

Selective Inference for Deep Learning Model-driven Hypotheses

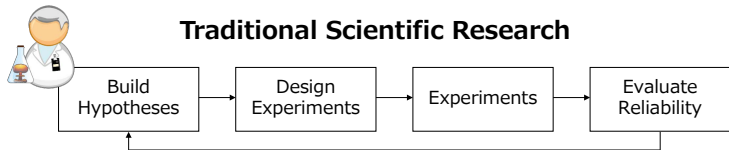
Ichiro Takeuchi

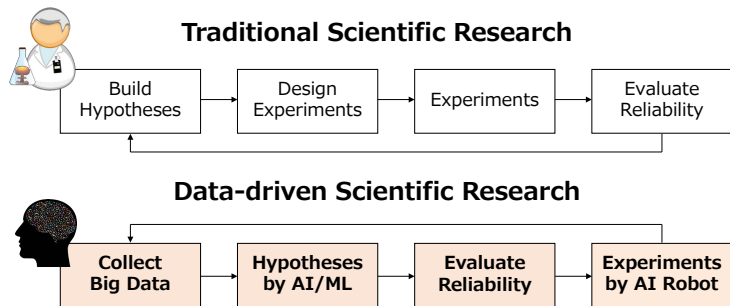
RIKEN AIP (Data-driven Biomedical Science Team)

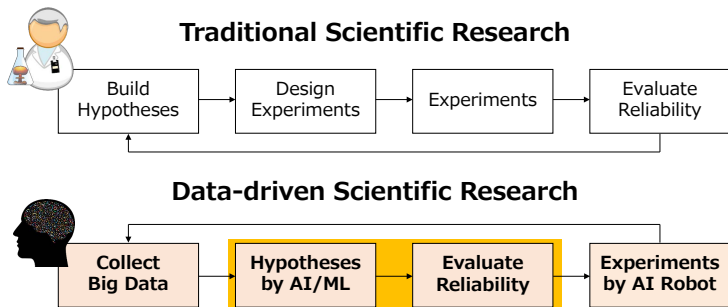
This is joint work with V.N.L Duy, D. Miwa, and S.Iwazaki.

Brief self-introduction

- ▶ Professor at Nagoya Institute of Technology
- ▶ Team leader of data-driven biomedical science team at RIKEN AIP
- ▶ Mission: develop AI and ML methods for data-driven science and their applications to biology, medicine, and material science







- ▶ ELSI (Ethical, Legal and Social Issues)
 - ▶ Fairness
 - ▶ SDG
- ▶ Interpretability
 - ▶ Visualization
 - ▶ Rule Extraction
- ▶ Reliability
 - ▶ Robustness
 - ▶ **Statistical Significance**

Reliability in AI

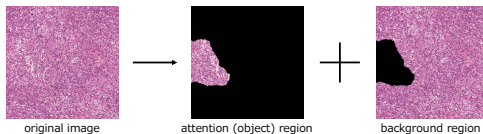
- ▶ Robustness: the complexity of AI increases the risk that a small change in the data leads to a big change in the result.



Goodfellow et al. (ICLR2015) Fig.1

- ▶ Statistical significance: The flexibility of AI increases the risk of finding false positive (FP) results (which seems meaningful but is just an artifact).

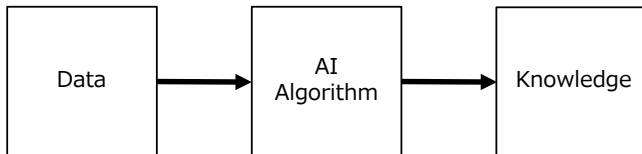
An example of medical image segmentation



(Traditional) Naive p -value = 0.000 (statistically significant)

Uncertainty quantification

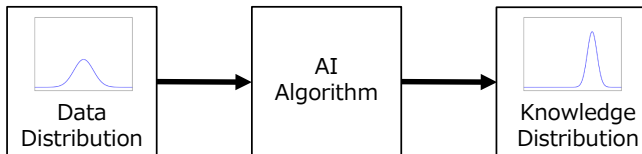
- ▶ For evaluating statistical reliability of the knowledge obtained by AI, **uncertainty quantification** of the knowledge is needed.



- ▶ Frequentist approach (sampling distribution)
 - ▶ **Exact inference (deriving exact sampling distribution)**
 - ▶ Randomized inference
 - ▶ Asymptotic inference
- ▶ Bayesian approach (posterior distribution)
 - ▶ Exact Bayesian inference
 - ▶ MCMC
 - ▶ Variational inference
- ▶ Uncertainty quantification approaches in deep neural network (DNN)
 - ▶ Dropout
 - ▶ Ensemble learning
 - ▶ Bayesian NN

Uncertainty quantification

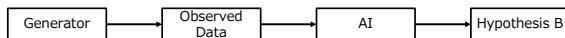
- ▶ For evaluating statistical reliability of the knowledge obtained by AI, **uncertainty quantification** of the knowledge is needed.



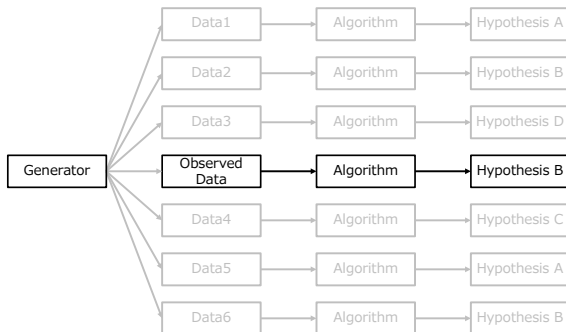
- ▶ Frequentist approach (sampling distribution)
 - ▶ **Exact inference (deriving exact sampling distribution)**
 - ▶ Randomized inference
 - ▶ Asymptotic inference
- ▶ Bayesian approach (posterior distribution)
 - ▶ Exact Bayesian inference
 - ▶ MCMC
 - ▶ Variational inference
- ▶ Uncertainty quantification approaches in deep neural network (DNN)
 - ▶ Dropout
 - ▶ Ensemble learning
 - ▶ Bayesian NN

Probabilistic data generation model

- ▶ (Frequentist) statistical inference framework (parallel world interpretation)

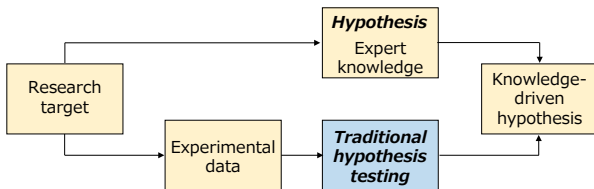


- ▶ (Frequentist) statistical inference framework (parallel world interpretation)

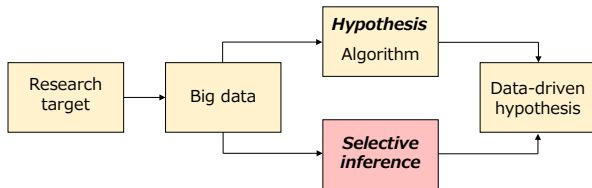


Knowledge-driven science and data-driven science

(Traditional) Knowledge-driven science



Data-driven science



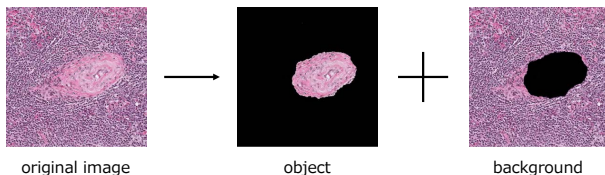
Outline

- ▶ Part 1: Hypothesis selection bias and multiple comparison
- ▶ Part 2: Conditional Selective Inference (SI)
- ▶ Part 3: Conditional SI for deep neural network (DNN)

Part 1: Hypothesis Selection Bias and Multiple Comparison

Problem 1: medical image segmentation

- ▶ Goal: Identify the attention (object) region in a medical image by segmentation



- ▶ An image is represented as an n -dimensional random vector of pixel values $X \in \mathbb{R}^n$ as

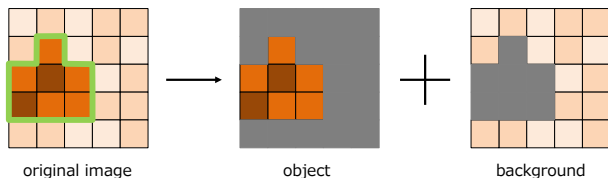
$$\underbrace{X}_{\text{(random) vector}} = \underbrace{M}_{\text{true vector}} + \underbrace{\varepsilon}_{\text{noise vector}}, \quad \underbrace{\varepsilon \sim N(0, \Sigma)}_{\text{Normally-distributed noise}}$$

- ▶ Segmentation algorithm \mathcal{A}

$$\underbrace{\mathcal{A}}_{\text{algorithm}} : \underbrace{X}_{\text{(random) image}} \mapsto \left\{ \underbrace{O_X}_{\text{pixels in object}}, \underbrace{B_X}_{\text{pixels in background}} \right\}$$

Problem 1: medical image segmentation

- ▶ Goal: Identify the attention (object) region in a medical image by segmentation



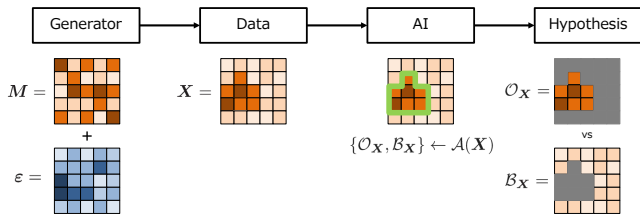
- ▶ An image is represented as an n -dimensional random vector of pixel values $\mathbf{X} \in \mathbb{R}^n$ as

$$\underbrace{\mathbf{X}}_{\text{(random) vector}} = \underbrace{\mathbf{M}}_{\text{true vector}} + \underbrace{\boldsymbol{\varepsilon}}_{\text{noise vector}}, \quad \underbrace{\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \Sigma)}_{\text{Normally-distributed noise}}$$

- ▶ Segmentation algorithm \mathcal{A}

$$\underbrace{\mathcal{A}}_{\text{algorithm}} : \underbrace{\mathbf{X}}_{\text{(random) image}} \mapsto \left\{ \underbrace{\mathcal{O}_{\mathbf{X}}}_{\text{pixels in object}}, \underbrace{\mathcal{B}_{\mathbf{X}}}_{\text{pixels in background}} \right\}$$

▶ Data-driven hypothesis



Statistical hypothesis testing

- ▶ Statistical hypothesis testing
 - ▶ Null hypothesis

$$H_0 : \underbrace{\frac{1}{|\mathcal{O}_X|} \sum_{i \in \mathcal{O}_X} M_i}_{\text{mean pixel value in object}} = \underbrace{\frac{1}{|\mathcal{B}_X|} \sum_{i \in \mathcal{B}_X} M_i}_{\text{mean pixel value in background}}$$

- ▶ Alternative hypothesis

$$H_1 : \underbrace{\frac{1}{|\mathcal{O}_X|} \sum_{i \in \mathcal{O}_X} M_i}_{\text{mean pixel value in object}} \neq \underbrace{\frac{1}{|\mathcal{B}_X|} \sum_{i \in \mathcal{B}_X} M_i}_{\text{mean pixel value in background}}$$

- ▶ Test statistic: Difference of mean pixel values between object and background regions

$$\Delta_X := \frac{1}{|\mathcal{O}_X|} \sum_{i \in \mathcal{O}_X} X_i - \frac{1}{|\mathcal{B}_X|} \sum_{i \in \mathcal{B}_X} X_i$$

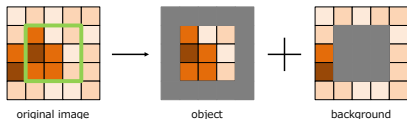
- ▶ Statistical significance (two-sided p -value)

$$p = \Pr \left(\underbrace{|\Delta_X|}_{\text{random variable}} \geq \underbrace{|\Delta_x|}_{\text{observation}} \right)$$

Knowledge-driven vs. data-driven hypotheses

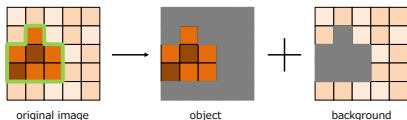
- ▶ Knowledge-driven hypothesis: object/background regions do not depend on the data \Rightarrow (traditional) z -test or t -test

$$\underbrace{\frac{1}{|\mathcal{O}|} \sum_{i \in \mathcal{O}} M_i}_{\text{fixed mean pixel value in object}} = \text{or } \neq \underbrace{\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} M_i}_{\text{fixed mean pixel value in background}}$$



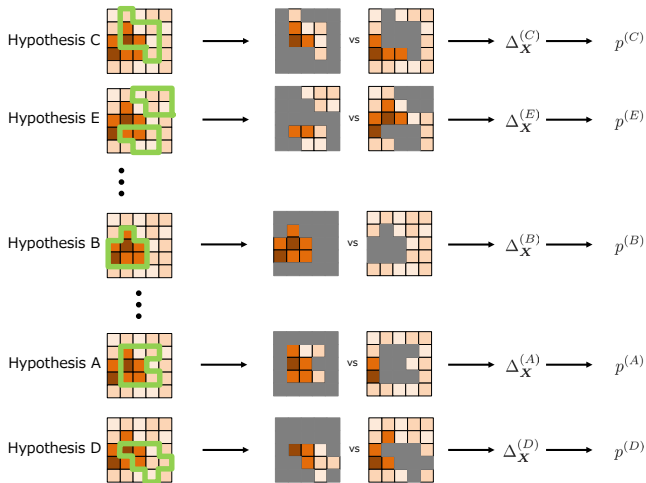
- ▶ Data-driven hypothesis: object/background regions are determined by the data \Rightarrow data/algorithm dependent

$$\underbrace{\frac{1}{|\mathcal{O}_x|} \sum_{i \in \mathcal{O}_x} M_i}_{\text{selected mean pixel value in object}} = \text{or } \neq \underbrace{\frac{1}{|\mathcal{B}_x|} \sum_{i \in \mathcal{B}_x} M_i}_{\text{selected mean pixel value in background}}$$



Multiple comparison, hypothesis selection, and selection bias

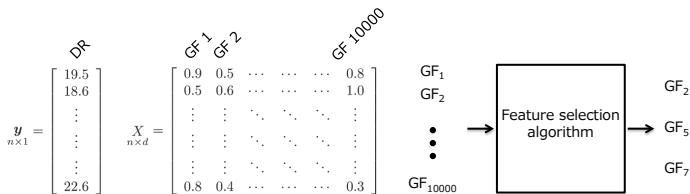
- The data-driven hypothesis is interpreted as the result of multiple comparison with all possible $2^{\#\text{pixels}}$ segmentation results.



- Correction of the selection bias is indispensable in multiple comparison.

Problem 2: feature selection in linear models

- ▶ Goal: select a subset of 10,000 genes that are useful for predicting drug effects
- ▶ High-dimensional data and feature selection



- ▶ Linear model with the selected features by least-square method

$$\hat{y}_i = \hat{\beta}_2 x_{i2} + \hat{\beta}_5 x_{i5} + \hat{\beta}_7 x_{i7},$$

where

$$\begin{bmatrix} \hat{\beta}_2 \\ \hat{\beta}_5 \\ \hat{\beta}_7 \end{bmatrix} = \underset{\beta_2, \beta_5, \beta_7}{\operatorname{argmin}} \sum_{i=1}^n (y_i - (\beta_2 x_{i2} + \beta_5 x_{i5} + \beta_7 x_{i7}))^2$$

Problem 2: feature selection in linear models (problem formulation)

- ▶ Data ($n = 50$, $d = 10000$ in the example)

$$\mathbf{X} \in \mathbb{R}^{n \times d}, \mathbf{Y} \in \mathbb{R}^n$$

- ▶ Probabilistic model

$$\underbrace{\mathbf{Y}}_{\text{drug effect}} = \underbrace{\boldsymbol{\mu}(\mathbf{X})}_{\text{true drug effect}} + \underbrace{\boldsymbol{\varepsilon}}_{\text{noise}}, \quad \underbrace{\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \Sigma)}_{\text{Normally-distributed noise}}$$

- ▶ Feature selection algorithm \mathcal{A}

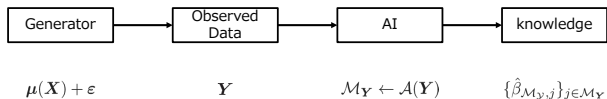
$$\mathcal{A} : \mathbf{Y} \mapsto \mathcal{M}_{\mathbf{Y}},$$

where $\mathcal{M}_{\mathbf{Y}}$ is the set of selected features ($\mathcal{M}_{\mathbf{y}} = \{2, 5, 7\}$ in the example)

- ▶ Linear model with the selected features by least-square method

$$\hat{\boldsymbol{\beta}}_{\mathcal{M}_{\mathbf{Y}}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^{|\mathcal{M}_{\mathbf{Y}}|}}{\operatorname{argmin}} \|\mathbf{Y} - \mathbf{X}_{\mathcal{M}_{\mathbf{Y}}}^{\top} \boldsymbol{\beta}\|^2 = (\mathbf{X}_{\mathcal{M}_{\mathbf{Y}}}^{\top} \mathbf{X}_{\mathcal{M}_{\mathbf{Y}}})^{-1} \mathbf{X}_{\mathcal{M}_{\mathbf{Y}}}^{\top} \mathbf{Y}$$

► Data-driven hypothesis



Statistical hypothesis testing

- ▶ Statistical hypothesis testing
 - ▶ Population least-square solution

$$\beta_{\mathcal{M}_Y} = (X_{\mathcal{M}_Y}^\top X_{\mathcal{M}_Y})^{-1} X_{\mathcal{M}_Y}^\top \underbrace{\mu(\mathbf{X})}_{\text{true drug effect}}$$

- ▶ Null hypothesis

$$H_0 : \underbrace{\beta_{\mathcal{M}_Y, j}}_{\text{effect of the selected feature } j} = 0$$

- ▶ Alternative hypothesis

$$H_1 : \underbrace{\beta_{\mathcal{M}_Y, j}}_{\text{effect of the selected feature } j} \neq 0$$

- ▶ Test-statistic

$$\hat{\beta}_{\mathcal{M}_Y, j} = (X_{\mathcal{M}_Y}^\top X_{\mathcal{M}_Y})^{-1} X_{\mathcal{M}_Y}^\top \underbrace{\mathbf{Y}}_{\text{observed drug effect}}$$

- ▶ Statistical significance (two-sided p -value)

$$p = \Pr \left(\underbrace{|\hat{\beta}_{\mathcal{M}_Y, j}|}_{\text{random var.}} \geq \underbrace{|\hat{\beta}_{\mathcal{M}_y, j}|}_{\text{observation}} \right)$$

Knowledge-driven hypotheses and data-driven hypotheses

- ▶ Knowledge-driven hypothesis: the set of features are selected without looking at the data \Rightarrow (traditional) z -test or t -test

$$\underbrace{\beta_{\mathcal{M},j}}_{\text{effect of the selected feature } j \text{ for the fixed model}} = \text{or } \neq 0$$

- ▶ Data-driven hypothesis: the set of features are selected by the data \Rightarrow data/algorithm dependent

$$\underbrace{\beta_{\mathcal{M}_Y,j}}_{\text{effect of the selected feature } j \text{ for the selected model}} = \text{or } \neq 0$$

Multiple comparison, hypothesis selection, and selection bias

- ▶ This data-driven hypothesis is interpreted as the result of multiple comparison with $2^{\# \text{features}} \times \# \text{selected feature}$ hypotheses.

$$\text{Hypothesis C } \mathcal{M}^{(C)} = \{3, 7\} \longrightarrow \begin{array}{l} \beta_{\mathcal{M}^{(C)},3} \neq 0 \\ \beta_{\mathcal{M}^{(C)},7} \neq 0 \end{array} \longrightarrow \begin{array}{l} \hat{\beta}_{\mathcal{M}^{(C)},3} \\ \hat{\beta}_{\mathcal{M}^{(C)},7} \end{array} \longrightarrow \begin{array}{l} p_3^{(C)} \\ p_7^{(C)} \end{array}$$

$$\text{Hypothesis E } \mathcal{M}^{(E)} = \{2\} \longrightarrow \beta_{\mathcal{M}^{(E)},2} \neq 0 \longrightarrow \hat{\beta}_{\mathcal{M}^{(E)},2} \longrightarrow p_2^{(E)}$$

⋮

$$\text{Hypothesis B } \mathcal{M}^{(B)} = \{2, 5, 7\} \longrightarrow \begin{array}{l} \beta_{\mathcal{M}^{(B)},2} \neq 0 \\ \beta_{\mathcal{M}^{(B)},5} \neq 0 \\ \beta_{\mathcal{M}^{(B)},7} \neq 0 \end{array} \longrightarrow \begin{array}{l} \hat{\beta}_{\mathcal{M}^{(B)},2} \\ \hat{\beta}_{\mathcal{M}^{(B)},5} \\ \hat{\beta}_{\mathcal{M}^{(B)},7} \end{array} \longrightarrow \begin{array}{l} p_2^{(B)} \\ p_5^{(B)} \\ p_7^{(B)} \end{array}$$

⋮

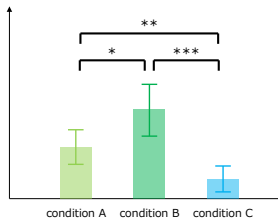
$$\text{Hypothesis A } \mathcal{M}^{(A)} = \{9\} \longrightarrow \beta_{\mathcal{M}^{(A)},9} \neq 0 \longrightarrow \hat{\beta}_{\mathcal{M}^{(A)},9} \longrightarrow p_9^{(A)}$$

$$\text{Hypothesis D } \mathcal{M}^{(D)} = \{3, 8\} \longrightarrow \begin{array}{l} \beta_{\mathcal{M}^{(D)},3} \neq 0 \\ \beta_{\mathcal{M}^{(D)},8} \neq 0 \end{array} \longrightarrow \begin{array}{l} \hat{\beta}_{\mathcal{M}^{(D)},3} \\ \hat{\beta}_{\mathcal{M}^{(D)},8} \end{array} \longrightarrow \begin{array}{l} p_3^{(D)} \\ p_8^{(D)} \end{array}$$

- ▶ Correction of the selection bias is indispensable in multiple comparison.

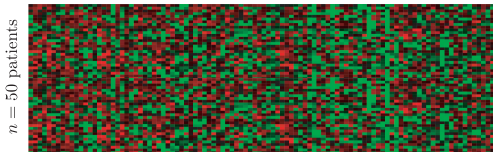
Multiple comparison

- ▶ In the context of traditional multiple hypothesis testing, only a handful of tests are considered.



- ▶ In the context of genetic data analysis (2000~), large-scale multiple comparison with tens of thousands of hypotheses were considered.

$d = 10000$ genes



- ▶ The number of all possible hypotheses that AI/ML can produce is much more than the existing methods can handle.

Three approaches for multiple comparison correction

- ▶ **Family-wise error rate (FWER) control**: controlling the probability of finding a false positive (FP) $< \alpha$ (e.g., 0.05)
- ▶ **False discover rate (FDR)**: controlling the expected proportion of discoveries that are false $< \alpha$ (e.g., 0.05)
- ▶ **Conditional selective inference (SI)**: controlling the probability of finding a FP conditional on the hypothesis selection event $< \alpha$ (e.g., 0.05)

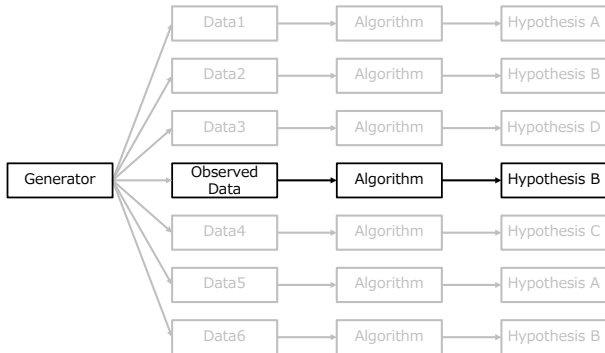
Summary of part 1

- ▶ Knowledge obtained by AI/ML algorithm is considered as data-driven hypotheses.
- ▶ Statistical reliability of data-driven hypotheses cannot be properly evaluated with traditional statistical inference due to the selection bias.
- ▶ This problem can be interpreted as a huge-scale multiple comparison problem where the one is selected from all possible hypotheses that AI/ML can produce.

Part 2: Conditional Selective Inference (SI)

Basic idea of conditional SI

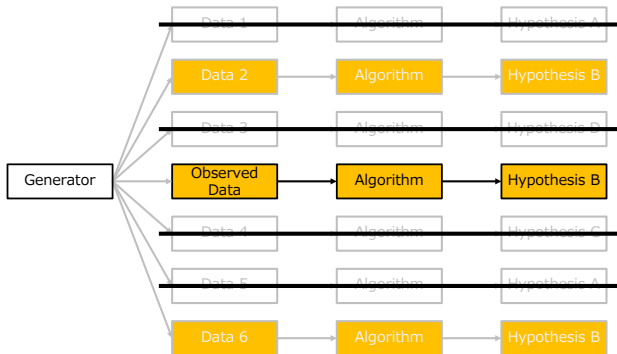
- ▶ The key idea of conditional SI is to consider only the cases (parallel worlds) where the same hypothesis is selected.



- ▶ Intuitively, by considering only the randomness where the same hypothesis is selected, the hypothesis selection bias disappears.

Basic idea of conditional SI

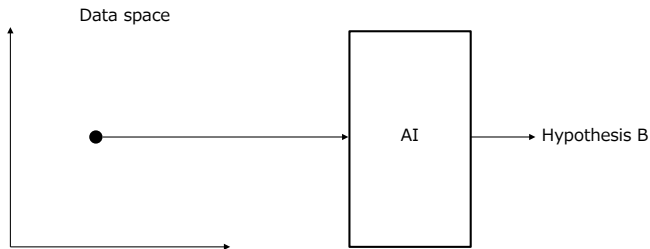
- ▶ The key idea of conditional SI is to consider only the cases (parallel worlds) where the same hypothesis is selected.



- ▶ Intuitively, by considering only the randomness where the same hypothesis is selected, the hypothesis selection bias disappears.

Basic idea of conditional SI

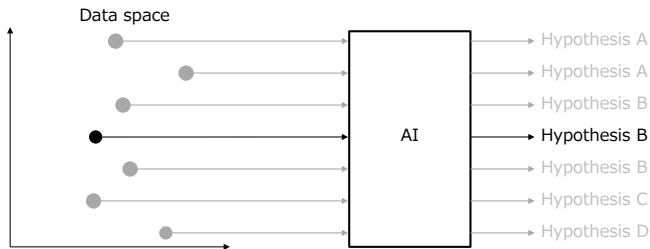
- ▶ The key idea of conditional SI is to consider only the cases (parallel worlds) where the same hypothesis is selected.



- ▶ Intuitively, by considering only the randomness in the subset of the data space, the hypothesis selection bias disappears.

Basic idea of conditional SI

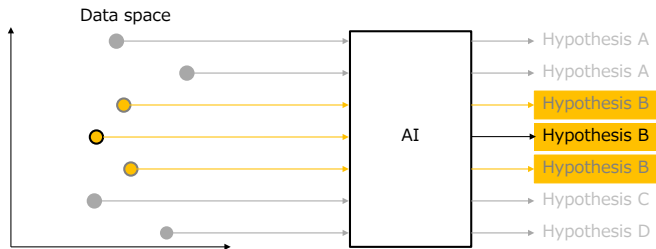
- ▶ The key idea of conditional SI is to consider only the cases (parallel worlds) where the same hypothesis is selected.



- ▶ Intuitively, by considering only the randomness in the subset of the data space, the hypothesis selection bias disappears.

Basic idea of conditional SI

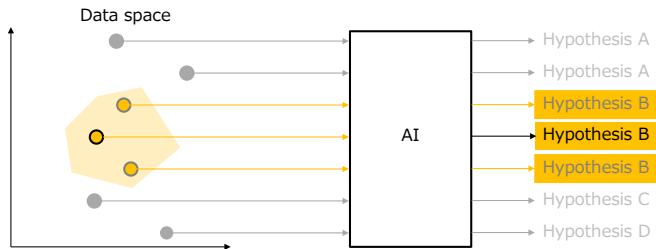
- ▶ The key idea of conditional SI is to consider only the cases (parallel worlds) where the same hypothesis is selected.



- ▶ Intuitively, by considering only the randomness in the subset of the data space, the hypothesis selection bias disappears.

Basic idea of conditional SI

- ▶ The key idea of conditional SI is to consider only the cases (parallel worlds) where the same hypothesis is selected.



- ▶ Intuitively, by considering only the randomness in the subset of the data space, the hypothesis selection bias disappears.

Conditional SI for medical image segmentation problem

- ▶ Ordinary statistical significance (p -value)

$$p = \Pr \left(\underbrace{|\Delta_{\mathbf{X}}|}_{\text{random var.}} \geq \underbrace{|\Delta_{\mathbf{x}}|}_{\text{observation}} \right)$$

- ▶ Conditional statistical significance (selective p -value)

$$p = \Pr \left(\underbrace{|\Delta_{\mathbf{X}}|}_{\text{random var.}} \geq \underbrace{|\Delta_{\mathbf{x}}|}_{\text{observation}} \mid \underbrace{\{\mathcal{O}_{\mathbf{X}}, \mathcal{B}_{\mathbf{X}}\} = \{\mathcal{O}_{\mathbf{x}}, \mathcal{B}_{\mathbf{x}}\}}_{\text{the same object/background are selected}} \right)$$

- ▶ The main challenge of conditional SI is to characterize the selection event and compute the conditional probability.

A simple example of conditional SI for image segmentation

- ▶ A simple segmentation algorithm based on a threshold θ

$$\mathcal{O}_X = \{X_i \geq \theta\},$$

$$\mathcal{B}_X = \{X_i < \theta\}$$

- ▶ Selection event

$$X_i \geq \theta, i \in \mathcal{O}_x,$$

$$X_i < \theta, i \in \mathcal{B}_x$$

- ▶ Selective p -value

$$\begin{aligned} p &= \Pr \left(\underbrace{|\Delta_X|}_{\text{random var.}} \geq \underbrace{|\Delta_x|}_{\text{observation}} \mid \underbrace{\{\mathcal{O}_X, \mathcal{B}_X\} = \{\mathcal{O}_x, \mathcal{B}_x\}}_{\text{the same object/background are selected}} \right) \\ &= \Pr \left(\underbrace{|\Delta_X|}_{\text{random var.}} \geq \underbrace{|\Delta_x|}_{\text{observation}} \mid \underbrace{X_i \geq \theta, i \in \mathcal{O}_x, X_i < \theta, i \in \mathcal{B}_x}_{\text{the same object/background are selected}} \right) \end{aligned}$$

Conditional SI for feature selection

- ▶ Naive p -value:

$$p = \Pr \left(\underbrace{|\hat{\beta}_{\mathcal{M}_{\mathbf{Y}},j}|}_{\text{random var.}} \geq \underbrace{|\hat{\beta}_{\mathcal{M}_{\mathbf{y}},j}|}_{\text{observation}} \right)$$

- ▶ Selective p -value:

$$p = \Pr \left(\underbrace{|\hat{\beta}_{\mathcal{M}_{\mathbf{Y}},j}|}_{\text{random var.}} \geq \underbrace{|\hat{\beta}_{\mathcal{M}_{\mathbf{y}},j}|}_{\text{observation}} \mid \underbrace{\mathcal{M}_{\mathbf{Y}} = \mathcal{M}_{\mathbf{y}}}_{\text{the same set of features are selected}} \right)$$

A simple example of conditional SI for feature selection

- ▶ Marginal screening: select k features whose correlation between $\mathbf{Y} \in \mathbb{R}^n$ and $\mathbf{x}_j \in \mathbb{R}^n$ are large:

$$\underbrace{\mathbf{x}_{(1)}^\top \mathbf{Y} \geq \mathbf{x}_{(2)}^\top \mathbf{Y} \geq \mathbf{x}_{(3)}^\top \mathbf{Y}}_{\text{selected (when } k=3\text{)}} \geq \underbrace{\mathbf{x}_{(4)}^\top \mathbf{Y} \geq \mathbf{x}_{(5)}^\top \mathbf{Y} \geq \dots \geq \mathbf{x}_{(D)}^\top \mathbf{Y}}_{\text{not selected (when } k=3\text{)}}$$

Note that the correlation is represented as an inner product when variables are standardized.

- ▶ Hypothesis selection event

$$\begin{array}{ccc} \mathbf{x}_{(1)}^\top \mathbf{Y} \geq \mathbf{x}_{(4)}^\top \mathbf{Y} & \mathbf{x}_{(2)}^\top \mathbf{Y} \geq \mathbf{x}_{(4)}^\top \mathbf{Y} & \mathbf{x}_{(3)}^\top \mathbf{Y} \geq \mathbf{x}_{(4)}^\top \mathbf{Y} \\ \vdots & \vdots & \vdots \\ \mathbf{x}_{(1)}^\top \mathbf{Y} \geq \mathbf{x}_{(D)}^\top \mathbf{Y} & \mathbf{x}_{(2)}^\top \mathbf{Y} \geq \mathbf{x}_{(D)}^\top \mathbf{Y} & \mathbf{x}_{(3)}^\top \mathbf{Y} \geq \mathbf{x}_{(D)}^\top \mathbf{Y}, \end{array}$$

- ▶ Selective p -value

$$\begin{aligned} p &= \Pr \left(\underbrace{|\hat{\beta}_{\mathcal{M}_{\mathbf{Y}},j}|}_{\text{random var.}} \geq \underbrace{|\hat{\beta}_{\mathcal{M}_{\mathbf{y}},j}|}_{\text{observation}} \mid \underbrace{\mathcal{M}_{\mathbf{Y}} = \mathcal{M}_{\mathbf{y}}}_{\text{the same set of features are selected}} \right) \\ &= \Pr \left(\underbrace{|\hat{\beta}_{\mathcal{M}_{\mathbf{Y}},j}|}_{\text{random var.}} \geq \underbrace{|\hat{\beta}_{\mathcal{M}_{\mathbf{y}},j}|}_{\text{observation}} \mid \underbrace{\left\{ \mathbf{x}_{(\ell)}^\top \mathbf{Y} \geq \mathbf{x}_{(m)}^\top \mathbf{Y} \right\}_{(\ell,m) \in \{1,\dots,k\} \times \{k+1,\dots,d\}}}_{\text{the same set of features are selected}} \right) \end{aligned}$$

Hypothesis selection event

Segmentation problem

Hypothesis selection event

$$X_i \geq \theta \text{ if } x_i \in \mathcal{O}_{\mathbf{x}}, X_i < \theta \text{ if } x_i \in \mathcal{B}_{\mathbf{x}}$$

Feature selection problem

Hypothesis selection event

$$\{\mathbf{x}^{(\ell)\top} \mathbf{Y} \geq \mathbf{x}^{(m)\top} \mathbf{Y}\}_{(\ell,m) \in \{1, \dots, k\} \times \{k+1, \dots, d\}}$$

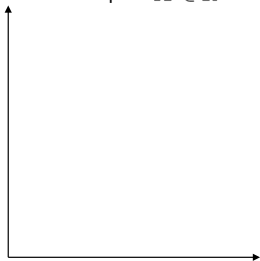
Hypothesis selection event

Segmentation problem

Hypothesis selection event

$$X_i \geq \theta \text{ if } x_i \in \mathcal{O}_x, X_i < \theta \text{ if } x_i \in \mathcal{B}_x$$

Data space $\mathbf{X} \in \mathbb{R}^n$

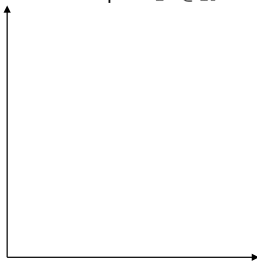


Feature selection problem

Hypothesis selection event

$$\{\mathbf{x}^{(\ell)\top} \mathbf{Y} \geq \mathbf{x}^{(m)\top} \mathbf{Y}\}_{(\ell,m) \in \{1,\dots,k\} \times \{k+1,\dots,d\}}$$

Data space $\mathbf{Y} \in \mathbb{R}^n$

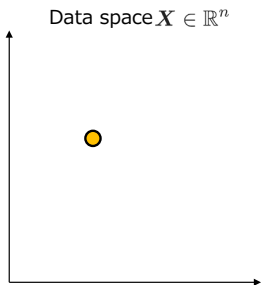


Hypothesis selection event

Segmentation problem

Hypothesis selection event

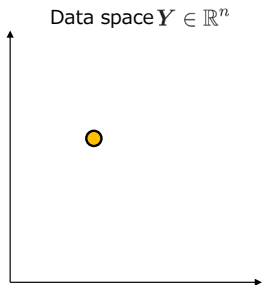
$$X_i \geq \theta \text{ if } x_i \in \mathcal{O}_x, X_i < \theta \text{ if } x_i \in \mathcal{B}_x$$



Feature selection problem

Hypothesis selection event

$$\{\mathbf{x}^{(\ell)\top} \mathbf{Y} \geq \mathbf{x}^{(m)\top} \mathbf{Y}\}_{(\ell, m) \in \{1, \dots, k\} \times \{k+1, \dots, d\}}$$

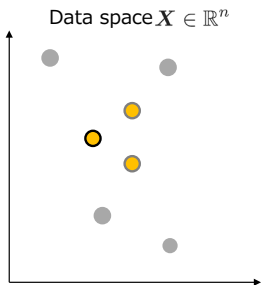


Hypothesis selection event

Segmentation problem

Hypothesis selection event

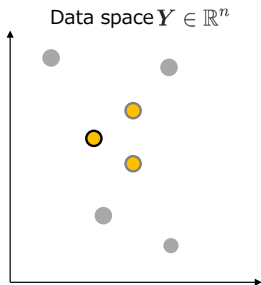
$$X_i \geq \theta \text{ if } x_i \in \mathcal{O}_x, X_i < \theta \text{ if } x_i \in \mathcal{B}_x$$



Feature selection problem

Hypothesis selection event

$$\{\mathbf{x}^{(\ell)\top} \mathbf{Y} \geq \mathbf{x}^{(m)\top} \mathbf{Y}\}_{(\ell, m) \in \{1, \dots, k\} \times \{k+1, \dots, d\}}$$

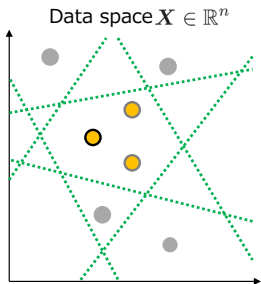


Hypothesis selection event

Segmentation problem

Hypothesis selection event

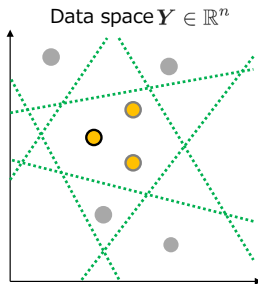
$$X_i \geq \theta \text{ if } x_i \in \mathcal{O}_x, X_i < \theta \text{ if } x_i \in \mathcal{B}_x$$



Feature selection problem

Hypothesis selection event

$$\{\mathbf{x}^{(\ell)\top} \mathbf{Y} \geq \mathbf{x}^{(m)\top} \mathbf{Y}\}_{(\ell,m) \in \{1,\dots,k\} \times \{k+1,\dots,d\}}$$

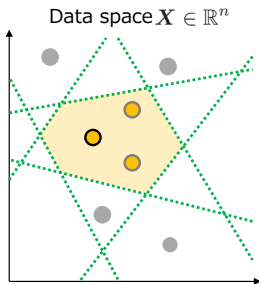


Hypothesis selection event

Segmentation problem

Hypothesis selection event

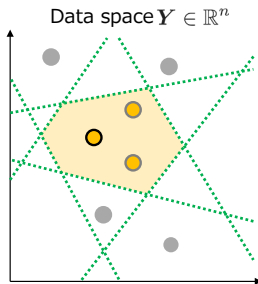
$$X_i \geq \theta \text{ if } x_i \in \mathcal{O}_x, X_i < \theta \text{ if } x_i \in \mathcal{B}_x$$



Feature selection problem

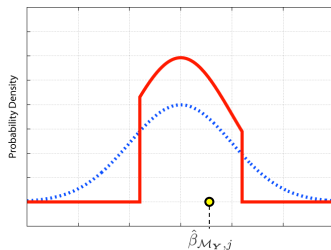
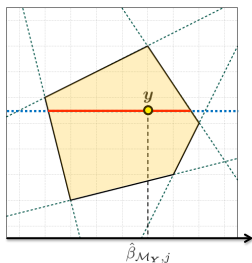
Hypothesis selection event

$$\{\mathbf{x}^{(\ell)\top} \mathbf{Y} \geq \mathbf{x}^{(m)\top} \mathbf{Y}\}_{(\ell,m) \in \{1,\dots,k\} \times \{k+1,\dots,d\}}$$



Polyhedral lemma (Lee+16)

- ▶ Conditional SI has been actively studied after the seminal work by (Lee+2016)
- ▶ Conditional SI for Lasso feature selection was studied in (Lee+2016)
- ▶ Brief summary of polyhedral lemma: If
 1. the selection event is represented by a polyhedron (a set of linear inequalities) in the data space, and
 2. the test-statistic is linear function of the data,then the exact selective p -values can be computed based on truncated Normal distribution.
- ▶ By further conditioning on the sufficient statistic of the nuisance component, we have the sampling distribution on a line in the data space truncated by a polyhedron.



Conditional SI for the selected features by Lasso (Lee+16)

- ▶ Consider Lasso is an algorithm to select a set of features and their signs

$$\mathcal{A}^{\text{Lasso}} : \mathbf{Y} \mapsto \{\mathcal{M}, \mathbf{s}\},$$

where \mathcal{M} is the set of selected features and \mathbf{s} is the set of their signs.

- ▶ Test-statistic of the selected features

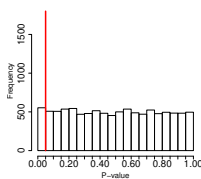
$$\hat{\beta}_{\mathcal{M}_{\mathbf{Y}},j} = (X_{\mathcal{M}_{\mathbf{Y}}}^{\top} X_{\mathcal{M}_{\mathbf{Y}}})^{-1} X_{\mathcal{M}_{\mathbf{Y}}}^{\top} \mathbf{Y} = \boldsymbol{\eta}^{\top} \mathbf{Y}$$

- ▶ Selective p -value

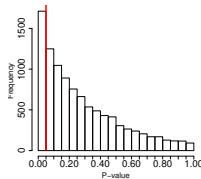
$$p = \Pr \left(\left| \hat{\beta}_{\mathcal{M}_{\mathbf{Y}},j} \right| \geq \left| \hat{\beta}_{\mathcal{M}_{\mathbf{y}},j} \right| \mid \underbrace{\mathcal{M}_{\mathbf{Y}} = \mathcal{M}_{\mathbf{y}}}_{\text{features}}, \underbrace{\mathbf{s}_{\mathbf{Y}} = \mathbf{s}_{\mathbf{y}}}_{\text{signs}}, \underbrace{P_{\boldsymbol{\eta}}^{\perp} \mathbf{Y} = P_{\boldsymbol{\eta}}^{\perp} \mathbf{y}}_{\text{nuisance component}} \right)$$

- ▶ Selective p -values follow uniform distribution

$$\Pr_{H_0} (p \leq \alpha \mid \mathcal{M}_{\mathbf{Y}} = \mathcal{M}_{\mathbf{y}}, \mathbf{s}_{\mathbf{Y}} = \mathbf{s}_{\mathbf{y}}) = \alpha \quad \forall \alpha \in (0, 1)$$



Selective p -value distribution



Naive p -value distribution

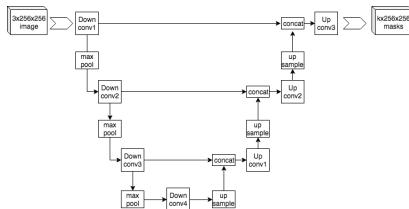
Summary of part 2

- ▶ Conditional SI has been actively studied as a promising approach for hypothesis selection bias correction.
- ▶ By polyhedral lemma, if the selection event is represented by a polyhedron, the selective p -values can be computed.
- ▶ Selection event depends on each algorithm — it is challenging to apply conditional SI to complicated algorithms such as DNN.

Part 3: Conditional SI for DNN-driven Hypotheses

Image segmentation by DNN

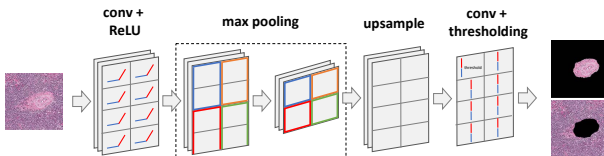
- U-net is one of the most-commonly used CNN for image segmentation task:



Basic structure of U-Net (Wikipedia)

U-net is fully convolutional network and has U-shape.

- Basic components of CNN



- CNN is a complicated function as a whole, but it consists of several basic simple components.

Selection event for DNN

- ▶ CNN-based segmentation algorithm

$$\mathcal{A}^{\text{CNN}} : \mathbf{X} \mapsto \{\mathcal{O}_{\mathbf{X}}, \mathcal{B}_{\mathbf{X}}\}$$

- ▶ Conditional SI for CNN-based segmentation

$$p = \Pr \left(|\Delta_{\mathbf{X}}| \geq |\Delta_{\mathbf{x}}| \mid \underbrace{\mathcal{A}^{\text{CNN}}(\mathbf{X}) = \mathcal{A}^{\text{CNN}}(\mathbf{x})}_{\text{the same CNN outputs are selected}}, \underbrace{P_{\eta}^{\perp} \mathbf{X} = P_{\eta}^{\perp} \mathbf{x}}_{\text{nuisance component}} \right)$$

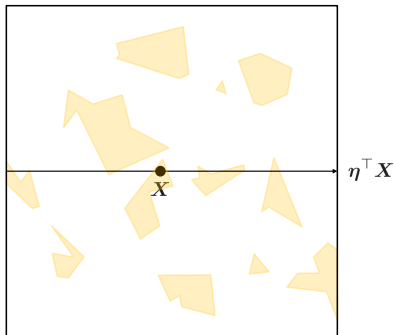
- ▶ Q. Can we characterize the complicated selection event of CNN?

$$\{\mathcal{O}_{\mathbf{X}}, \mathcal{B}_{\mathbf{X}}\} = \{\mathcal{O}_{\mathbf{x}}, \mathcal{B}_{\mathbf{x}}\} \Leftrightarrow \mathcal{A}^{\text{CNN}}(\mathbf{X}) = \mathcal{A}^{\text{CNN}}(\mathbf{x})$$

Unfortunately, the selection event cannot be represented as a polyhedron.

Parametric programming approach

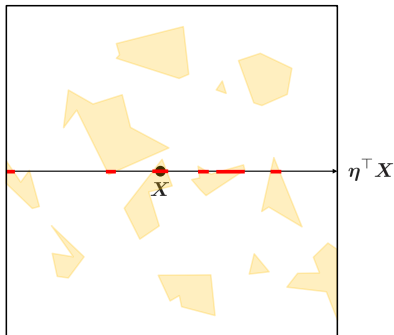
- ▶ The selection event of CNN segmentation algorithm is complicated:



- ▶ Our idea is to consider solving a sequence of segmentation problems for a parametrized data in the direction of the test-statistic:
 1. Consider multiple finer selection events with additional conditions;
 2. Run the segmentation algorithm for each finer selection event on the line;
 3. Identify the truncation region at which the same result is obtained;
 4. Combine the probability mass of multiply truncated Normal distributions;

Parametric programming approach

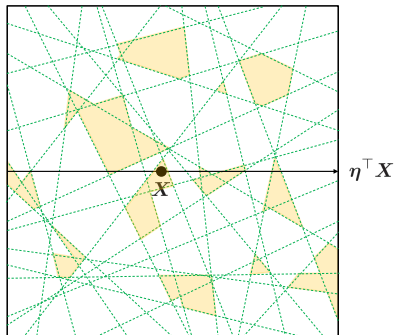
- ▶ The selection event of CNN segmentation algorithm is complicated:



- ▶ Our idea is to consider solving a sequence of segmentation problems for a parametrized data in the direction of the test-statistic:
 1. Consider multiple finer selection events with additional conditions;
 2. Run the segmentation algorithm for each finer selection event on the line;
 3. Identify the truncation region at which the same result is obtained;
 4. Combine the probability mass of multiply truncated Normal distributions;

Parametric programming approach

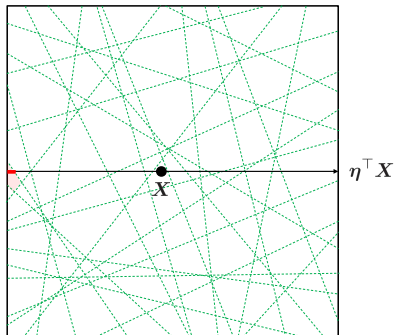
- ▶ The selection event of CNN segmentation algorithm is complicated:



- ▶ Our idea is to consider solving a sequence of segmentation problems for a parametrized data in the direction of the test-statistic:
 1. Consider multiple finer selection events with additional conditions;
 2. Run the segmentation algorithm for each finer selection event on the line;
 3. Identify the truncation region at which the same result is obtained;
 4. Combine the probability mass of multiply truncated Normal distributions;

Parametric programming approach

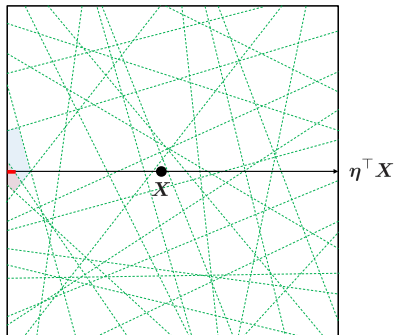
- ▶ The selection event of CNN segmentation algorithm is complicated:



- ▶ Our idea is to consider solving a sequence of segmentation problems for a parametrized data in the direction of the test-statistic:
 1. Consider multiple finer selection events with additional conditions;
 2. Run the segmentation algorithm for each finer selection event on the line;
 3. Identify the truncation region at which the same result is obtained;
 4. Combine the probability mass of multiply truncated Normal distributions;

Parametric programming approach

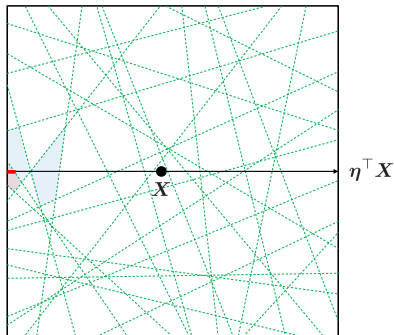
- ▶ The selection event of CNN segmentation algorithm is complicated:



- ▶ Our idea is to consider solving a sequence of segmentation problems for a parametrized data in the direction of the test-statistic:
 1. Consider multiple finer selection events with additional conditions;
 2. Run the segmentation algorithm for each finer selection event on the line;
 3. Identify the truncation region at which the same result is obtained;
 4. Combine the probability mass of multiply truncated Normal distributions;

Parametric programming approach

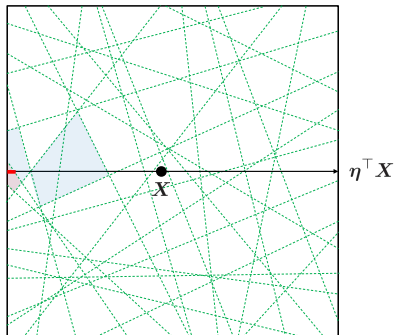
- ▶ The selection event of CNN segmentation algorithm is complicated:



- ▶ Our idea is to consider solving a sequence of segmentation problems for a parametrized data in the direction of the test-statistic:
 1. Consider multiple finer selection events with additional conditions;
 2. Run the segmentation algorithm for each finer selection event on the line;
 3. Identify the truncation region at which the same result is obtained;
 4. Combine the probability mass of multiply truncated Normal distributions;

Parametric programming approach

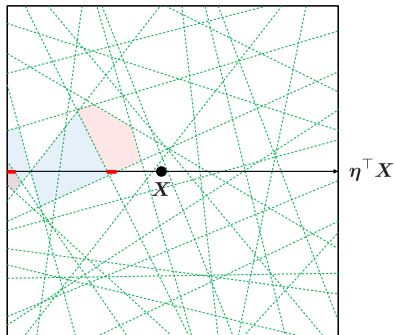
- ▶ The selection event of CNN segmentation algorithm is complicated:



- ▶ Our idea is to consider solving a sequence of segmentation problems for a parametrized data in the direction of the test-statistic:
 1. Consider multiple finer selection events with additional conditions;
 2. Run the segmentation algorithm for each finer selection event on the line;
 3. Identify the truncation region at which the same result is obtained;
 4. Combine the probability mass of multiply truncated Normal distributions;

Parametric programming approach

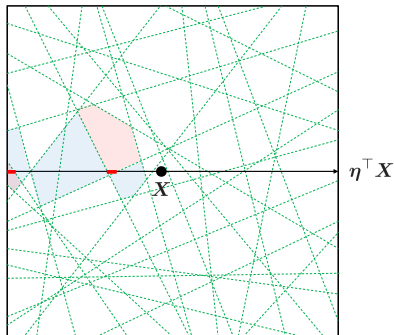
- ▶ The selection event of CNN segmentation algorithm is complicated:



- ▶ Our idea is to consider solving a sequence of segmentation problems for a parametrized data in the direction of the test-statistic:
 1. Consider multiple finer selection events with additional conditions;
 2. Run the segmentation algorithm for each finer selection event on the line;
 3. Identify the truncation region at which the same result is obtained;
 4. Combine the probability mass of multiply truncated Normal distributions;

Parametric programming approach

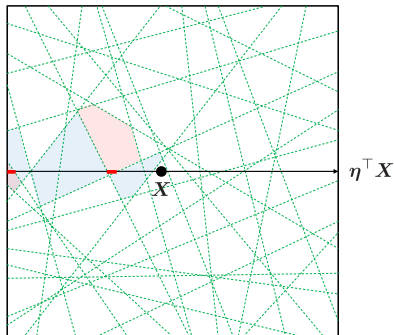
- ▶ The selection event of CNN segmentation algorithm is complicated:



- ▶ Our idea is to consider solving a sequence of segmentation problems for a parametrized data in the direction of the test-statistic:
 1. Consider multiple finer selection events with additional conditions;
 2. Run the segmentation algorithm for each finer selection event on the line;
 3. Identify the truncation region at which the same result is obtained;
 4. Combine the probability mass of multiply truncated Normal distributions;

Parametric programming approach

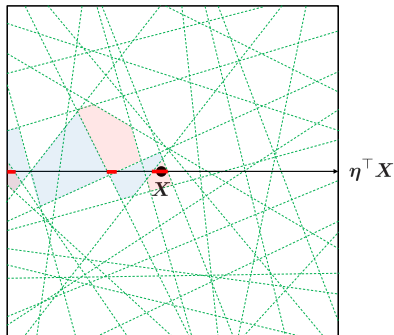
- ▶ The selection event of CNN segmentation algorithm is complicated:



- ▶ Our idea is to consider solving a sequence of segmentation problems for a parametrized data in the direction of the test-statistic:
 1. Consider multiple finer selection events with additional conditions;
 2. Run the segmentation algorithm for each finer selection event on the line;
 3. Identify the truncation region at which the same result is obtained;
 4. Combine the probability mass of multiply truncated Normal distributions;

Parametric programming approach

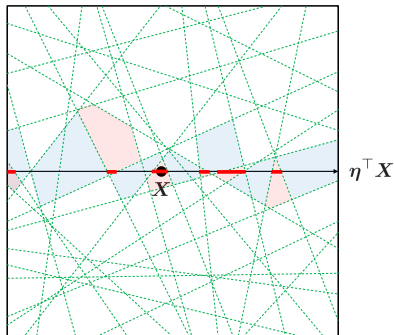
- ▶ The selection event of CNN segmentation algorithm is complicated:



- ▶ Our idea is to consider solving a sequence of segmentation problems for a parametrized data in the direction of the test-statistic:
 1. Consider multiple finer selection events with additional conditions;
 2. Run the segmentation algorithm for each finer selection event on the line;
 3. Identify the truncation region at which the same result is obtained;
 4. Combine the probability mass of multiply truncated Normal distributions;

Parametric programming approach

- ▶ The selection event of CNN segmentation algorithm is complicated:

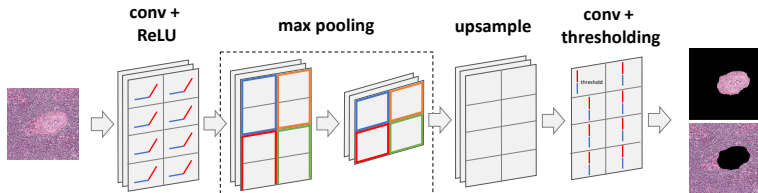


- ▶ Our idea is to consider solving a sequence of segmentation problems for a parametrized data in the direction of the test-statistic:
 1. Consider multiple finer selection events with additional conditions;
 2. Run the segmentation algorithm for each finer selection event on the line;
 3. Identify the truncation region at which the same result is obtained;
 4. Combine the probability mass of multiply truncated Normal distributions;

Additional conditioning by finer events regarding basic components of CNN

- By additionally conditioning on finer selection events regarding the basic components of CNN, the selection event is characterized as a union of polhedra.

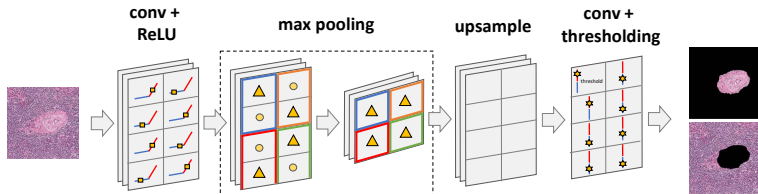
Component	Operations
Convolution	linear
ReLU transfer function	piecewise-linear
Max-pooling	comparison
Upsampling	linear
Thresholding	comparison



Additional conditioning by finer events regarding basic components of CNN

- By additionally conditioning on finer selection events regarding the basic components of CNN, the selection event is characterized as a union of polyhedra.

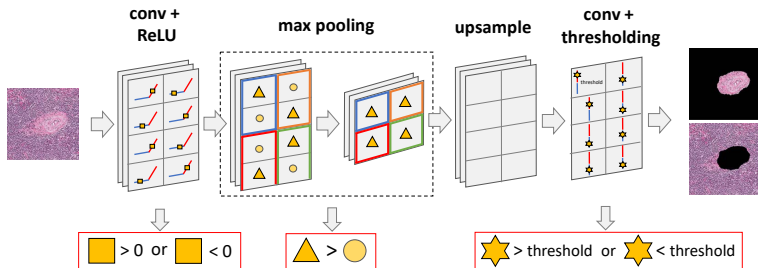
Component	Operations
Convolution	linear
ReLU transfer function	piecewise-linear
Max-pooling	comparison
Upsampling	linear
Thresholding	comparison



Additional conditioning by finer events regarding basic components of CNN

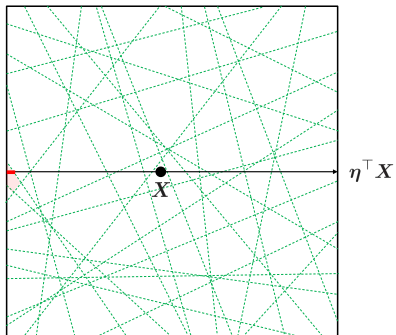
- By additionally conditioning on finer selection events regarding the basic components of CNN, the selection event is characterized as a union of polhedra.

Component	Operations
Convolution	linear
ReLU transfer function	piecewise-linear
Max-pooling	comparison
Upsampling	linear
Thresholding	comparison



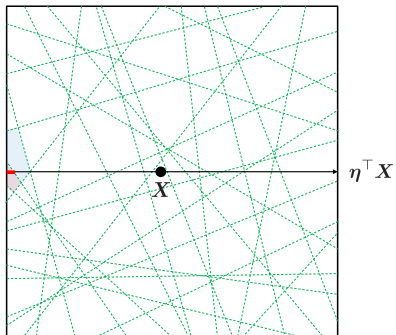
A summary of the divide-and-conquer approach

- ▶ Parametric programming approach looks like:



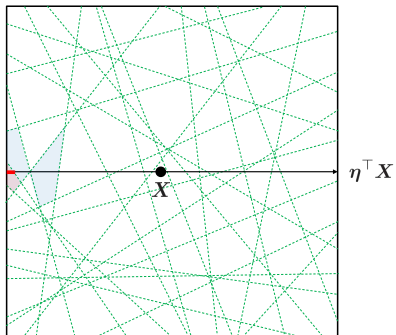
A summary of the divide-and-conquer approach

- ▶ Parametric programming approach looks like:



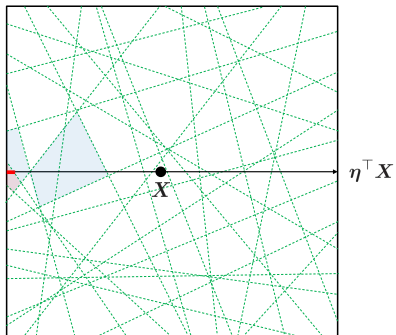
A summary of the divide-and-conquer approach

- ▶ Parametric programming approach looks like:



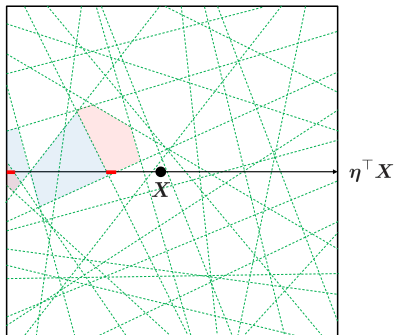
A summary of the divide-and-conquer approach

- ▶ Parametric programming approach looks like:



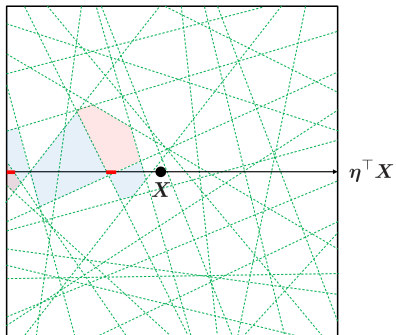
A summary of the divide-and-conquer approach

- ▶ Parametric programming approach looks like:



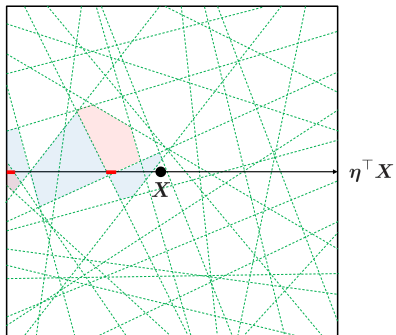
A summary of the divide-and-conquer approach

- ▶ Parametric programming approach looks like:



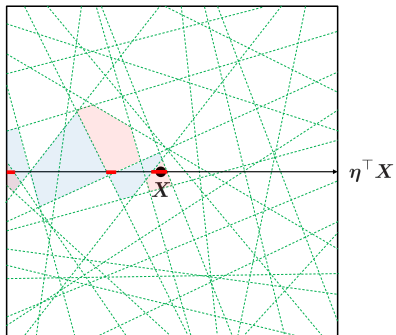
A summary of the divide-and-conquer approach

- ▶ Parametric programming approach looks like:



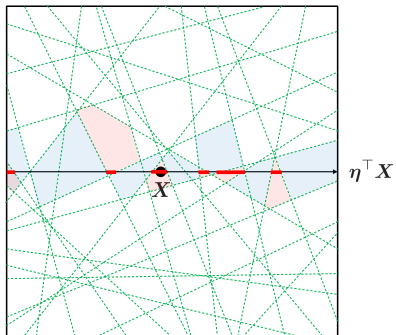
A summary of the divide-and-conquer approach

- ▶ Parametric programming approach looks like:



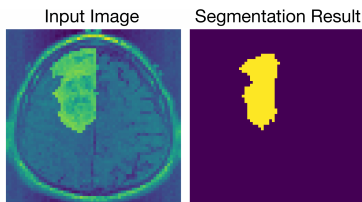
A summary of the divide-and-conquer approach

- ▶ Parametric programming approach looks like:



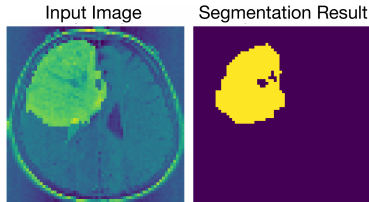
Examples of brain tumor image segmentation

- ▶ Positive cases (with true brain tumors)



naive p : 0.000
true pos.

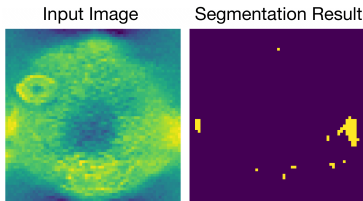
selective p : 0.000
true pos.



naive p : 0.000
true pos.

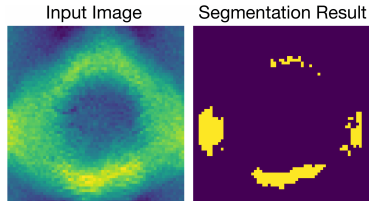
selective p : 0.000
true pos.

- ▶ Negative cases (without true brain tumors)



naive p : 0.000
false pos.

selective p : 0.670
true neg.



naive p : 0.000
false pos.

selective p : 0.451
true neg.

Summary

- ▶ Knowledge obtained by AI is a data-driven hypothesis, which cannot be properly evaluated by traditional statistical inference.
- ▶ Conditional SI is a promising approach for exact inference for data-driven hypotheses.
- ▶ The main technical challenge of conditional SI is how to characterize the selection event.
- ▶ Polyhedral lemma enables us to handle selection event represented as a polyhedron.
- ▶ Our parametric programming approach can be used for more complicate selection event such as DNN-driven hypotheses.
- ▶ We applied this parametric programming approach to several other problems such change point detection, outlier detection, clustering etc.

References

- ▶ J. Lee et al. Exact post-selection inference, with application to the lasso. *Annals of Statistics* 2016.
- ▶ R. Tibshirani et al. Exact post-selection inference for sequential regression procedures. *Journal of American Statistical Association* 2016.
- ▶ W. Fithian et al. Optimal inference after model selection. *arXiv* 2014.
- ▶ J. Taylor et al. Statistical learning and selective inference. *PNAS* 2015.
- ▶ X. Tian, J. Taylor. Asymptotics of selective inference. *arXiv* 2015.
- ▶ S. Suzumura et al. Selective inference for sparse higher-order interaction models. *ICML2017*.
- ▶ K. Tanizaki et al. Computing valid p-values for image segmentation by selective inference. *CVPR2020*.
- ▶ VNL. Duy et al. Computing valid p-value for optimal changepoint by selective inference using dynamic programming. *NeurIPS2020*.
- ▶ VNL. Duy et al. Parametric programming approach for more powerful and general Lasso selective inference. *AISTATS2021*.
- ▶ VNL. Duy et al. Quantifying Statistical Significance of Neural Network-based Image Segmentation by Selective Inference. *arXiv* 2020.
- ▶ K. Sugiyam et al. More powerful and general selective inference for stepwise feature selection using homotopy method. *ICML2021*.
- ▶ D. Das et al. Fast and more powerful selective inference for sparse high-order interaction model. *AAAI2022*.