

Security and Limitations of Cyber-Physical Systems

Lecture 3: Defense Mechanisms

ROYAL INSTITUTE OF TECHNOLOGY

Henrik Sandberg André Teixeira

Department of Automatic Control ACCESS Linnaeus Centre, KTH Royal Institute of Technology Stockholm, Sweden

> Linköping University August 24-26, 2015



TECHNOLOGY

Course Outline

- Monday (8:30-10:30):
 - Lecture 1 (HS): Introduction, data attacks against nondynamic systems, power network monitoring, security index, graph min cut
- Tuesday
 - 8:30-12:30:
 - Lecture 2 (HS): Attack space for cyber-physical systems: DoS, undetectable, stealth, covert, bias, replay attacks
 - Lecture 3 (AT): Defense mechanisms, risk management, anomaly detectors, watermarking
 - 15:30-16:30:
 - Exercise session (Graph min cut, security index)
- Wednesday (8:30-10:30):
 - Lecture 4 (HS): Physical limits of control implementations ²
 - Exercise session



Key References for Lecture

- [1] André Teixeira, Kin Sou, Henrik Sandberg, and Karl H. Johansson. "Secure Control Systems: A Quantitative Risk Management Approach". IEEE Control System Magazine, vol. 35, no. 1, pp. 24-25, Feb. 2015
- [2] Inseok Hwang, Sungwan Kim, Youdan Kim, C.E. Seah: "A Survey of Fault Detection, Isolation, and Reconfiguration Methods," *Control Systems Technology, IEEE Transactions on*, vol.18, no.3, pp.636,653, May 2010
- [3] Yilin Mo, Sean Weerakkody, Bruno Sinopoli: "Physical Authentication of Control Systems". IEEE Control Systems Magazine, vol. 35, no. 1, pp. 93-109, February 2015.
- [4] André Teixeira, Iman Shames, Henrik Sandberg, and Karl H. Johansson, "Revealing Stealthy Attacks in Control Systems". In Proc. 50th Annual Allerton Conference on Communication, Control, and Computing, Allerton, IL, USA, 2012



Lecture 3

- Risk management
 - security metrics, quantifying impact [1]

Anomaly detectors (for detectable attacks) [2]

- Watermarking (against undetectable attacks)
 - Replay attacks [3]
 - Zero dynamics attacks [4]



F TECHNOLOGY

Motivation



- Too costly to secure the entire system against all attack scenarios
- What scenarios to prioritize?
- What components to protect?
- How to mitigate undetectable attacks?



Actuators

 $a \cdot$

Sensors

Physical Plant

Communication Network

Distributed Controllers

 c_1



Goals of Lecture 3

 Analyze and compare existing attack scenarios in a risk management framework

• Understand the basics of anomaly detection

• Design methods to reveal undetectable attacks



Defining Risk

ROYAL INSTITUTE OF TECHNOLOGY

Risk = (Scenario , Likelihood, Impact)

- Scenario
 - How to describe the system under attack?
 - (recall Lecture 2)
- Likelihood
 - How much effort does a given attack require?
 - (compare to security index, Lecture 1)
- Impact
 - What are the consequences of an attack?
 - (relate to control objectives)







OF TECHNOLOGY

Risk Management

- Main steps in risk management
 - Scope definition
 - Models, Scenarios, objectives
 - Risk Analysis
 - Threat Identification
 - Likelihood Assessment
 - Impact Assessment
 - Risk Treatment
 - Prevention, Detection, Mitigation





Risk Management

- How to model adversaries and attacks?
 Lecture 2
- How to measure likelihood (attack effort)?
 - Lectures 1 and 2
- How to compute impact?
 - Lecture 2
 - This lecture
- How to design protection and detection mechanisms?
 - This lecture





TECHNOLOGY

Risk Analysis for Dynamical Systems

• Analysis of effort and impact of stealthy attacks

- Cases considered here:
 1. Minimum resource attacks
 2. Maximum impact attacks
 - 3. Maximum impact bounded resource attacks



 Considered attacks are in open loop. No disclosure resources explicitly used (works due to linearity of systems)



Networked Control System

ROYAL INSTITUTE OF TECHNOLOGY



- Physical Plant $\mathcal{P}: \begin{cases} x_{k+1} = Ax_k + B\tilde{u}_k + Gw_k + Ff_k \\ y_k = Cx_k + v_k \end{cases}$
- Feedback Controller $\mathcal{F}: \begin{cases} z_{k+1} = A_c z_k + B_c \tilde{y}_k \\ u_k = C_c z_k + D_c \tilde{y}_k \end{cases}$
- Anomaly Detector $\mathcal{D}: \begin{cases} s_k = A_e s_k + B_e u_k + K_e \tilde{y}_k \\ r_k = C_e s_k + D_e u_k + E_e \tilde{y}_k \end{cases}$



Networked Control System under Attack

ROYAL INSTITUTE OF TECHNOLOGY

- Cyber-Physical Attacks
 - Disclosure Component

 $\mathcal{I}_k := \mathcal{I}_{k-1} \cup \begin{bmatrix} \Upsilon^u & 0\\ 0 & \Upsilon^y \end{bmatrix} \begin{bmatrix} u_k\\ y_k \end{bmatrix}$

- Disruptive component

 $a_k = [f_k^\top \quad b_k^{u\top} \quad b_k^{y\top}]^\top$



- Closed-loop Dynamics $\eta_k = \begin{bmatrix} x_k^\top & z_k^\top \end{bmatrix}^\top$ $\eta_{k+1} = \mathbf{A}\eta_k + \mathbf{B}a_k + \mathbf{G} \begin{bmatrix} w_k \\ v_k \end{bmatrix}$ $\tilde{y}_k = \mathbf{C}\eta_k + \mathbf{D}a_k + \mathbf{H} \begin{bmatrix} w_k \\ v_k \end{bmatrix}$
 - Anomaly Detector $\xi_{k} = \begin{bmatrix} \eta_{k}^{\top} & s_{k}^{\top} \end{bmatrix}^{\top}$ $\xi_{k+1} = \mathbf{A}_{e}\xi_{k} + \mathbf{B}_{e}a_{k} + \mathbf{G}_{e} \begin{bmatrix} w_{k} \\ v_{k} \end{bmatrix}$ $r_{k} = \mathbf{C}_{e}\xi_{k} + \mathbf{D}_{e}a_{k} + \mathbf{H}_{e} \begin{bmatrix} w_{k} \\ v_{k} \end{bmatrix}$ - Alarm triggered if: $\|r_{k}\| \ge \delta_{r} + \delta_{\alpha}, \quad \delta_{\alpha} \in \mathbb{R}^{+}$



1. Minimum Resource Attack: Dynamical Case

Dynamical anomaly detector for closed-loop system:

$$\xi_{k+1} = \mathbf{A}_{\mathbf{e}}\xi_k + \mathbf{B}_{\mathbf{e}}a_k + \mathbf{G}_{\mathbf{e}}w_k$$
$$r_k = \mathbf{C}_{\mathbf{e}}\xi_k + \mathbf{D}_{\mathbf{e}}a_k + \mathbf{H}_{\mathbf{e}}v_k$$

Lift to time interval [0, N] with zero-initial conditions and no noise:





1. Minimum Resource Attack: Dynamical Case



- Minimize disruption resources (#channels attacked)
- No alarms (threshold δ_{α})
- Reach attack goals \mathcal{G} (compare to security index)



1. Minimum Resource Attack: Formulate as MILP (1) $\min_{\mathbf{a}} ||h_p(\mathbf{a})||_0$

Note that

 $\|h_p(\mathbf{a})\|_0 \le \epsilon$

 $h_p(\mathbf{a}) = [\|\mathbf{a}_{(1)}\|_{\ell_p}, \dots, \|\mathbf{a}_{(i)}\|_{\ell_p}, \dots, \|\mathbf{a}_{(q_a)}\|_{\ell_p}]$ $\|\mathbf{r}\|_q = \|\mathcal{T}_r \mathbf{a}\|_q \le \delta_\alpha$ $\mathbf{a} \in \mathcal{G}$

can equivalently be formulated as

$$\begin{array}{rcl} \mathbf{a}_{(i)} &\leq & M_h \mathbf{z}_i \mathbf{1} & \forall i = 1, \dots, q_a \\ -\mathbf{a}_{(i)} &\leq & M_h \mathbf{z}_i \mathbf{1} & \forall i = 1, \dots, q_a \\ & \sum_{i=1}^{q_a} \mathbf{z}_i &\leq & \epsilon \\ & \mathbf{z}_i &\in & \{0,1\} & \forall i = 1, \dots, q_a. \end{array}$$
where M_h is a large constant ("infinity")



1. Minimum Resource Attack: Formulate as MILP (2)

 $\min_{\mathbf{a}, \epsilon} \epsilon$

such that

$$h_p(\mathbf{a}) = [\|\mathbf{a}_{(1)}\|_{\ell_p}, \dots, \|\mathbf{a}_{(i)}\|_{\ell_p}, \dots, \|\mathbf{a}_{(q_a)}\|_{\ell_p}]$$
$$\|h_p(\mathbf{a})\|_0 \le \epsilon$$
$$\|\mathbf{r}\|_q = \|\mathcal{T}_r \mathbf{a}\|_q \le \delta_\alpha$$
$$\mathbf{a} \in \mathcal{G}$$

- Minimize disruption resources (#channels attacked)
- No alarms (threshold δ_{α})
- Reach attack goals \mathcal{G} (compare to security index)
- MILP if $p=q=\infty$



OYAL INSTITUTI

2. Maximum Impact Attack: Dynamical Case

Dynamics of plant and controller:

$$\eta_{k+1} = \mathbf{A}\eta_k + \mathbf{B}a_k + \mathbf{G}w_k$$
$$x_k = \mathbf{C}\eta_k + \mathbf{D}a_k + \mathbf{H}v_k$$

Lifting to time interval [0, N] with zero-initial conditions and no noise:





2. Maximum Impact Attack: Dynamical Case

$$\begin{split} \max_{\mathbf{a}} \|\mathcal{T}_{x}\mathbf{a}\|_{p} \\ \text{such that} \\ \|\mathbf{r}\|_{q} = \|\mathcal{T}_{r}\mathbf{a}\|_{q} \leq \delta_{\alpha} \end{split}$$

- Maximize impact (push $\|\mathbf{x}\|_p$ far away from equilibrium)
- No alarms (threshold δ_{α})
- Not a convex optimization problem!
- Closed-form solution when p = q = 2 (use Courant-Fischer Theorem) See Exercise 6



2. Maximum Impact Attack: Dynamical Case

$$\begin{split} \max_{\mathbf{a}} \|\mathcal{T}_{x}\mathbf{a}\|_{p} \\ \text{such that} \\ \|\mathbf{r}\|_{q} = \|\mathcal{T}_{r}\mathbf{a}\|_{q} \leq \delta_{\alpha} \end{split}$$

Theorem: Bounded solution iff $ker(\mathcal{T}_r) \subseteq ker(\mathcal{T}_x)$ See Exercise 7

Theorem (p\psi\psi\psi\psi\psi\psi\psi\psi\psi\psi\psi\psi): Assume bounded solution, then $\mathbf{a}^{\star} = \frac{\delta_{\alpha}}{\|\mathcal{T}_{r}\mathbf{v}_{\max}\|_{2}}\mathbf{v}_{\max}, \quad \|\mathcal{T}_{\mathbf{x}}\mathbf{a}^{\star}\|_{2} = \sqrt{\lambda_{\max}}\delta_{\alpha}$ $0 = (\lambda_{\max}\mathcal{T}_{r}^{\top}\mathcal{T}_{r} - \mathcal{T}_{x}^{\top}\mathcal{T}_{x})\mathbf{v}_{\max} \quad (\lambda_{\max}/\mathbf{v}_{\max} \text{ max generalized eigenpair})$

What happens for infinite time-horizons?



2. Maximum-Impact Attack Infinite Horizon



- Maximum-impact stealthy attack:
 - Maximize "energy" of the state signal
 - Keep the residual signal "small"



- If the system has unstable zerodynamics:
 - There exists an *exponentially increasing* input that attains a "small" output

$$\{a_k\}_{k=0}^{\infty}: r_k \approx 0, \ \forall k$$

 $\|a_k\| \to \infty, \quad \|x_k\| \to \infty$

See Exercise 7



3. Maximum Impact Bounded Resource Attack

 $\max_{\mathbf{a}} \|\mathcal{T}_{x}\mathbf{a}\|_{p}$ such that $\|\mathbf{r}\|_{q} = \|\mathcal{T}_{r}\mathbf{a}\|_{q} \le \delta_{\alpha}$ $\|h_{p}(\mathbf{a})\|_{0} \le \epsilon$

- Maximize impact (push $\|\mathbf{x}\|_p$ far away from equilibrium)
- No alarms (threshold δ_{α})
- Use no more than ϵ channels

Jointly considers impact and likelihood

• $p = q = \infty$ can be formulated as MILP (see slide 15)



Numerical Example

ROYAL INSTITUTE OF TECHNOLOGY



$$\begin{split} \dot{h_4} & \dot{h_1} = -\frac{a_1}{A_1}\sqrt{2gh_1} + \frac{a_3}{A_1}\sqrt{2gh_3} + \frac{\gamma_1k_1}{A_1}u_1, \\ \dot{h_2} = -\frac{a_2}{A_2}\sqrt{2gh_2} + \frac{a_4}{A_2}\sqrt{2gh_4} + \frac{\gamma_2k_2}{A_2}u_2, \\ \dot{h_3} = -\frac{a_3}{A_3}\sqrt{2gh_3} + \frac{(1-\gamma_2)k_2}{A_3}u_2, \\ \dot{h_4} = -\frac{a_4}{A_4}\sqrt{2gh_4} + \frac{(1-\gamma_1)k_1}{A_4}u_1, \end{split}$$

- Wireless LQG controller
- 4 channels: 2 actuators and 2 measurements
- Minimum phase or non-minimum phase depending on $\gamma_1,\,\gamma_2$



OF TECHNOLOGY

Numerical Example (Non-Min Phase)

Values of $\|\mathbf{x}\|_p$ for maximum impact formulation with $p = q = 2, \ \delta_{\alpha} = 0.15$

	$\Pi = F \times -7 \Pi^{-1}$			
	1	2	3	4
Minimum phase	1.15	140.39	∞	∞
Non-minimum phase	2.80	689.43	∞	∞





Numerical Example (Non-Min Phase)





Numerical Example (MILP)

ROYAL INSTITUTE OF TECHNOLOGY





OF TECHNOLOGY

What components to protect?

- Attacks on 3 or more components have very high impact
 - Must protect 2 components
- Attacks on 2 compoents still have **high risk**
- What to protect?
 - Protecting ({u₁, u₂}) yields the lowest risk





Numerical Example

- Maximum Impact Bounded Resource attack illustrated
- 2 channels allowed: MILP selects the actuators
- 3-4 channels allowed: Unbounded impact (any attack on actuators can be hidden by corrupting 2 measurements)
- Infinity norm criteria ($p=q=\infty$) yields more aggressive attack (bounds saturated)
- Not surprisingly, non-min phase plant more sensitive



- Tools for quantitative trade-off analysis between attacker's impact and resources: Important for defense prioritization
- For dynamical systems there are *temporal* as well as *spatial* (*channel*) *constraints* for attacker to fulfill
 - Enforced through lifting models
- Closed-form solutions and mixed integer linear programming formulations



Lecture 3

- ROYAL INSTITUTE OF TECHNOLOGY
 - Risk management
 - security metrics, quantifying impact [1]

Anomaly detectors (for detectable attacks) [2]

- Watermarking (against undetectable attacks)
 - Replay attacks [3]
 - Zero dynamics attacks [4]



Model-Based Fault Diagnosis

ROYAL INSTITUTE OF TECHNOLOGY



• Basic ideas:

- Compute an expected output (using model information)
- Evaluate the difference between the real and expected outputs (residual)



F TECHNOLOGY

Fault-Diagnosis Objectives

- Fault Detection: detect faults
 - Generate a residual sensitive to faults
 - E.g.: Kalman filter $\hat{x}_{k+1} = A\hat{x}_k + Bu_k + K(y_k C\hat{x}_k)$ $r_k = y_k - C\hat{x}_k$



- Fault Isolation: locate the faulty components
 - Generate a set of structured residuals

- E.g.:
$$\begin{bmatrix} r_1(s) \\ r_2(s) \end{bmatrix} = \begin{bmatrix} G_{11}(s) & 0 \\ 0 & G_{22}(s) \end{bmatrix} \begin{bmatrix} f_1(s) \\ f_2(s) \end{bmatrix}$$

- Fault Identification: estimate the fault signal
 - Use a state estimator
 - Simple example: $\hat{x}_{k+1} A\hat{x}_k Bu_k = F\hat{f}_k$



• All fail in the presence of undetectable attacks



Beyond Fault Diagnosis

ROYAL INSTITUTE OF TECHNOLOGY



Often requires redundancy (extra sensors / actuators)
See [2] for more details and references



Lecture 3

- ROYAL INSTITUTE OF TECHNOLOGY
 - Risk management
 - security metrics, quantifying impact [1]

Anomaly detectors (for detectable attacks) [2]

- Watermarking (against undetectable attacks)
 - Replay attacks [3]
 - Zero dynamics attacks [4]



Replay Attack – Phase II [3]

ROYAL INSTITUTE OF TECHNOLOGY

- No more data is recorded $\mathcal{I}_k := \mathcal{I}_{k-1}, \quad k > r$
- The previously recorded data is replayed

$$\begin{bmatrix} b_k^u \\ b_k^y \end{bmatrix} = \begin{bmatrix} \Upsilon^u(u_{k-T} - u_k) \\ \Upsilon^y(y_{k-T} - y_k) \end{bmatrix}$$

 Physical attack is also performed

 $f_k = g_f(\emptyset, \mathcal{I}_{k_r})$

No system knowledge is needed:



$$\mathcal{K} = \{\hat{P}, \, \hat{F}, \, \hat{D}\} = \emptyset$$

[Mo and Sinopoli, Allerton, 2009]



Replay Attack - Experiment

ROYAL INSTITUTE OF TECHNOLOGY





- Attack Goal: Empty tank 4
- Replay attack on sensor 2
- Physical attack on tank 4
- ³⁵⁰ Tank 4 is **emptied**
 - Physical attack ends at t=180s
 - Replay attack ends at t=280s
 - The attack is **not detected**
 - Why is it undetectable?
 - Can we "make" it detectable?



Sensor Replay Attack – Analysis [3]

ROYAL INSTITUTE OF TECHNOLOGY

• Residual generated by a Kalman Filter (+LQG controller $u_k = L\hat{x}_k$) $\hat{x}_{k+1|k} = A\hat{x}_k + Bu_k$ (no fault) $\hat{x}_k = \hat{x}_{k|k-1} + K(y_k - C\hat{x}_{k|k-1})$ (no fault) $r_k = y_k - C\hat{x}_k$ $x_k - \hat{x}_k \sim \mathcal{N}(0, P_k)$ $r_k \sim \mathcal{N}(0, CP_kC^\top + R)$

 Model data replay as a virtual time-shifted plant

$$x_k^v = x_{k-T}, \quad y_k^v = y_{k-T}, \quad r_k^v = r_{k-T}$$

• Residual under sensor replay attack $r_{k+1} = r_{k+1}^v - C\mathcal{A}(\hat{x}_{0|-1} - \hat{x}_{0|-1}^v)$ $\mathcal{A} = (A + BL)(I - KC)$

- Attack is stealthy if $\ensuremath{\mathcal{A}}$ is stable
 - Attacked residual converges to healthy residual $r_k \to r_k^v \sim \mathcal{N}(0, CP_k C^\top + R)$
 - Relies on the (virtual) plant and Kalman filter having the same control policy





Sensor Replay Attack – Watermarking [3]

- The plant **proactively changes** the control policy
 - Adds noise to control input: $u_k = L\hat{x}_k + \zeta_k$
 - Noise is randomly generated, but known
- Residual under sensor replay attack

$$r_{k+1} = r_{k+1}^v - C\mathcal{A}(\hat{x}_{0|-1} - \hat{x}_{0|-1}^v) - C\sum_{i=0}^n \mathcal{A}^{k-i}B(\zeta_i - \zeta_i^v)$$
$$\mathcal{A} = (A + BL)(I - KC)$$

- Distribution of residual under attack changes
 - Nominal distribution: $r_k \sim \mathcal{N}(0, CP_kC^\top + R)$
 - Distribution with attack: $r_k \sim \mathcal{N}(\mu_{k-1}, CP_kC^\top + R + \Sigma)$
- Attack can be detected by comparing the two distributions
 - See Kullback-Leibler divergence and Neyman-Pearson test
 - Detection enabled by **asymmetries** between the time-shifted and real-time models





Lecture 3

- ROYAL INSTITUTE OF TECHNOLOGY
 - Risk management
 - security metrics, quantifying impact [1]

Anomaly detectors (for detectable attacks) [2]

- Watermarking (against undetectable attacks)
 - Replay attacks [3]
 - Zero dynamics attacks [4]



Zero-Dynamics Attack Model

ROYAL INSTITUTE OF TECHNOLOGY



• Physical Plant under attack

$$\mathcal{P}^a: \begin{cases} x_{k+1} = Ax_k + Ba_k \\ y_k = Cx_k \end{cases}, \quad x_0 = 0$$

- Attack policy
 - Computed using A, B, C
 - Open-loop policy
- Attack Goals and Constraints
 - Reach an unsafe state
 - Remain stealthy



Testbed for Networked Control System Security

ROYAL INSTITUTE OF TECHNOLOGY



Quadruple-tank process has an unstable zero if $0 < \gamma_1 + \gamma_2 < 1$

[Johansson, IEEE TCST, 2000]



Experimental Result (Lecture 2)





- Attack Goal: Empty tank 3
- Zero-dynamics attack on both actuators
- Tank 3 becomes empty
- The attack is **detected**



Experimental Result – Why?



Smooth increase

- What causes it?
- Does it compromise the attack's stealthiness?
- Abrupt increase
 - How can it be induced so that attacks are revealed?





OF TECHNOLOGY

Revealing Zero Dynamics Attacks [4]



Revisit zero dynamics

 Output behavior with initial condition mismatch

Revealing zero dynamics attacks



Zero Dynamics

OF TECHNOLOGY

• Physical plant $\mathcal{P}: \begin{cases} x_{k+1} = Ax_k + Ba_k \\ y_k = Cx_k \end{cases}$

- Output-zeroing problem:
 - Find z_0 and F such that $x_{k+1} = (A + BF)x_k$, $x_0 = z_0$ $0 = Cx_k$

 $\begin{array}{l} \textbf{Zero Dynamics} \\ x_{k+1} = (A + BF)x_k \\ 0 = Cx_k, \end{array}$ with $x_0 \in \mathcal{V}^* \subset \ker(\mathbf{C})$ and F such that $(A + BF)\mathcal{V}^* \subseteq \mathcal{V}^*.$



F TECHNOLOGY

Zero Dynamics Attack



• Physical plant under attack $\mathcal{P}: \begin{cases} x_{k+1} = Ax_k + Ba_k \\ y_k = Cx_k \end{cases}, \quad x_0 = 0$

Undetectable attacks

 $\begin{aligned} x_{k+1} &= Ax_k + Ba_k\\ 0 &= Cx_k \end{aligned}, \quad x_0 = 0 \end{aligned}$

Zero dynamics attack policy

What happens when $x_0 \neq z_0$?



Initial Condition Mismatch

ROYAL INSTITUTE OF TECHNOLOGY

> **Theorem 1.** The output produced by the zero dynamics attack $z_{k+1} = (A + BF)z_k$, $z_0 \in \mathcal{V}^*$ $a_k = Fz_k$, $z_0 \in \mathcal{V}^*$ is described by $e_{k+1} = Ae_k$ $y_k = Ce_k$, $e_0 = -z_0$

- Attack is not undetectable if z_0 belongs to the observable subspace of (A, C) (i.e. z_0 yields a zero output)
- If *A* is stable:
 - the resulting output energy is finite;
 - the output can be made arbitrarily small by scaling down the initial condition z_0



Revealing Zero Dynamics Attacks

Definition: A zero dynamics attack is revealed if $y_k \neq 0$

Proposed approach (watermarking):

- change the system dynamics from $\Sigma = (A, B, C)$ to $\tilde{\Sigma} = (\tilde{A}, \tilde{B}, \tilde{C})$

Zero-dynamics attacks are stealthy with respect to the system

$$\begin{bmatrix} x_{k+1} \\ z_{k+1} \end{bmatrix} = \begin{bmatrix} A & BF \\ 0 & A + BF \end{bmatrix} \begin{bmatrix} x_k \\ z_k \end{bmatrix}$$
$$y_k = \begin{bmatrix} C & 0 \end{bmatrix} \begin{bmatrix} x_k \\ z_k \end{bmatrix}$$
for all $x_0 = z_0 \in \mathcal{V}^*$

Every zero-dynamics attack is revealed if the system

$$\begin{bmatrix} x_{k+1} \\ z_{k+1} \end{bmatrix} = \begin{bmatrix} \tilde{A} & \tilde{B}F \\ 0 & A+BF \end{bmatrix} \begin{bmatrix} x_k \\ z_k \end{bmatrix}$$
$$y_k = \begin{bmatrix} \tilde{C} & 0 \end{bmatrix} \begin{bmatrix} x_k \\ z_k \end{bmatrix}$$
is observable for all $x_0 = z_0 \in \mathcal{V}^*$



Modifying the input matrix B (1)

- Consider $\tilde{B} = BW$
- Observation: attacks remain undetectable w.r.t to W if and only if the unobservable trajectories are not perturbed

Proof sketch: Check the conditions when $x_0 = z_0 \in \mathcal{V}^*$ is unobservable:

$$\begin{bmatrix} \lambda I - A & -BF - B(W - I)F \\ 0 & \lambda I - (A + BF) \\ C & 0 \end{bmatrix} \begin{bmatrix} z_0 \\ z_0 \end{bmatrix} = 0 \quad \text{, where } (\lambda I - (A + BF))z_0 = 0$$

 $B(W-I)Fz_0 = 0$

• **Revealing attacks:** Choose W such that $B(W - I)Fz_0 \neq 0$

Theorem. All the zero dynamics attacks associated with a given $z_0 \in \mathcal{V}^*$ remain stealthy with respect to $\tilde{\Sigma} = (A, \tilde{B}, C)$ if and only if $\mathcal{V}^* \subseteq \ker(B(W - I)F)$



Modifying the input matrix B (2)

ROYAL INSTITUTE OF TECHNOLOGY

- Simply changing B to B affects the system performance under no attack
- Coordinated input scaling:
 - Similar to encryption



Theorem. Let $z_k = x_k = z \in \mathcal{V}^*$ and apply $W = \alpha I$ at time k. The output trajectory is described by

$$e_{k+1} = Ae_k$$

$$y_k = Ce_k, \ e_0 = (1 - \alpha)z$$

- **Revealing attack:** choose α such that y_k is "large" enough.
 - Does not affect the system dynamics



Example – modifying B

ROYAL INSTITUTE OF TECHNOLOGY

• Solution to reveal attacks: input scaling $W = \alpha I$



• Example: choose $\alpha = 0.987$

- Attack begins at k=0
 - Initial condition mismatch
- Input scaling applied at k = 100
 - the attack is revealed
- Stable A results in finite output energy





Modifying the output matrix C

- Consider $\tilde{C} = \begin{bmatrix} C \\ \Delta C \end{bmatrix}$
- **Observation:** attacks remain *undetectable* w.r.t to \tilde{C} if and only if $\Sigma = (A, B, C)$ and $\tilde{\Sigma} = (A, B, \tilde{C})$ share common unobservable trajectories
- Revealing attacks: add measurements so that $\mathcal{X} = \tilde{\mathcal{V}}^{\star} \cap \mathcal{V}^{\star}$ becomes empty
 - system dynamics are not affected
 - Requires at most $\,\dim(\mathcal{V}^\star)$ new measurements

Theorem. There exists a $z_0 \in \mathcal{V}^*$ generating an undetectable attack to $\tilde{\Sigma} = (A, B, \tilde{C})$ if and only if (A + BF) has an eigenvector in $\mathcal{V}^* \cap \ker \tilde{C}$



Modifying the system matrix A

- Consider $\tilde{A} = A + \Delta A$
- **Observation:** attacks remain *undetectable* w.r.t to \tilde{A} if and only if the unobservable trajectories are *not perturbed* (similar to \tilde{B})

Proof sketch: Check the conditions for which $x_0 = z_0 \in \mathcal{V}^*$ is unobservable:

$$\begin{bmatrix} \lambda I - A - \Delta A & -BF \\ 0 & \lambda I - (A + BF) \\ C & 0 \end{bmatrix} \begin{bmatrix} z_0 \\ z_0 \end{bmatrix} = 0 \quad \text{, where } (\lambda I - (A + BF)) z_0 = 0$$

- **Revealing attacks:** choose ΔA such that $\mathcal{V}^* \cap \ker(\Delta A) = \emptyset$
 - Affects the system dynamics and may also require re-designing the controller

Theorem. All the zero-dynamics attacks associated with a given $z_0 \in \mathcal{V}^*$ remain stealthy with respect to $\tilde{\Sigma} = (\tilde{A}, B, C)$ if and only if $\mathcal{V}^* \subseteq \ker(\Delta A)$



Example – modifying A

ROYAL INSTITUTI OF TECHNOLOGY

• Solution to reveal attacks: ΔA such that $\mathcal{V}^* \cap \ker(\Delta A) = \emptyset$

 $V^{\star} = \begin{vmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \\ -1 & \mathbf{0} \\ \mathbf{0} & 1 \end{vmatrix} \qquad \Delta A = \begin{bmatrix} \mathbf{0} & \Delta \end{bmatrix}$





Tank 2 y₂

Tank 1 y1



Summary



- Zero-dynamics attacks are robust to initial condition mismatch
- Proposed methods to reveal attacks by
 - Changing C: Adding measurements
 - Changing A: Modifying the open-loop dynamics
 - Changing B: Cooperatively scaling the input signals
- Adding measurements and scaling input signals does not affect the system performance



Summary of Lecture 3

- Risk management
 - security metrics, quantifying impact [1]
 - Tools for quantitative trade-off analysis between attacker's impact and resources: Important for defense prioritization
- Basics of Fault Diagnosis (for detectable attacks) [2]
- Watermarking (against undetectable attacks)
 - Induce asymmetries in the attacker/system models
 - Additive signals against sensor replay attacks [3]
 - Model perturbations against zero dynamics attacks [4]