

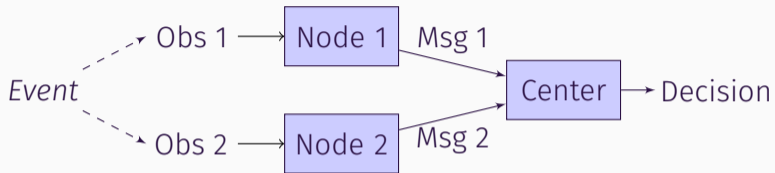
AN INFORMATION THEORETIC FRAMEWORK FOR DISTRIBUTED LEARNING ALGORITHMS

Xiangxiang Xu Shao-Lun Huang

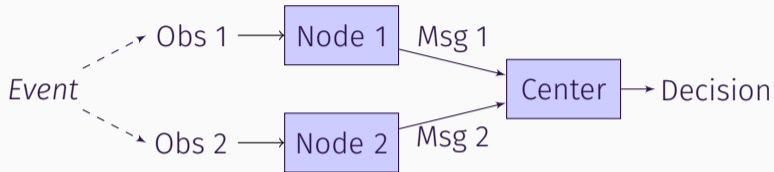
ISIT 2021

Tsinghua-Berkeley Shenzhen Institute, Tsinghua University

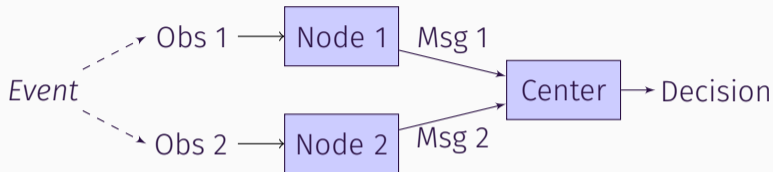
BACKGROUND: DISTRIBUTED LEARNING



BACKGROUND: DISTRIBUTED LEARNING



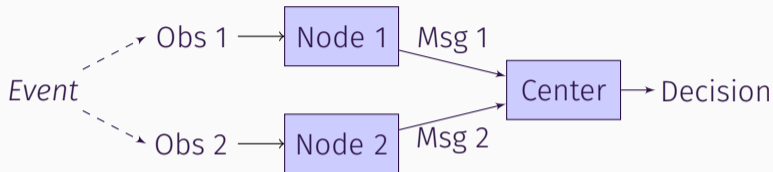
Related: Federated Learning, Sensor Fusion, ...



Related: Federated Learning, Sensor Fusion, ...

Information-Theoretic Characterizations

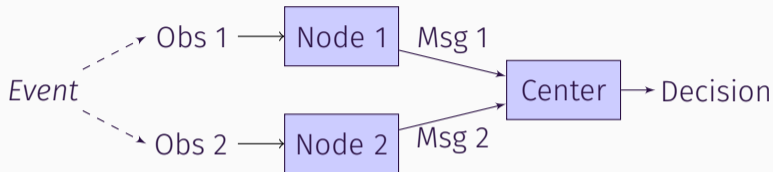
- Multiterminal Statistical Inference/Estimation (Ahlsvede, Csiszár, Han, ...)



Related: Federated Learning, Sensor Fusion, ...

Information-Theoretic Characterizations

- Multiterminal Statistical Inference/Estimation (Ahlsvede, Csiszár, Han, ...)
- CEO problem (Berger et al.)

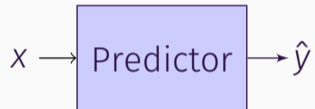


Related: Federated Learning, Sensor Fusion, ...

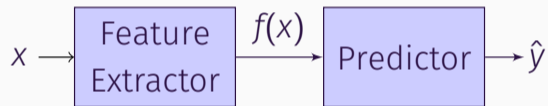
Information-Theoretic Characterizations

- Multiterminal Statistical Inference/Estimation (Ahlsvede, Csiszár, Han, ...)
- CEO problem (Berger et al.)
- Multi-user information theory ...

Use data X to predict label Y

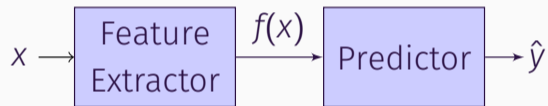


Use data X to predict label Y



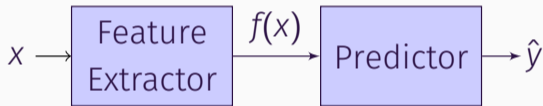
- Extract low-dim feature for prediction

Use data X to predict label Y

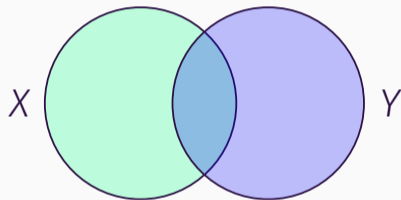


- Extract low-dim feature for prediction
- f should be “informative” about Y

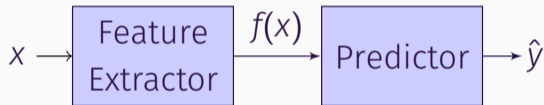
Use data X to predict label Y



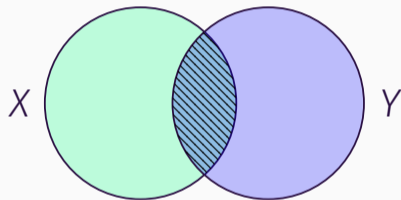
- Extract low-dim feature for prediction
- f should be “informative” about Y



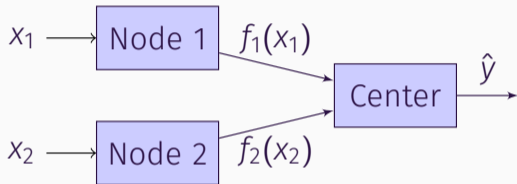
Use data X to predict label Y



- Extract low-dim feature for prediction
- f should be “informative” about Y

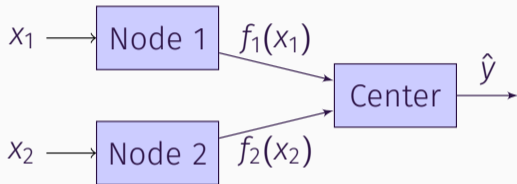


Prediction by Distributed Nodes



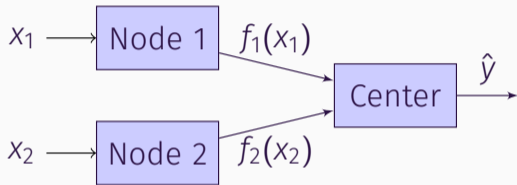
- Feature extraction required due to communication constraints

Prediction by Distributed Nodes

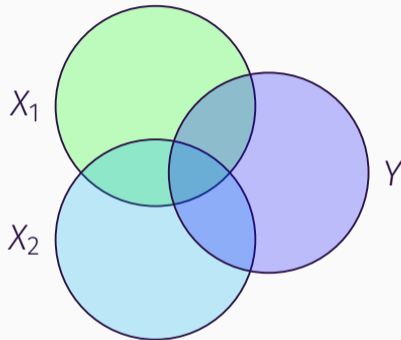


- Feature extraction required due to communication constraints

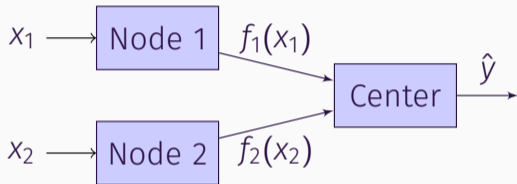
Prediction by Distributed Nodes



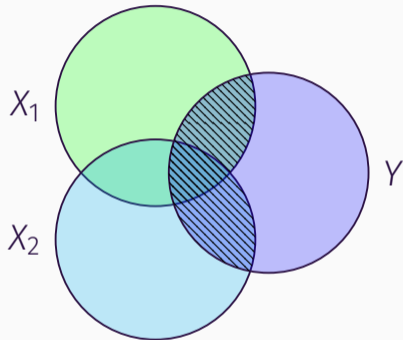
- Feature extraction required due to communication constraints
- (f_1, f_2) should be “informative” for prediction



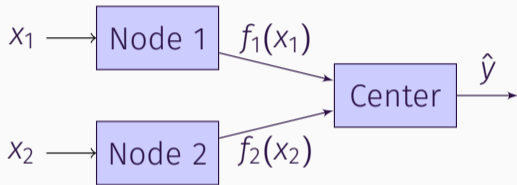
Prediction by Distributed Nodes



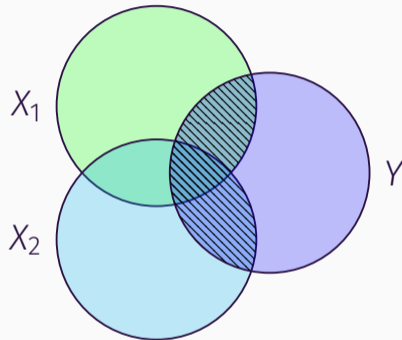
- Feature extraction required due to communication constraints
- (f_1, f_2) should be “informative” for prediction



Prediction by Distributed Nodes

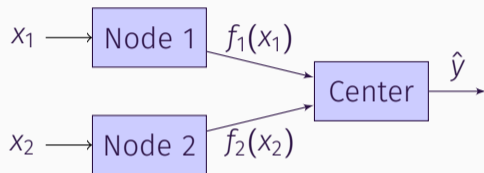


- Feature extraction required due to communication constraints
- (f_1, f_2) should be “informative” for prediction

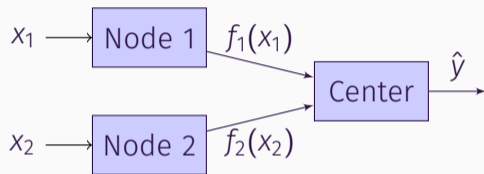


- Overlap between distributed nodes

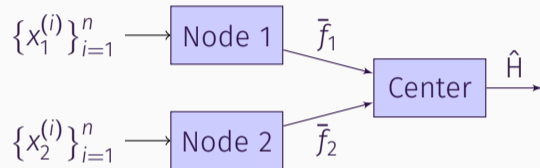
Distributed Classification



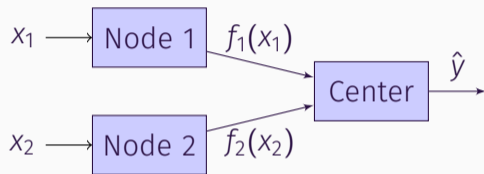
Distributed Classification



Distributed Hypothesis Testing

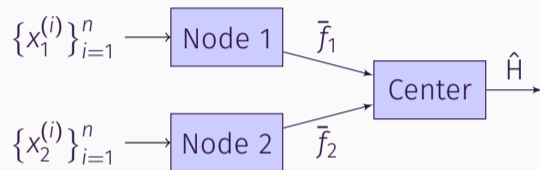


Distributed Classification



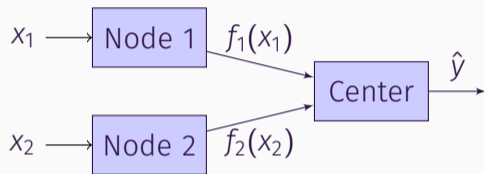
- Binary label $Y \in \{0, 1\}$

Distributed Hypothesis Testing



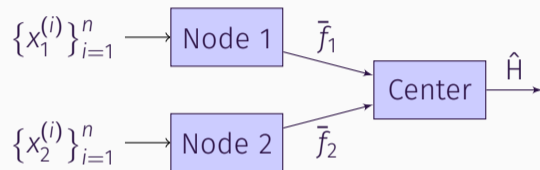
- Binary hypothesis $H \in \{0, 1\}$

Distributed Classification



- Binary label $Y \in \{0, 1\}$
- Sample pairs $\sim P_{\underline{X}|Y=y}$

Distributed Hypothesis Testing



- Binary hypothesis $H \in \{0, 1\}$
- Sample pairs $\stackrel{\text{i.i.d.}}{\sim} P_{\underline{X}}^{(H)}$

$$\underline{X} \triangleq \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$$

HYPOTHESIS TESTING WITH DISTRIBUTED NODES



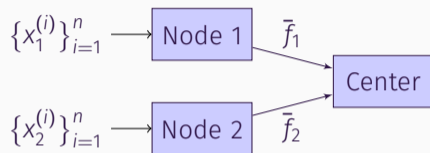
- (X_1, X_2) samples i.i.d. drawn from $P_{\underline{X}}^{(H)}$

HYPOTHESIS TESTING WITH DISTRIBUTED NODES



- (X_1, X_2) samples i.i.d. drawn from $P_{\underline{X}}^{(H)}$
- Node 1 observes only X_1 samples
- Node 2 observes only X_2 samples

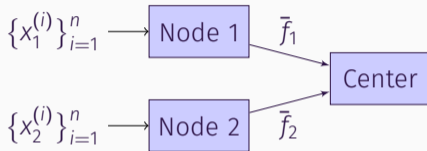
HYPOTHESIS TESTING WITH DISTRIBUTED NODES



- (X_1, X_2) samples i.i.d. drawn from $P_{\underline{X}}^{(H)}$
- Node 1 observes only X_1 samples
- Node 2 observes only X_2 samples
- Each node k sends center the statistic

$$\bar{f}_k \triangleq \frac{f_k(x_k^{(1)}) + \cdots + f_k(x_k^{(n)})}{n}$$

HYPOTHESIS TESTING WITH DISTRIBUTED NODES

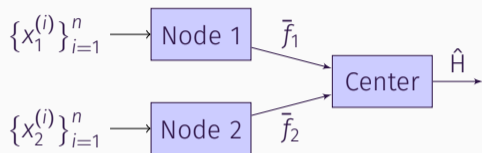


- (X_1, X_2) samples i.i.d. drawn from $P_{\underline{X}}^{(H)}$
- Node 1 observes only X_1 samples
- Node 2 observes only X_2 samples
- Each node k sends center the statistic

$$\bar{f}_k \triangleq \frac{f_k(x_k^{(1)}) + \cdots + f_k(x_k^{(n)})}{n}$$

- *zero-rate* communication

HYPOTHESIS TESTING WITH DISTRIBUTED NODES



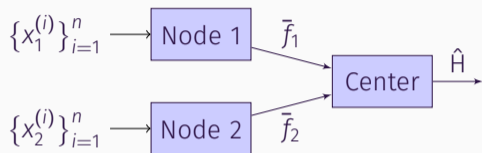
- (X_1, X_2) samples i.i.d. drawn from $P_{\underline{X}}^{(H)}$
- Node 1 observes only X_1 samples
- Node 2 observes only X_2 samples
- Each node k sends center the statistic

$$\bar{f}_k \triangleq \frac{f_k(x_k^{(1)}) + \cdots + f_k(x_k^{(n)})}{n}$$

- The center decides \hat{H} by MAP:

$$\frac{\mathbb{P}\{\bar{f}_1, \bar{f}_2 | H = 1\}}{\mathbb{P}\{\bar{f}_1, \bar{f}_2 | H = 0\}} \stackrel{\hat{H}=1}{\underset{\hat{H}=0}{\geq}} \frac{P_H(0)}{P_H(1)}$$

HYPOTHESIS TESTING WITH DISTRIBUTED NODES



- (X_1, X_2) samples i.i.d. drawn from $P_{\underline{X}}^{(H)}$
- Node 1 observes only X_1 samples
- Node 2 observes only X_2 samples
- Each node k sends center the statistic

Error exponent E of the decision:

$$\mathbb{P} \{ \hat{H} \neq H \} \doteq \exp(-nE)$$

Design Goal: maximize E

$$\bar{f}_k \triangleq \frac{f_k(x_k^{(1)}) + \cdots + f_k(x_k^{(n)})}{n}$$

- The center decides \hat{H} by MAP:

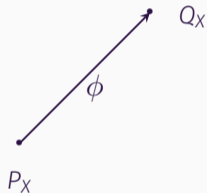
$$\frac{\mathbb{P} \{ \bar{f}_1, \bar{f}_2 | H = 1 \}}{\mathbb{P} \{ \bar{f}_1, \bar{f}_2 | H = 0 \}} \stackrel{\hat{H}=1}{\underset{\hat{H}=0}{\geq}} \frac{P_H(0)}{P_H(1)}$$

Local Geometric Structure

• Q_x

•
 P_x

Local Geometric Structure

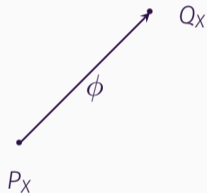


Distribution
 Q_X

\leftrightarrow

Information Vector
 $\phi(x) \triangleq \frac{Q_X(x) - P_X(x)}{\sqrt{P_X(x)}}$

Local Geometric Structure



Distribution
 Q_X

\leftrightarrow

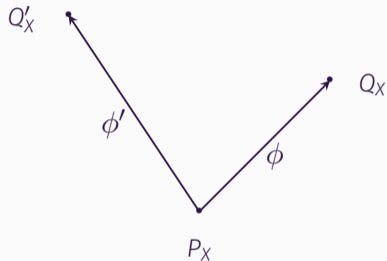
Information Vector
 $\phi(x) \triangleq \frac{Q_X(x) - P_X(x)}{\sqrt{P_X(x)}}$

\leftrightarrow

Feature Function
 $f(x) \triangleq \frac{\phi(x)}{\sqrt{P_X(x)}}$

Local Geometric Structure

$$\bullet Q_X \leftrightarrow \phi, Q'_X \leftrightarrow \phi'$$



Distribution
 Q_X

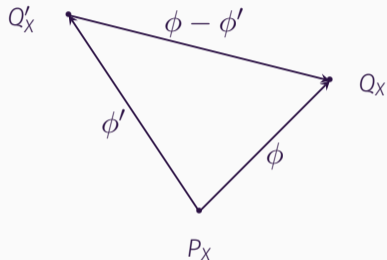
\leftrightarrow

Information Vector
 $\phi(x) \triangleq \frac{Q_X(x) - P_X(x)}{\sqrt{P_X(x)}}$

\leftrightarrow

Feature Function
 $f(x) \triangleq \frac{\phi(x)}{\sqrt{P_X(x)}}$

Local Geometric Structure



- $Q_X \leftrightarrow \phi, Q'_X \leftrightarrow \phi'$

- LLR (log-likelihood ratio):

$$\log \frac{Q'_X(x)}{Q_X(x)} \leftrightarrow \phi - \phi'$$

Distribution

Q_X

\leftrightarrow

Information Vector

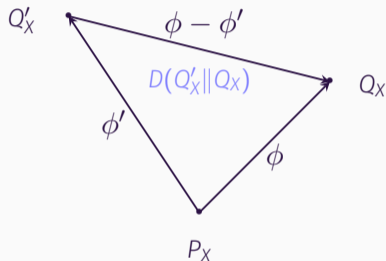
$$\phi(x) \triangleq \frac{Q_X(x) - P_X(x)}{\sqrt{P_X(x)}}$$

\leftrightarrow

Feature Function

$$f(x) \triangleq \frac{\phi(x)}{\sqrt{P_X(x)}}$$

Local Geometric Structure



- $Q_X \leftrightarrow \phi, Q'_X \leftrightarrow \phi'$
- LLR (log-likelihood ratio):

$$\log \frac{Q'_X(x)}{Q_X(x)} \leftrightarrow \phi - \phi'$$

- K-L divergence:

$$D(Q'_X || Q_X) = \frac{1}{2} \|\phi - \phi'\|^2$$

Distribution

Q_X

\leftrightarrow

Information Vector

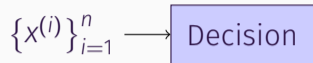
$$\phi(x) \triangleq \frac{Q_X(x) - P_X(x)}{\sqrt{P_X(x)}}$$

\leftrightarrow

Feature Function

$$f(x) \triangleq \frac{\phi(x)}{\sqrt{P_X(x)}}$$

HYPOTHESIS TESTING REVISITED: LOCAL GEOMETRIC ANALYSIS



- n samples $\{x^{(1)}, \dots, x^{(n)}\} \stackrel{\text{i.i.d.}}{\sim} P_X^{(H)}$

HYPOTHESIS TESTING REVISITED: LOCAL GEOMETRIC ANALYSIS



- n samples $\{x^{(1)}, \dots, x^{(n)}\} \stackrel{\text{i.i.d.}}{\sim} P_X^{(H)}$
- LLRT with opt. feature ℓ

$$\bar{\ell} \stackrel{\hat{H}=1}{\underset{\hat{H}=0}{\geq}} 0, \quad \ell(x) \triangleq \log \frac{P_X^{(1)}(x)}{P_X^{(0)}(x)}$$

HYPOTHESIS TESTING REVISITED: LOCAL GEOMETRIC ANALYSIS



- n samples $\{x^{(1)}, \dots, x^{(n)}\} \stackrel{\text{i.i.d.}}{\sim} P_X^{(H)}$
- LLRT with opt. feature ℓ

$$\bar{\ell} \underset{\hat{H}=0}{\overset{\hat{H}=1}{\geq}} 0, \quad \ell(x) \triangleq \log \frac{P_X^{(1)}(x)}{P_X^{(0)}(x)}$$

- Opt. error exponent E^* :
= Chernoff Information between $P_X^{(0)}$ and $P_X^{(1)}$

HYPOTHESIS TESTING REVISITED: LOCAL GEOMETRIC ANALYSIS

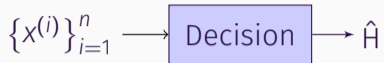


- n samples $\{x^{(1)}, \dots, x^{(n)}\} \stackrel{\text{i.i.d.}}{\sim} P_X^{(H)} \leftrightarrow \psi^{(H)}$
- LLRT with opt. feature ℓ

$$\bar{\ell} \stackrel{\hat{H}=1}{\underset{\hat{H}=0}{\geq}} 0, \quad \ell(x) \triangleq \log \frac{P_X^{(1)}(x)}{P_X^{(0)}(x)}$$

- Opt. error exponent E^* :
= Chernoff Information between $P_X^{(0)}$ and $P_X^{(1)}$

HYPOTHESIS TESTING REVISITED: LOCAL GEOMETRIC ANALYSIS



- n samples $\{x^{(1)}, \dots, x^{(n)}\} \stackrel{\text{i.i.d.}}{\sim} P_X^{(H)} \leftrightarrow \psi^{(H)}$
- LLRT with opt. feature $\ell \leftrightarrow \phi_\ell = \psi^{(1)} - \psi^{(0)}$

$$\bar{\ell} \stackrel{\hat{H}=1}{\underset{\hat{H}=0}{\geq}} 0, \quad \ell(x) \triangleq \log \frac{P_X^{(1)}(x)}{P_X^{(0)}(x)}$$

- Opt. error exponent E^* :
= Chernoff Information between $P_X^{(0)}$ and $P_X^{(1)}$

HYPOTHESIS TESTING REVISITED: LOCAL GEOMETRIC ANALYSIS



- n samples $\{x^{(1)}, \dots, x^{(n)}\} \stackrel{\text{i.i.d.}}{\sim} P_X^{(H)} \leftrightarrow \psi^{(H)}$
- LLRT with opt. feature $\ell \leftrightarrow \phi_\ell = \psi^{(1)} - \psi^{(0)}$

$$\bar{\ell} \stackrel{\hat{H}=1}{\underset{\hat{H}=0}{\geq}} 0, \quad \ell(x) \triangleq \log \frac{P_X^{(1)}(x)}{P_X^{(0)}(x)}$$

- Opt. error exponent E^* :
 - = Chernoff Information between $P_X^{(0)}$ and $P_X^{(1)}$
 - = $\|\phi_\ell\|^2/8$

HYPOTHESIS TESTING REVISITED: LOCAL GEOMETRIC ANALYSIS



- n samples $\{x^{(1)}, \dots, x^{(n)}\} \stackrel{\text{i.i.d.}}{\sim} P_X^{(H)} \leftrightarrow \psi^{(H)}$
- LLRT with opt. feature $\ell \leftrightarrow \phi_\ell = \psi^{(1)} - \psi^{(0)}$

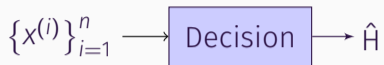
$$\bar{\ell} \stackrel{\hat{H}=1}{\underset{\hat{H}=0}{\gtrless}} 0, \quad \ell(x) \triangleq \log \frac{P_X^{(1)}(x)}{P_X^{(0)}(x)}$$

- Opt. error exponent E^* :
 - = Chernoff Information between $P_X^{(0)}$ and $P_X^{(1)}$
 - = $\|\phi_\ell\|^2/8$

Decision with feature $f \leftrightarrow \phi$

- Decision rule: $\bar{f} \stackrel{\hat{H}=1}{\underset{\hat{H}=0}{\gtrless}} 0$

HYPOTHESIS TESTING REVISITED: LOCAL GEOMETRIC ANALYSIS



- n samples $\{x^{(1)}, \dots, x^{(n)}\} \stackrel{\text{i.i.d.}}{\sim} P_X^{(H)} \leftrightarrow \psi^{(H)}$
- LLRT with opt. feature $\ell \leftrightarrow \phi_\ell = \psi^{(1)} - \psi^{(0)}$

$$\bar{\ell} \stackrel{\hat{H}=1}{\underset{\hat{H}=0}{\geq}} 0, \quad \ell(x) \triangleq \log \frac{P_X^{(1)}(x)}{P_X^{(0)}(x)}$$

- Opt. error exponent E^* :
 - = Chernoff Information between $P_X^{(0)}$ and $P_X^{(1)}$
 - = $\|\phi_\ell\|^2/8$

Decision with feature $f \leftrightarrow \phi$

- Decision rule: $\bar{f} \stackrel{\hat{H}=1}{\underset{\hat{H}=0}{\geq}} 0$
- Exponent $E = \frac{1}{8} \cdot \frac{\langle \phi, \phi_\ell \rangle^2}{\|\phi\|^2}$

HYPOTHESIS TESTING REVISITED: LOCAL GEOMETRIC ANALYSIS



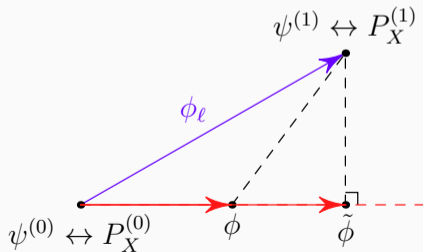
- n samples $\{x^{(1)}, \dots, x^{(n)}\} \stackrel{\text{i.i.d.}}{\sim} P_X^{(H)} \leftrightarrow \psi^{(H)}$
- LLRT with opt. feature $l \leftrightarrow \phi_l = \psi^{(1)} - \psi^{(0)}$

$$\bar{l} \stackrel{\hat{H}=1}{\underset{\hat{H}=0}{\geq}} 0, \quad l(x) \triangleq \log \frac{P_X^{(1)}(x)}{P_X^{(0)}(x)}$$

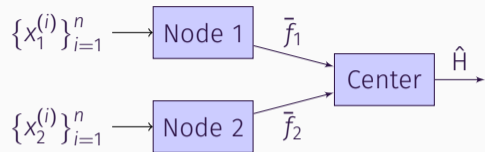
- Opt. error exponent E^* :
 = Chernoff Information between $P_X^{(0)}$ and $P_X^{(1)}$
 = $\|\phi_l\|^2/8$

Decision with feature $f \leftrightarrow \phi$

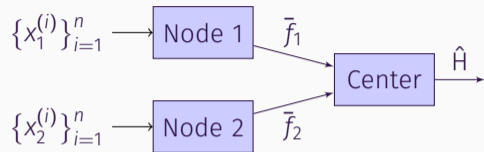
- Decision rule: $\bar{f} \stackrel{\hat{H}=1}{\underset{\hat{H}=0}{\geq}} 0$
- Exponent $E = \frac{1}{8} \cdot \frac{\langle \phi, \phi_l \rangle^2}{\|\phi\|^2} = \frac{1}{8} \|\tilde{\phi}\|^2$



DISTRIBUTED HYPOTHESIS TESTING: LOCAL GEOMETRIC ANALYSIS

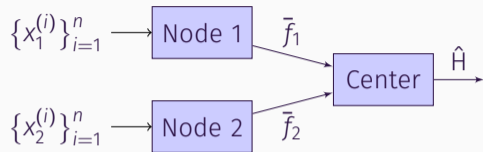


DISTRIBUTED HYPOTHESIS TESTING: LOCAL GEOMETRIC ANALYSIS



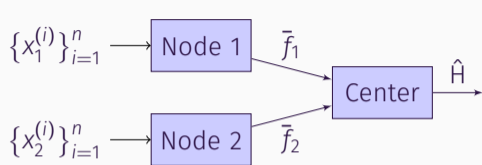
- f_1^*, f_2^* are both one-dim

DISTRIBUTED HYPOTHESIS TESTING: LOCAL GEOMETRIC ANALYSIS

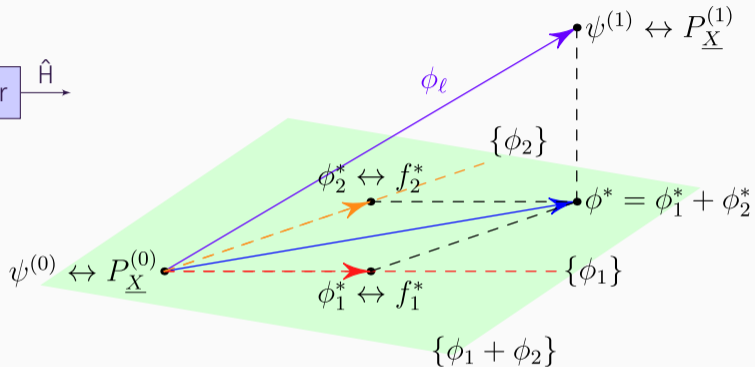


- f_1^*, f_2^* are both one-dim
- $f_1^* \leftrightarrow \phi_1^*, f_2^* \leftrightarrow \phi_2^*$

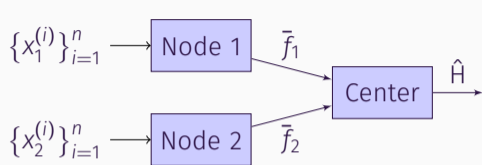
DISTRIBUTED HYPOTHESIS TESTING: LOCAL GEOMETRIC ANALYSIS



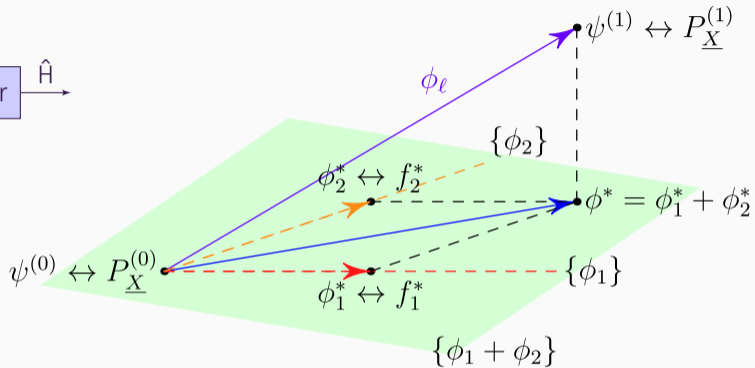
- f_1^*, f_2^* are both one-dim
- $f_1^* \leftrightarrow \phi_1^*, f_2^* \leftrightarrow \phi_2^*$
- $\phi^* = \phi_1^* + \phi_2^*$:
Proj. of ϕ_ℓ onto $\{\phi_1 + \phi_2\}$
- $\{\phi_1 + \phi_2\} \leftrightarrow \{f_1 + f_2\}$, features
of the form $f_1(x_1) + f_2(x_2)$



DISTRIBUTED HYPOTHESIS TESTING: LOCAL GEOMETRIC ANALYSIS

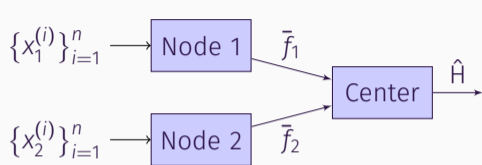


- f_1^*, f_2^* are both one-dim
- $f_1^* \leftrightarrow \phi_1^*, f_2^* \leftrightarrow \phi_2^*$
- $\phi^* = \phi_1^* + \phi_2^*$:
Proj. of ϕ_ℓ onto $\{\phi_1 + \phi_2\}$
- $\{\phi_1 + \phi_2\} \leftrightarrow \{f_1 + f_2\}$, features of the form $f_1(x_1) + f_2(x_2)$

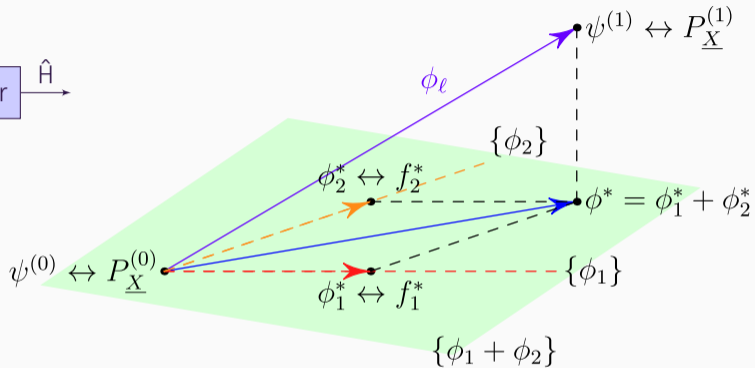


Decision based only on X_1 samples:
 \iff Proj. of ϕ_ℓ onto $\{\phi_1\}$

DISTRIBUTED HYPOTHESIS TESTING: LOCAL GEOMETRIC ANALYSIS



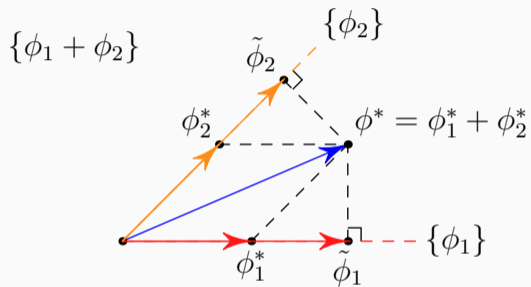
- f_1^*, f_2^* are both one-dim
- $f_1^* \leftrightarrow \phi_1^*, f_2^* \leftrightarrow \phi_2^*$
- $\phi^* = \phi_1^* + \phi_2^*$:
Proj. of ϕ_ℓ onto $\{\phi_1 + \phi_2\}$
- $\{\phi_1 + \phi_2\} \leftrightarrow \{f_1 + f_2\}$, features of the form $f_1(x_1) + f_2(x_2)$



Decision based only on X_1 samples:

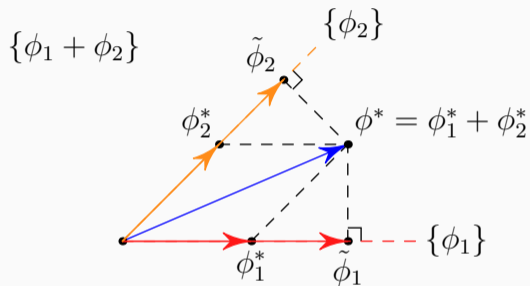
- \iff Proj. of ϕ_ℓ onto $\{\phi_1\}$
- \iff Proj. of ϕ^* onto $\{\phi_1\}$

INFORMATION DECOMPOSITION AMONG DISTRIBUTED NODES



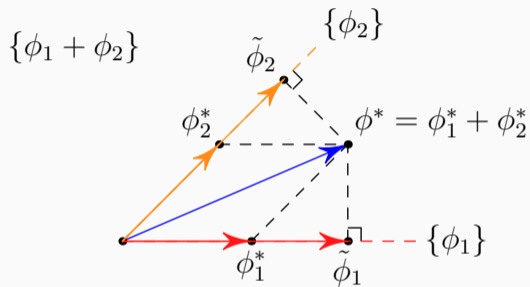
- $\{\phi_1\} \leftrightarrow \{\text{feature } f_1 \text{ of } x_1\}$
- $\{\phi_2\} \leftrightarrow \{\text{feature } f_2 \text{ of } x_2\}$
- $\{\phi_1 + \phi_2\} \leftrightarrow \{f_1 + f_2\}$

INFORMATION DECOMPOSITION AMONG DISTRIBUTED NODES



- $\{\phi_1\} \leftrightarrow \{\text{feature } f_1 \text{ of } x_1\}$
- $\{\phi_2\} \leftrightarrow \{\text{feature } f_2 \text{ of } x_2\}$
- $\{\phi_1 + \phi_2\} \leftrightarrow \{f_1 + f_2\}$
- $\tilde{\phi}_1$: proj. of ϕ^* onto $\{\phi_1\}$
 \leftrightarrow opt. feature in X_1 -based decision
- $\tilde{\phi}_2$: proj. of ϕ^* onto $\{\phi_2\}$
 \leftrightarrow opt. feature in X_2 -based decision
- Performances: $\|\tilde{\phi}_1\|^2, \|\tilde{\phi}_2\|^2, \|\phi^*\|^2$

INFORMATION DECOMPOSITION AMONG DISTRIBUTED NODES

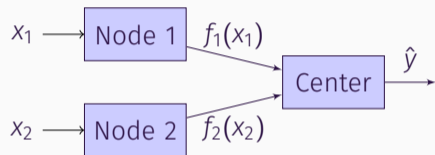


- $\{\phi_1\} \leftrightarrow \{\text{feature } f_1 \text{ of } x_1\}$
- $\{\phi_2\} \leftrightarrow \{\text{feature } f_2 \text{ of } x_2\}$
- $\{\phi_1 + \phi_2\} \leftrightarrow \{f_1 + f_2\}$
- $\tilde{\phi}_1$: proj. of ϕ^* onto $\{\phi_1\}$
 \leftrightarrow opt. feature in X_1 -based decision
- $\tilde{\phi}_2$: proj. of ϕ^* onto $\{\phi_2\}$
 \leftrightarrow opt. feature in X_2 -based decision
- Performances: $\|\tilde{\phi}_1\|^2, \|\tilde{\phi}_2\|^2, \|\phi^*\|^2$

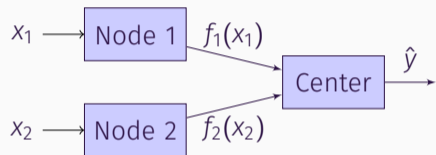
Due to the correlations between distributed nodes, orthogonal decomposition is in general not optimal.

CLASSIFICATION WITH DISTRIBUTED NODES

- Each node k transmits feature $f_k(x_k)$



CLASSIFICATION WITH DISTRIBUTED NODES

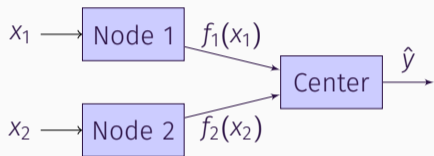


$$\underline{x} \triangleq \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \underline{f}(\underline{x}) \triangleq \begin{bmatrix} f_1(x_1) \\ f_2(x_2) \end{bmatrix}$$

- Each node k transmits feature $f_k(x_k)$
- The center learns a softmax classifier

$$\tilde{P}_{Y|\underline{X}}(y|\underline{x}) = \frac{\exp(\underline{f}^T(\underline{x})g(y))}{\sum_{y'} \exp(\underline{f}^T(\underline{x})g(y'))}$$

CLASSIFICATION WITH DISTRIBUTED NODES



$$\underline{x} \triangleq \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \underline{f}(\underline{x}) \triangleq \begin{bmatrix} f_1(x_1) \\ f_2(x_2) \end{bmatrix}$$

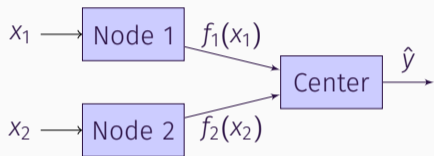
- Each node k transmits feature $f_k(x_k)$
- The center learns a softmax classifier

$$\tilde{P}_{Y|\underline{X}}(y|\underline{x}) = \frac{\exp(\underline{f}^T(\underline{x})g(y))}{\sum_{y'} \exp(\underline{f}^T(\underline{x})g(y))}$$

- Predict \hat{y} by MAP:

$$\hat{y} = \arg \max_{y'} \tilde{P}_{Y|\underline{X}}(y'|\underline{x})$$

CLASSIFICATION WITH DISTRIBUTED NODES



$$\underline{x} \triangleq \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \underline{f}(\underline{x}) \triangleq \begin{bmatrix} f_1(x_1) \\ f_2(x_2) \end{bmatrix}$$

- Each node k transmits feature $f_k(x_k)$
- The center learns a softmax classifier

$$\tilde{P}_{Y|\underline{X}}(y|\underline{x}) = \frac{\exp(\underline{f}^T(\underline{x})g(y))}{\sum_{y'} \exp(\underline{f}^T(\underline{x})g(y))}$$

- Predict \hat{y} by MAP:

$$\hat{y} = \arg \max_{y'} \tilde{P}_{Y|\underline{X}}(y'|\underline{x})$$

- **Goal:** minimize the logarithm loss

$$L \triangleq -\frac{1}{n} \sum_{i=1}^n \log \tilde{P}_{Y|\underline{X}}(y^{(i)}|\underline{x}^{(i)})$$

f_1^* and f_2^* are also optimal, under the correspondences

Hypothesis Testing	Classification
--------------------	----------------

$$P_{\underline{X}}^{(0)}$$

$$P_{\underline{X}|Y=0}$$

$$P_{\underline{X}}^{(1)}$$

$$P_{\underline{X}|Y=1}$$

$$P_H$$

$$P_Y$$

f_1^* and f_2^* are also optimal, under the correspondences

Hypothesis Testing	Classification
--------------------	----------------

$$P_X^{(0)}$$

$$P_{X|Y=0}$$

$$P_X^{(1)}$$

$$P_{X|Y=1}$$

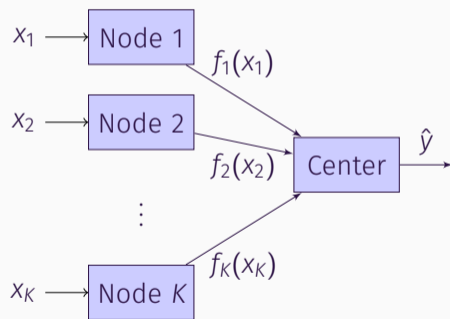
$$P_H$$

$$P_Y$$

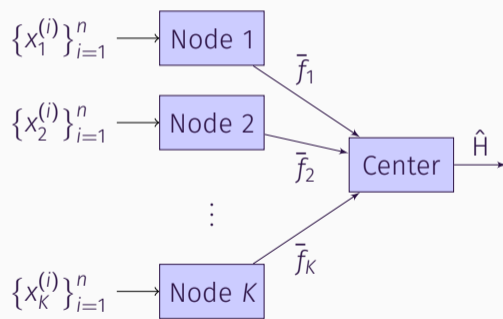
- Extension of [Huang et al., ISIT 2019]

GENERALIZATION: LEARNING WITH K DISTRIBUTED NODES

Distributed Classification

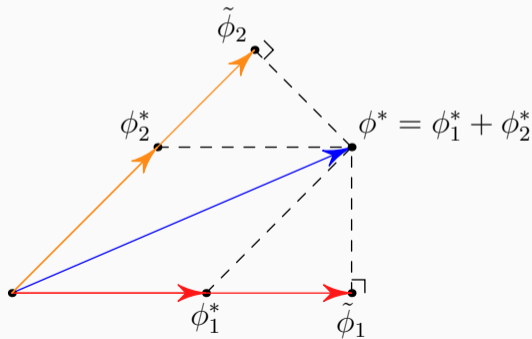
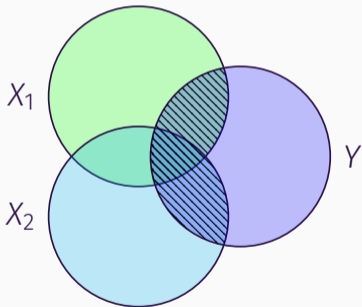


Distributed Hypothesis Testing



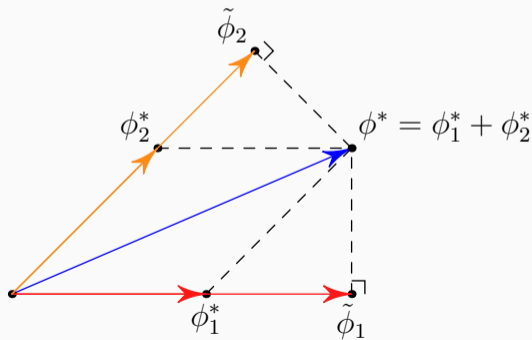
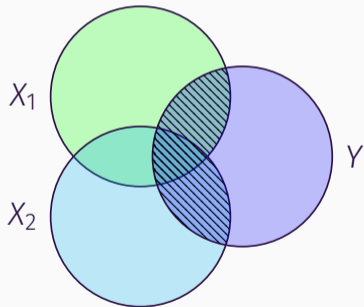
SUMMARY

- Information-theoretic framework



SUMMARY

- Information-theoretic framework
- Local geometric analyses: information decomposition structure



SUMMARY

- Information-theoretic framework
- Local geometric analyses: information decomposition structure
- Connection between information theory and machine learning

