = state;
51      end
```

# Non-liner models

The KF is based on the assumption of linear system.

In a general system, state transfer function can be written as:
$$\mathbf{x}_t = M(\mathbf{x}_{t-1}) + \mathbf{w}_t$$

Observation to the system:
$$\mathbf{y}_t = H(\mathbf{x}_t) + \mathbf{v}_t$$

where $\mathbf{x}$ and $\mathbf{y}$ are column vectors, and the subscriptions t and t+1 are the time step.

**Solutions for non-linear model**
    Extended Kalman filter, EKF
    Ensemble Kalman filter, EnKF
    Particle filter, PF

# Extended Kalman filter

**Hypothesis**

Model can be linearised by a Jacobians matrix.

# Ensemble Kalman filter

**Hypothesis**

The probability distribution of uncertainty can be represented by a finite number of simulations (ensemble).

The ensemble

$$\mathbf{X} = [\mathbf{x}^1 \quad \mathbf{x}^2 \quad ... \quad \mathbf{x}^N]$$

$\mathbf{X}$ is a M x N matrix, where M is the number of state in the state vector $\mathbf{x}$, N is the ensemble size.

**Forecast**

$$\mathbf{x}_t = M(\mathbf{x}_{t-1}) + \mathbf{w}_t$$

$$\mathbf{P}_t^- = \frac{1}{N-1} \mathbf{D}_t \mathbf{D}_t^T$$

where

$$\mathbf{D}_t^T = \left[ \mathbf{x}^1 - E[\mathbf{x}^1] \quad \mathbf{x}^2 - E[\mathbf{x}^2] \quad ... \quad \mathbf{x}^N - E[\mathbf{x}^N] \right]$$

**Analyse:**

$$\mathbf{y}_t = \mathbf{H}_t \mathbf{x}_t + \mathbf{v}_t$$

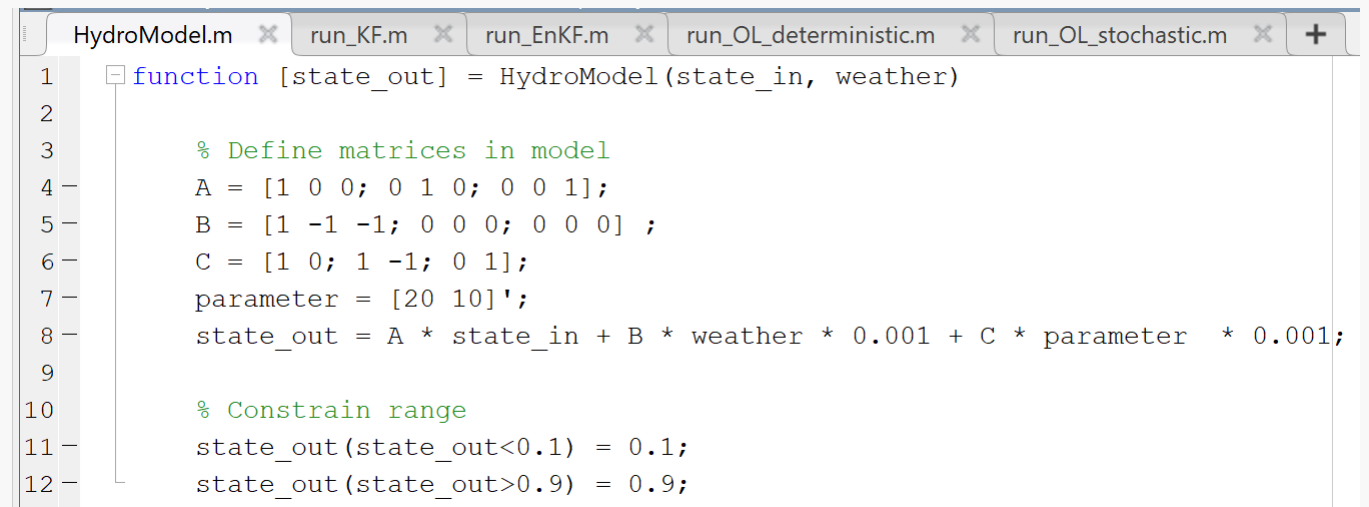$$\mathbf{x}_t^+ = \mathbf{x}_t^- + \mathbf{K}_t(\mathbf{y}_t - \mathbf{H}_t \mathbf{x}_t^-)$$

# Ensemble Kalman filter

## Formula

1. Define the model that runs one step forward.

$$\mathbf{x}_t = M(\mathbf{x}_{t-1})$$

**MATLAB code**

```matlab
HydroModel.m    run_KF.m    run_EnKF.m    run_OL_deterministic.m    run_OL_stochastic.m    +

1    function [state_out] = HydroModel(state_in, weather)
2
3        % Define matrices in model
4        A = [1 0 0; 0 1 0; 0 0 1];
5        B = [1 -1 -1; 0 0 0; 0 0 0] ;
6        C = [1 0; 1 -1; 0 1];
7        parameter = [20 10]';
8        state_out = A * state_in + B * weather * 0.001 + C * parameter  * 0.001;
9
10       % Constrain range
11       state_out(state_out<0.1) = 0.1;
12       state_out(state_out>0.9) = 0.9;
```

# Ensemble Kalman filter

## Formula

2. Set initial conditions

$$\begin{bmatrix} \theta(1) \\ \theta(2) \\ \theta(3) \end{bmatrix}_0 = \begin{bmatrix} 0.3 \\ 0.3 \\ 0.3 \end{bmatrix}$$

3. Define observational operator.

$$\begin{bmatrix} y(1) \\ y(2) \end{bmatrix}_t = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}_t \begin{bmatrix} \theta(1) \\ \theta(2) \\ \theta(3) \end{bmatrix}_t$$

$$\mathbf{y}_t \quad = \quad \mathbf{H}_t \cdot \mathbf{x}_t$$

4. Set uncertainty values for weather, initial conditions and observations.

## MATLAB code

```
HydroModel.m  ✕   run_KF.m  ✕   run_EnKF.m  ✕   run_OL_deterministic.m  ✕   run_OL_stochastic.m

1       %% Loading data
2   -   clear
3   -   close all
4   -   cd 'H:\My Drive\Slides\Introduction_DA\scripts'
5   -   weathers = load('data\ExampleData.csv');
6   -   observations = load('data\Obs.csv');
7
8       %% Model initial conditions
9   -   state= [0.3 0.3 0.3]';
10
11      %% EnKF settings.
12      % Define observational matrix H.
13  -   H = [1 0 0; 0, 1, 0];
14      % Set model and observational error variance by estimation.
15  -   R = [0.03 0; 0, 0.03].^2;
16      % Set uncertainty for weather, initial conditions and observations.
17  -   error_weather = [5 3 1]';
18  -   error_init = [0.1 0.1 0.1]';
19  -   error_obs = [0.03 0.03]';
```

# Kalman filter

## Formula

### MATLAB code

5. Generate ensembles. Construct a M x N state matrix, where M is the length of the state vector, and N is the ensemble size. Perturb initial conditions by adding Gaussian noise.

```matlab
21    %% Generate ensemble.
22    ensemble_size = 5;
23    % Perturb initial conditions.
24    X = repmat(state, 1, ensemble_size);
25    for i = 1:length(state)
26        X(i, :) = X(i, :) + randn(1, ensemble_size) * error_init(i);
27    end
```

6. Run the model **stochastically** in a loop.

- Perturb weather data by adding Gaussian noise. The matrix f a matrix with size of the length of the weather vector by the ensemble size.

- Run stochastic model. The estimated state is the mean of ensemble.

```matlab
HydroModel.m    run_KF.m    run_EnKF.m    run_OL_deterministic.m    run_OL_stochastic.m

40    %% Run the model in a loop.
41    for i = 1 : days
42        % Get daily data
43        weather = weathers (i, :)';
44        obs = observations (i, :)';
45        % Perturb weather forces. f is a matrix of no_of_weather by
46        % ensemble_size.
47        f = repmat(weather, 1, ensemble_size);
48        for j = 1:length(weather)
49            f(j, :) = f(j, :) + randn(1, ensemble_size) * error_weather(j);
50        end
51        % Run stochastic model. Note that X and f are matrices in EnKF rather
52        % than vectors in KF.
53        X = HydroModel(X, f);
54        X_prior_ensemble(:, :, i) = X;
55        X_prior(:, i) = mean(X, 2);
```

# Kalman filter

## Formula

7. If observation at the current timestep is available, run data assimilation.

- Perturb observations data by adding Gaussian noise. The matrix y a matrix with size of the length of the observation vector by the ensemble size.

- Calculate

$$\mathbf{D}_t^T = \begin{bmatrix} \mathbf{x}^1 - E[\mathbf{x}^1] & \mathbf{x}^2 - E[\mathbf{x}^2] & \dots & \mathbf{x}^N - E[\mathbf{x}^N] \end{bmatrix}$$

$$\mathbf{P}_t^- = \frac{1}{N-1} \mathbf{D}_t \mathbf{D}_t^T$$

$$\mathbf{K}_t = \mathbf{P}_t^- \mathbf{H}_t^T (\mathbf{H}_t \mathbf{P}_t^- \mathbf{H}_t^T + \mathbf{R}_t)^{-1}$$
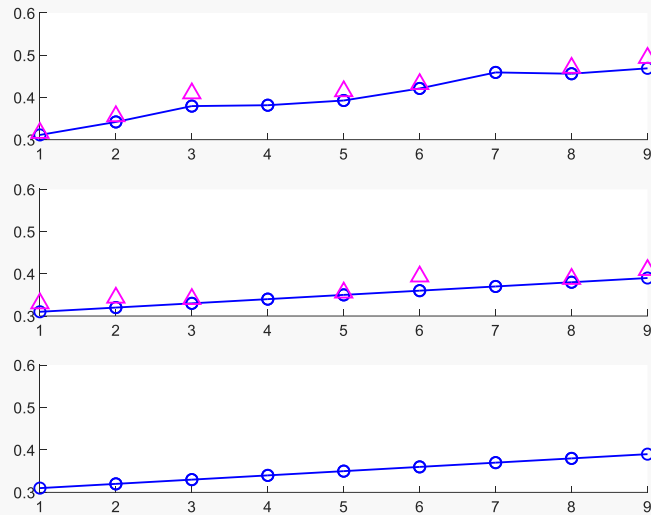
Once $\mathbf{K}_t$ is calculated, update the state and error covariance.

$$\mathbf{x}_t^+ = \mathbf{x}_t^- + \mathbf{K}_t (\mathbf{y}_t - \mathbf{H}_t \mathbf{x}_t^-)$$
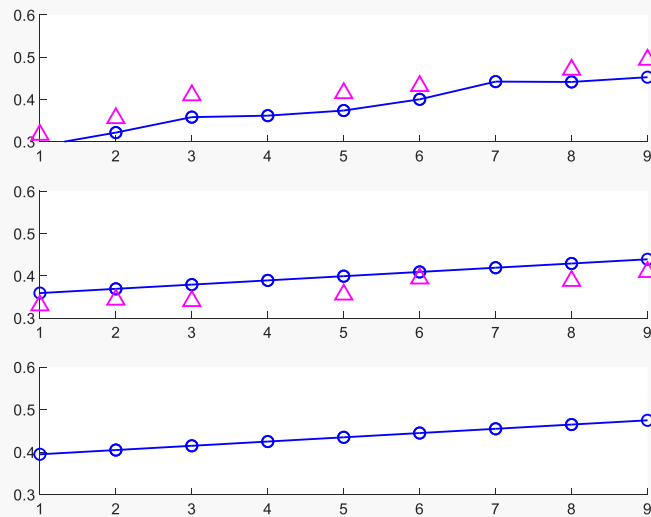
## MATLAB code

HydroModel.m  |  run_KF.m  |  run_EnKF.m  |  run_OL_deterministic.m  |  run_OL_stochastic.m

```matlab
57          % Check if observation is available.
58          % If yes, run data assimilation to update the state vector.
59          if ~isnan(obs)
60              % Perturb observation (y).
61              y = repmat(obs, 1, ensemble_size);
62              for j = 1:length(obs)
63                  y(j, :) = y(j, :) + randn(1, ensemble_size) * error_obs(j);
64              end
65              % Run EnKF
66              X_mean = mean(X, 2);
67              D = X - repmat(X_mean, 1, ensemble_size);
68              P = D * D' / (1 + ensemble_size);
69              K = P * H' * (H * P * H' + R) ^ (-1);
70              X = X + K * (y - H * X);
71              % Save or export the values of theta, K, P if necessary.
72              P_all(:,:,i) = P;
73              K_all(:,:,i) = K;
74              obs_ensemble(:, :, i) = y;
75              disp('State is updated.')
76          else
77              disp('Observation is not available!')
78          end
79          X_posterior_ensemble(:, :, i) = X;
80          X_posterior(:, i) = mean(X, 2);
81      end
```
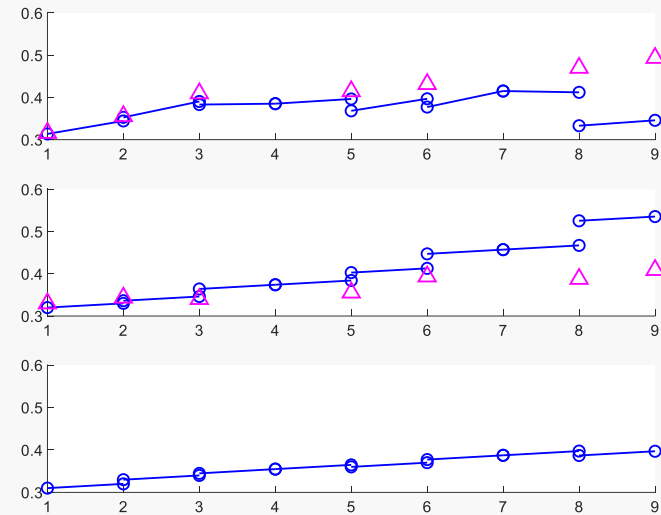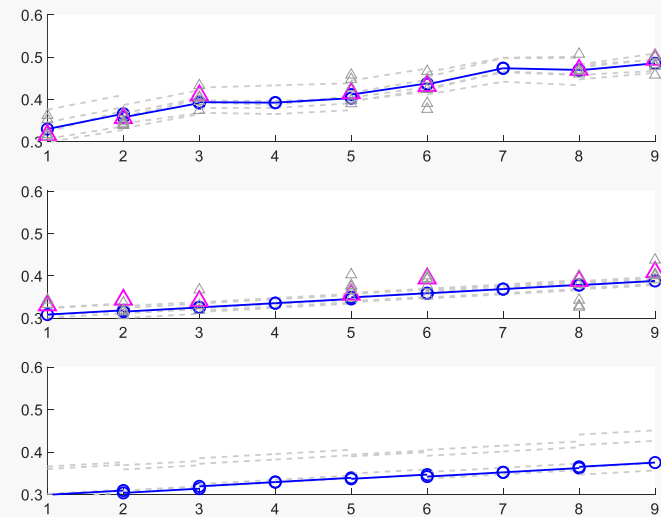
# Demo output



Open-loop, deterministic run

KF

Open-loop, stochastic run

EnKF