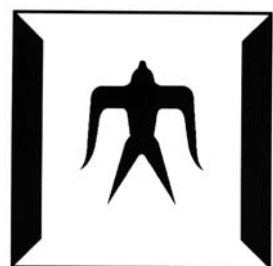


Pool-based Agnostic Experiment Design in Linear Regression



Masashi Sugiyama (Tokyo Tech.)
Shinichi Nakajima (Nikon)

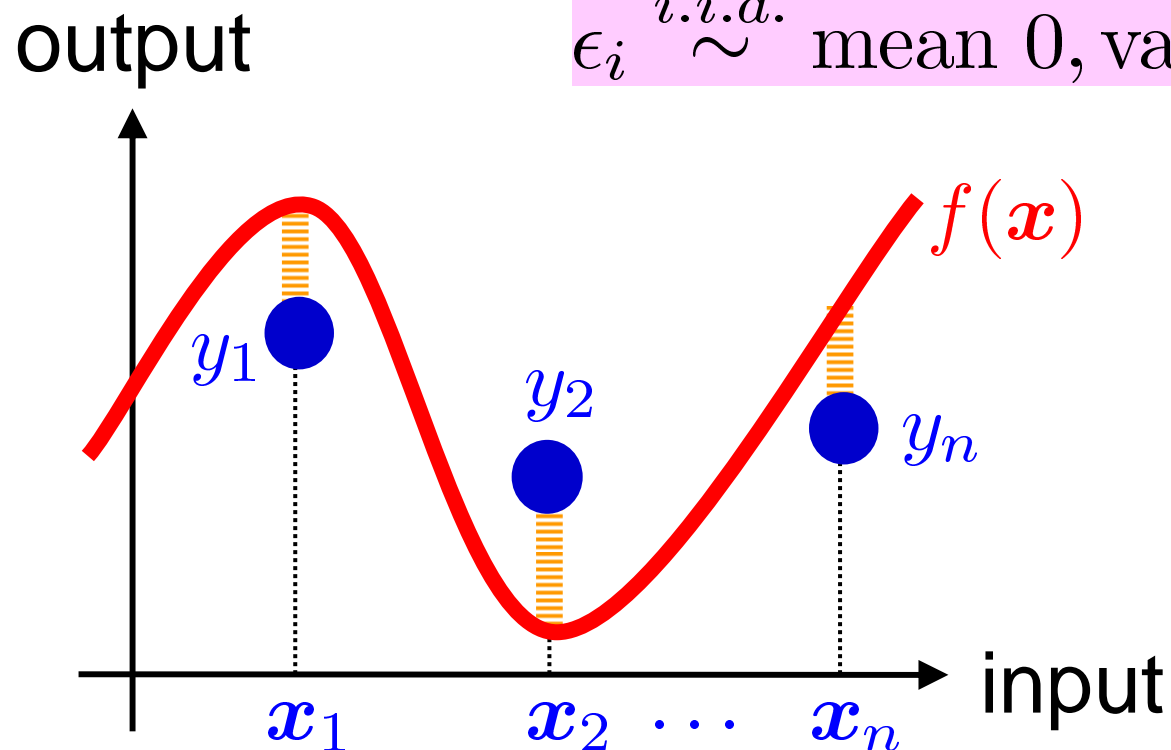
Linear Regression

2

- Learn a **real-valued** function $f(\mathbf{x})$ from input-output training samples $\{\mathbf{x}_i, y_i\}_{i=1}^n$.

$$y_i = f(\mathbf{x}_i) + \epsilon_i$$

$$\epsilon_i \stackrel{i.i.d.}{\sim} \text{mean } 0, \text{ variance } \sigma^2$$



Linear Regression (cont.)

3

- **Linear model** is used for learning:

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^b \alpha_i \varphi_i(\mathbf{x})$$

α_i : Parameter
 $\varphi_i(\mathbf{x})$: Basis function

- **Goal**: learn $\hat{\alpha}$ so that the **generalization error** is minimized

$$Gen = \mathbb{E}_{\epsilon} \int \left(f(\mathbf{x}) - \hat{f}(\mathbf{x}) \right)^2 p_{test}(\mathbf{x}) d\mathbf{x}$$

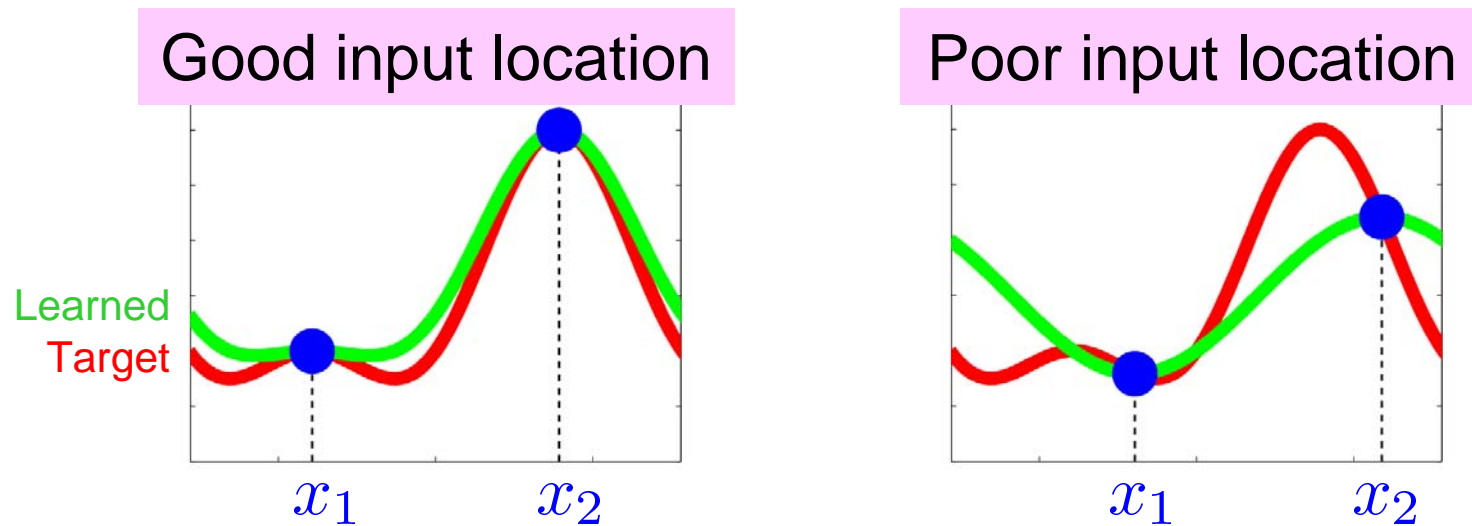
$$= \mathbb{E}_{\epsilon} \|f - \hat{f}\|_{p_{test}}^2$$

\mathbb{E}_{ϵ} : Expectation over noise

$p_{test}(\mathbf{x})$: Test input density

Experiment Design

- Quality of learned functions depends on **training input location** $\{\mathbf{x}_i\}_{i=1}^n$.



- Goal:** optimize training input location

$$\min_{\{\mathbf{x}_i\}_{i=1}^n} Gen$$

$$Gen = \mathbb{E}_{\epsilon} \|f - \hat{f}\|_{p_{test}}^2$$

Challenges

5

$$\min_{\{\mathbf{x}_i\}_{i=1}^n} Gen$$

$$Gen = \mathbb{E}_{\epsilon} \|f - \hat{f}\|_{p_{test}}^2$$

- Gen is unknown and needs to be estimated.
- In experiment design, we do not have **training output values** $\{y_i\}_{i=1}^n$ yet.
- Thus we cannot use, e.g., **cross-validation** which requires $\{\mathbf{x}_i, y_i\}_{i=1}^n$.
- Only training input positions $\{\mathbf{x}_i\}_{i=1}^n$ can be used in generalization error estimation!

Organization

6

Pool-based Agnostic Experiment Design in Linear Regression

1. Problem definition
2. **Basic strategy**
3. Proposed method
4. Experiments

Bias and Variance

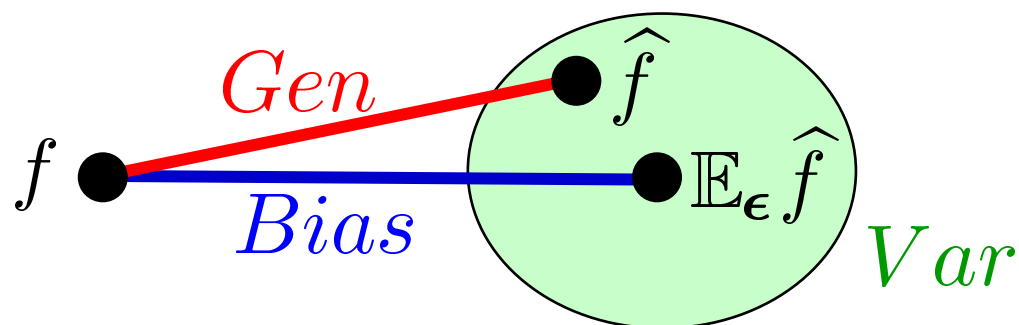
7

$$\underbrace{\mathbb{E}_{\epsilon} \|f - \hat{f}\|_{p_{test}}^2}_{Gen} = \underbrace{\|f - \mathbb{E}_{\epsilon} \hat{f}\|_{p_{test}}^2}_{Bias} + \underbrace{\mathbb{E}_{\epsilon} \|\mathbb{E}_{\epsilon} \hat{f} - \hat{f}\|_{p_{test}}^2}_{Var}$$

- *Bias* is **not estimable** without $\{y_i\}_{i=1}^n$.
- For linear learning $\hat{\alpha} = Ly$:

$$Var = \sigma^2 \text{tr}(ULL^{\top})$$

- Noise variance σ^2 is **not estimable** without $\{y_i\}_{i=1}^n$.
- $\text{tr}(ULL^{\top})$ is **computable** from $\{\mathbf{x}_i\}_{i=1}^n$.



$$U_{i,j} = \langle \varphi_i, \varphi_j \rangle_{p_{test}}$$

L : Learning matrix

$$\mathbf{y} = (y_1, \dots, y_n)^{\top}$$

Key Trick in Experiment Design 8

- Find a setup where $Bias = 0$ is guaranteed.
- Then

$$Gen = \underbrace{Bias}_{0} + \underbrace{Var}_{\sigma^2 \text{tr}(ULL^\top)} \propto \text{tr}(ULL^\top)$$

- Thus

$$\text{argmin } Gen = \text{argmin } \underbrace{\text{tr}(ULL^\top)}_{\substack{\text{computable before} \\ \text{observing } \{y_i\}_{i=1}^n}}$$

$$U_{i,j} = \langle \varphi_i, \varphi_j \rangle_{p_{test}}$$

L : Learning matrix

Traditional Method

9

(Fedorov 1972)

- Assume **model is correct**:

$$\exists \alpha^*, \hat{f}(\mathbf{x}; \alpha^*) = f(\mathbf{x})$$

$$\hat{f}(\mathbf{x}; \alpha) = \sum_{i=1}^b \alpha_i \varphi_i(\mathbf{x})$$

- Use **ordinary least squares (OLS)** estimation:

$$\min_{\alpha} \left[\sum_{i=1}^n \left(\hat{f}(\mathbf{x}_i) - y_i \right)^2 \right]$$

$$\hat{\alpha}_O = L_O \mathbf{y}$$

- Experiment design criterion:

$$\min_{\{\mathbf{x}_i\}_{i=1}^n} \text{tr}(\mathbf{U} \mathbf{L}_O \mathbf{L}_O^\top)$$

$$\mathbf{L}_O = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

$$\mathbf{X}_{i,j} = \varphi_j(\mathbf{x}_i)$$

$$\mathbf{y} = (y_1, \dots, y_n)^\top$$

$$\mathbf{U}_{i,j} = \langle \varphi_i, \varphi_j \rangle_{p_{test}}$$

Goal of This Work

10

- **Pros / cons of traditional method:**
 - + Generalization error estimation is exact.
 - + Easy to implement.
 - Correct-model assumption is not realistic.
 - Very poor performance when agnostic.
 - Test input density $p_{test}(\mathbf{x})$ is often unknown.
- We propose a new method that is
 - Still easy to implement,
 - Robust against agnosticity,
 - Able to work without $p_{test}(\mathbf{x})$.

Organization

11

Pool-based Agnostic Experiment Design in Linear Regression

1. Problem definition
2. Basic strategy
3. Proposed method
 1. Overcoming agnosticity
 2. Coping with pool-based setup
4. Experiments

Weak Agnostic Setup

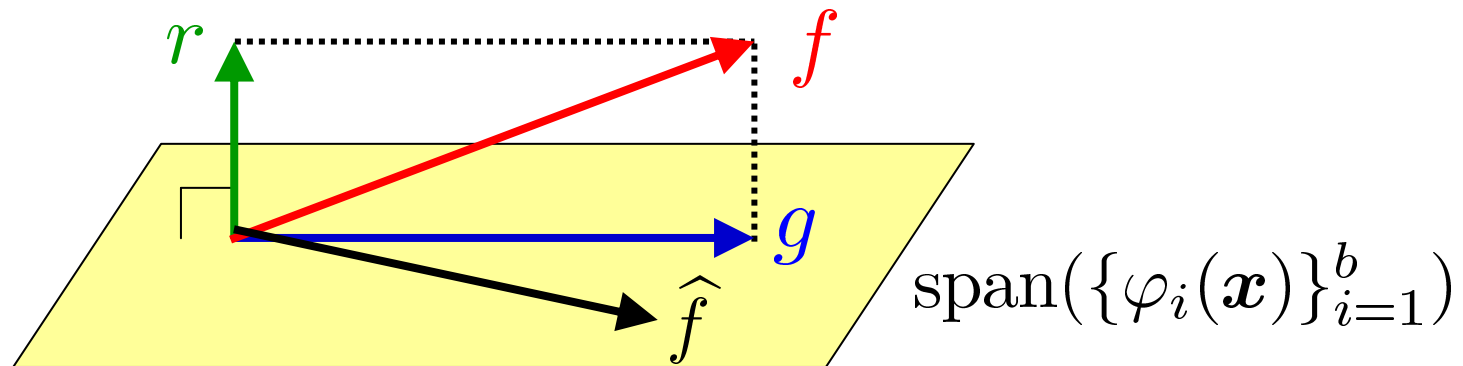
- The model is **not exactly correct**, but is **reasonably good**:

$$\exists \alpha^*, \|\hat{f}(\mathbf{x}; \alpha^*) - f(\mathbf{x})\| \approx 0$$

$$\hat{f}(\mathbf{x}; \alpha) = \sum_{i=1}^b \alpha_i \varphi_i(\mathbf{x})$$

- Decomposition of target function:

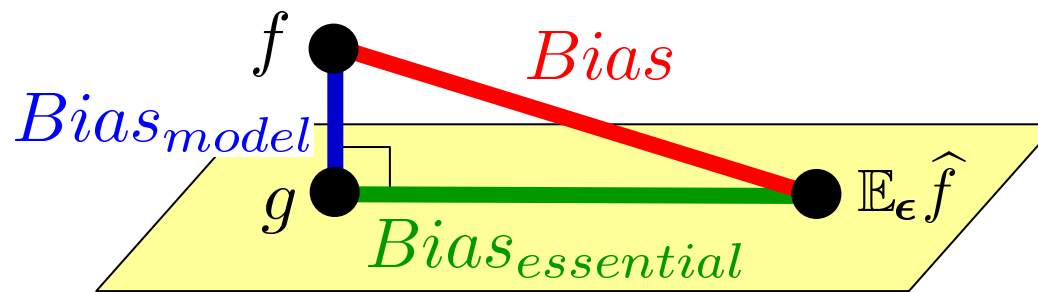
$$\underbrace{f(\mathbf{x})}_{\text{Target function}} = \underbrace{g(\mathbf{x})}_{\text{Approximable part}} + \underbrace{r(\mathbf{x})}_{\text{Residual part}}$$



Further Decomposition of Bias

13

$$\underbrace{\|f - \mathbb{E}_\epsilon \hat{f}\|_{p_{test}}^2}_{Bias} = \underbrace{\|f - g\|_{p_{test}}^2}_{Bias_{model}} + \underbrace{\|g - \mathbb{E}_\epsilon \hat{f}\|_{p_{test}}^2}_{Bias_{essential}}$$



- $Bias_{model}$ is constant and ignorable.
- But OLS cannot make $Bias_{essential}$ zero due to “covariate shift”: (Shimodaira JSPI2000)
 - Training / test inputs follow different distributions.
 - “Covariate” is another name for “input”.

Importance-Weighted LS (IWLS)¹⁴

$$\min_{\alpha} \left[\sum_{i=1}^n \underbrace{\frac{p_{test}(\mathbf{x}_i)}{p_{train}(\mathbf{x}_i)}}_{\text{Importance}} \left(\hat{f}(\mathbf{x}_i) - y_i \right)^2 \right] \quad \{\mathbf{x}_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} p_{train}(\mathbf{x})$$

- Even when agnostic: $\lim_{n \rightarrow \infty} Bias_{essential} = 0$
- When **weak agnostic**: $Bias_{essential} \ll Var$
- Solution is given by

$$\hat{\alpha}_W = L_W \mathbf{y}$$

$$L_W = (\mathbf{X}^\top \mathbf{D} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{D}$$

$$\mathbf{X}_{i,j} = \varphi_j(\mathbf{x}_i) \quad \mathbf{y} = (y_1, \dots, y_n)^\top$$

$$\mathbf{D} = \text{diag} \left(\frac{p_{test}(\mathbf{x}_1)}{p_{train}(\mathbf{x}_1)}, \dots, \frac{p_{test}(\mathbf{x}_n)}{p_{train}(\mathbf{x}_n)} \right)$$

Justification

15

(Sugiyama JMLR2006)

■ For IWLS

$$Gen = \underbrace{Bias_{model}}_{\text{constant}} + \underbrace{Bias_{essential}}_{\ll Var} + \underbrace{Var}_{\sigma^2 \text{tr}(\mathbf{U} \mathbf{L}_W \mathbf{L}_W^\top)}$$

■ Thus

$$\underset{p_{train}}{\text{argmin}} Gen \approx \underset{p_{train}}{\text{argmin}} \text{tr}(\mathbf{U} \mathbf{L}_W \mathbf{L}_W^\top)$$

computable before
observing $\{y_i\}_{i=1}^n$

$$\mathbf{L}_W = (\mathbf{X}^\top \mathbf{D} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{D}$$

$$\mathbf{X}_{i,j} = \varphi_j(\mathbf{x}_i) \quad \mathbf{U}_{i,j} = \langle \varphi_i, \varphi_j \rangle_{p_{test}}$$

$$\mathbf{D} = \text{diag} \left(\frac{p_{test}(\mathbf{x}_1)}{p_{train}(\mathbf{x}_1)}, \dots, \frac{p_{test}(\mathbf{x}_n)}{p_{train}(\mathbf{x}_n)} \right)$$

Pool-based Agnostic Experiment Design in Linear Regression

1. Problem definition
2. Basic strategy
3. Proposed method
 1. Overcoming agnosticity
 2. Coping with pool-based setup
4. Experiments

Pool-based Setup

■ Pool-based setup:

- The test input density $p_{test}(\mathbf{x})$ is **unknown**.

 **Importance weight** is not accessible.

$$D = \text{diag} \left(\frac{p_{test}(\mathbf{x}_1)}{p_{train}(\mathbf{x}_1)}, \dots, \frac{p_{test}(\mathbf{x}_n)}{p_{train}(\mathbf{x}_n)} \right)$$

- But **a pool of test input samples** is given.

$$\{\mathbf{x}'_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} p_{test}(\mathbf{x})$$

- Training input points are chosen **from the pool**.

$$\{\mathbf{x}_i\}_{i=1}^n \subset \{\mathbf{x}'_i\}_{i=1}^N$$

- We assume $N \gg n$.

Computing Importance Weight 18

- $\{b(\mathbf{x}'_i)\}_{i=1}^N$: **Resampling probability** of $\{\mathbf{x}'_i\}_{i=1}^N$

$$\sum_{i=1}^N b(\mathbf{x}'_i) = 1, \quad b(\mathbf{x}'_i) \geq 0$$

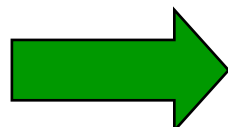
- Choose $\{\mathbf{x}_i\}_{i=1}^n$ following $\{b(\mathbf{x}'_i)\}_{i=1}^N$.

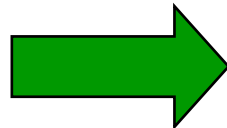
$$\{\mathbf{x}_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} \{b(\mathbf{x}'_i)\}_{i=1}^N$$

- Then D can be **exactly computed**:

$$p_{train}(\mathbf{x}_i) = p_{test}(\mathbf{x}_i)b(\mathbf{x}_i)$$

$$\{\mathbf{x}'_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} p_{test}(\mathbf{x})$$


$$\frac{p_{test}(\mathbf{x}_i)}{p_{train}(\mathbf{x}_i)} = \frac{1}{b(\mathbf{x}_i)}$$


$$D = \text{diag} \left(\frac{1}{b(\mathbf{x}_1)}, \dots, \frac{1}{b(\mathbf{x}_n)} \right)$$

Proposed Method

- Choose resampling function based on

$$\min_b \text{tr}(\hat{U} L_W L_W^\top)$$

$$\{\mathbf{x}'_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} p_{\text{test}}(\mathbf{x})$$

$$L_W = (X^\top D X)^{-1} X^\top D$$

$$\{\mathbf{x}_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} \{b(\mathbf{x}'_i)\}_{i=1}^N$$

$$X_{i,j} = \varphi_j(\mathbf{x}_i)$$

$$\hat{U}_{i,j} = \frac{1}{N} \sum_{i=1}^N \varphi_i(\mathbf{x}'_i) \varphi_j(\mathbf{x}'_i)$$

$$D = \text{diag} \left(\frac{1}{b(\mathbf{x}_1)}, \dots, \frac{1}{b(\mathbf{x}_n)} \right)$$

- Advantages:

- Robust against model misspecification.
- Easy to implement.

Organization

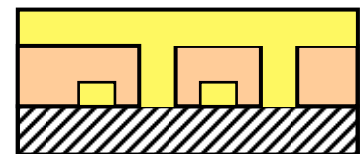
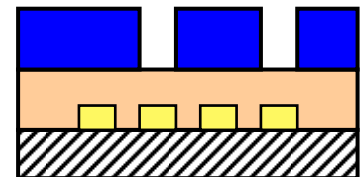
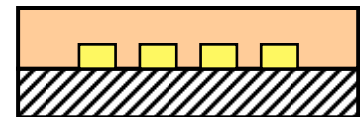
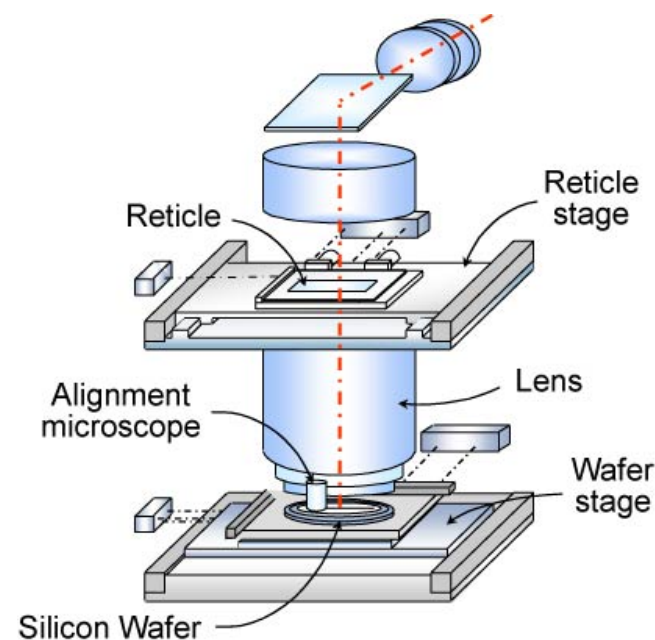
20

Pool-based Agnostic Experiment Design in Linear Regression

1. Problem definition
2. Basic strategy
3. Proposed method
4. Experiments

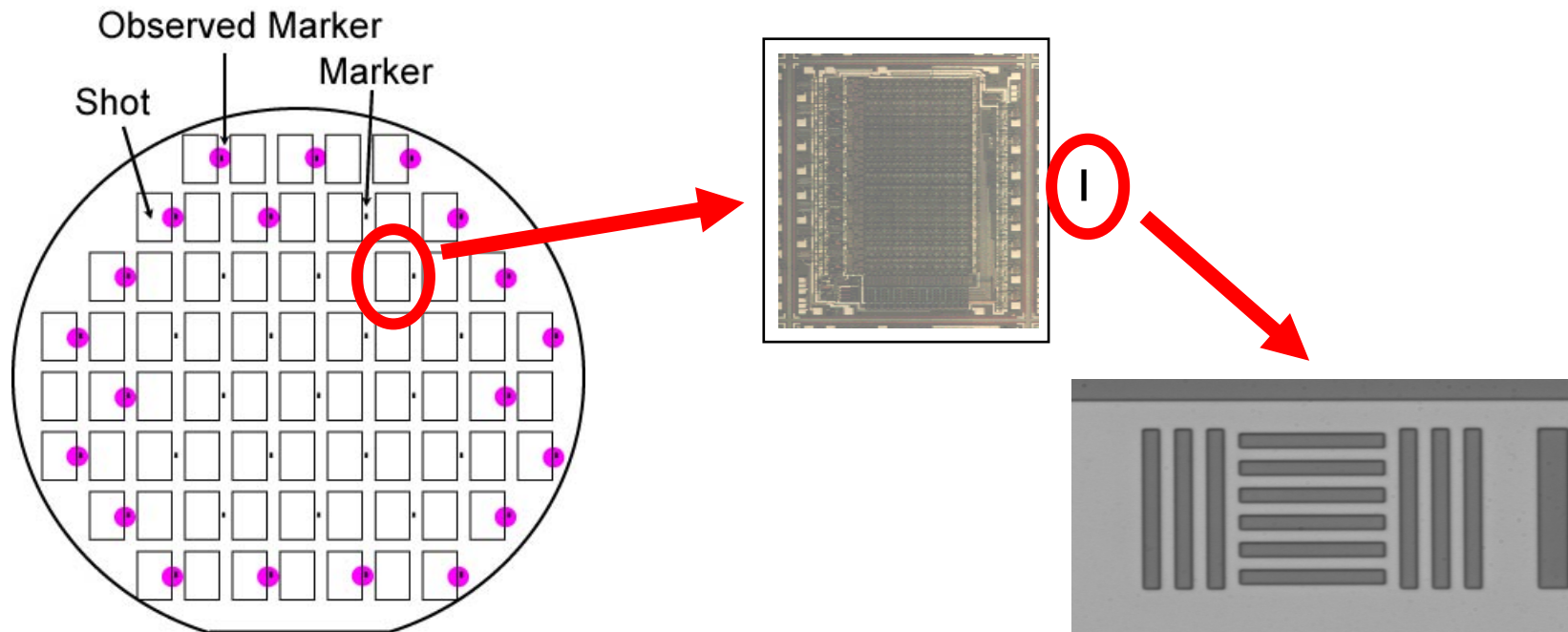
Wafer Alignment in Semiconductor Exposure Apparatus

- Recent silicon wafers have **layer structure**.
- Circuit patterns are exposed **multiple times**.
- **Exact alignment** of wafers is very important.



Markers on Wafer

- Wafer alignment process:
 - Measure marker location printed on wafers.
 - Shift and rotate the wafer to minimize the gap.
- For speeding up, **reducing the number of markers to measure** is very important.



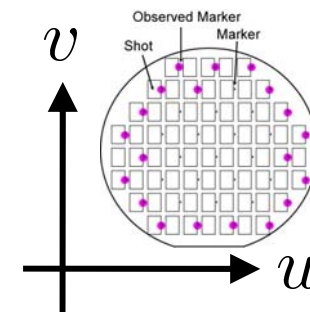
Non-linear Alignment Model

23

- When gap is only **shift and rotation**, linear model is exact:

$$\Delta u \text{ or } \Delta v = \theta_0 + \theta_1 u + \theta_2 v$$

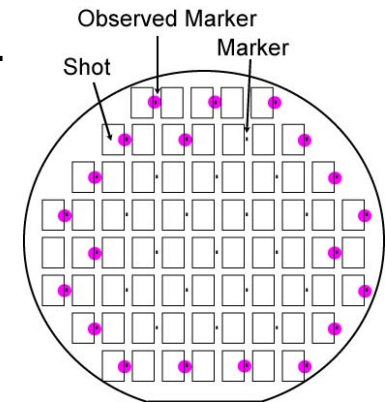
- However, **non-linear factors** exist, e.g.,
 - Warp
 - Biased characteristic of measurement apparatus
 - Different temperature conditions
- **Exactly modeling non-linear factors is very difficult in practice!**



Experimental Results

24

- 20 markers (out of 38) are chosen by experiment design methods.
- Gaps of all markers are predicted.
- Repeated for 220 different wafers.
- Mean (standard deviation) of the gap prediction error
- Red: Significantly better by 5% Wilcoxon test
- Blue: Worse than the baseline passive method



Model	Pool / Agnostic (Proposed)	Pool / Non-agnostic (Fedorov 1972)	“Outer” heuristic	Passive (Random)
Order 1	2.27(1.08)	2.37(1.15)	2.36(1.15)	2.32(1.11)
Order 2	1.93(0.89)	1.96(0.91)	2.13(1.08)	2.32(1.15)

Order 1: Δu or $\Delta v = \theta_0 + \theta_1 u + \theta_2 v$

Order 2: Δu or $\Delta v = \theta_0 + \theta_1 u + \theta_2 v + \theta_3 uv + \theta_4 u^2 + \theta_5 v^2$

■ Proposed method works the best!

Conclusions

- We proposed a **pool-based agnostic** experiment design method for linear regression.
- Proposed method is
 - **Robust against model misspecification,**
 - **Easy to implement.**
- Proposed method is promising in
 - Extensive benchmark simulations,
 - Real-world **wafer alignment task.**
- Come to our poster for technical details!