

# Optimally Combining Outcomes to Improve Prediction

David Benkeser

On behalf of the HBGD community & HBGDki

[benkeser@berkeley.edu](mailto:benkeser@berkeley.edu)

# Acknowledgments

## **Collaborators:**

Mark van der Laan, Alan Hubbard, Ben Arnold, Jack Colford, Andrew Mertens, Oleg Sofyrgin, Jonathan French, Aryeh Stein, Shasha Jumbe

## **Funding:**

Bill and Melinda Gates Foundation OPP1147962

# Outline

Motivation

Defining target parameters

- ▶ Measuring accuracy of predictions
- ▶ Optimal predictor and weights

Estimation

- ▶ Super learning
- ▶ Estimating weights

Evaluating predictions

- ▶ Estimation and inference
- ▶ Variable importance

Simulation

Data Analysis

Conclusions and future directions

# Outline

## Motivation

Defining target parameters

- ▶ Measuring accuracy of predictions
- ▶ Optimal predictor and weights

Estimation

- ▶ Super learning
- ▶ Estimating weights

Evaluating predictions

- ▶ Estimation and inference
- ▶ Variable importance

Simulation

Data Analysis

Conclusions and future directions

# Motivation

The observed data are  $n$  i.i.d. copies of  $O = (X, Y)$ .

- ▶  $X = D$  covariates
- ▶  $Y = J$  standardized outcomes

Together  $Y$  represents an unmeasured outcome of interest.

- ▶ Inflammation in CVD
- ▶ Immune response in HIV
- ▶ Subject area tests in cognitive development

Researchers are interested in prediction of unmeasured outcome using  $X$ .

# Motivation

Can we use neonatal information to **predict neurocognitive outcomes** later in life?

- ▶ Early identification of at-risk children.

What **covariates** are important for prediction?

- ▶ Informs what information to collect to screen children.

# Motivation: PROBIT

The **Promotion of Breastfeeding Intervention Trial** enrolled pregnant mothers in 1996–97 (Kramer et al, 2001).

<b>Group name</b>	<b>Variables</b>
Breastfeeding	breastfeeding encouraged
Socioeconomic	household animals
Parental	age, height, weight, education, siblings employment status
Birth	gestational age, Apgar score
Growth	WAZ, HAZ, HCAZ (0,1,2,3,6,9,12 months)
Other	mother smoked during pregnancy, mother drank during pregnancy
WASI score (age 6)	Matrix, Block, Vocabulary, Similarities

# Motivation: CLHNS

The **Cebu Longitudinal Health and Nutrition Survey** enrolled pregnant mothers in 1983–84 (Feranil et al, 2008).

<b>Group name</b>	<b>Variables</b>
Health care	health care access, preventive health care
Household	child:adult ratio, child dependency ratio, crowding index, urban score
Socioeconomic	total income, SES
Water/sanitation	sanitation, access to clean water
Parental	mother age, father age, mother height, mother education (years), father education (years), marital status, mother age first child, parity
Growth	WAZ, HAZ (0,6,12,18,24 months)
Other	mother smoked during pregnancy, child's sex, gestational age at birth
Achievement tests (age 11)	Math, Cebuano, English



# Motivation

How to **combine test scores** to measure “neurocognition”?

Give **equal weight** to all scores?

- ▶ What if some scores are noisy or not related to covariates?

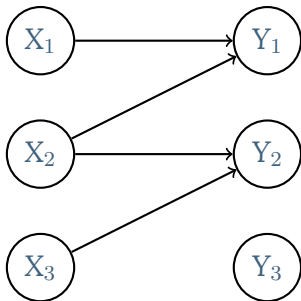
PCA or **factor analysis** to combine scores?

- ▶ Not related to scientific goal.

Use the combination that is **predicted most accurately**.

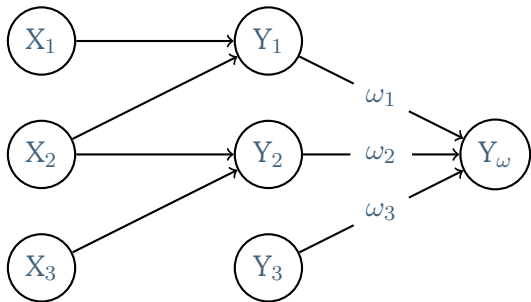
# Motivation

Consider this simple associative **directed acyclic graph**.



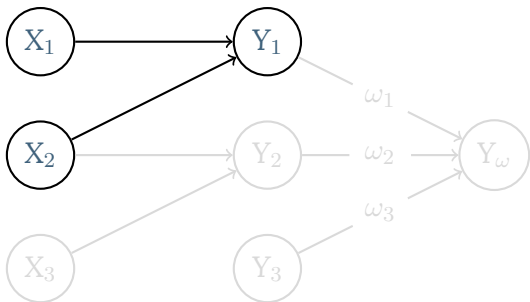
# Motivation

Let  $Y_\omega = \sum_{j=1}^J \omega_j Y_j$ , with  $\omega_j \geq 0$  for all  $j$  and  $\sum_{j=1}^J \omega_j = 1$ .



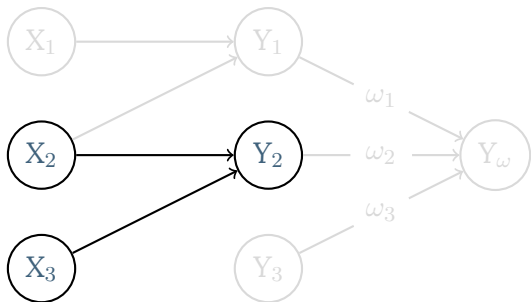
# Motivation

Predicting only  $Y_1$  ignores  $X_3$ 's association with  $Y$ .



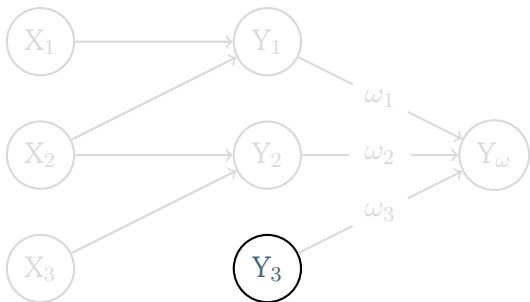
# Motivation

Predicting only  $Y_2$  ignores  $X_1$ 's association with  $Y$ .



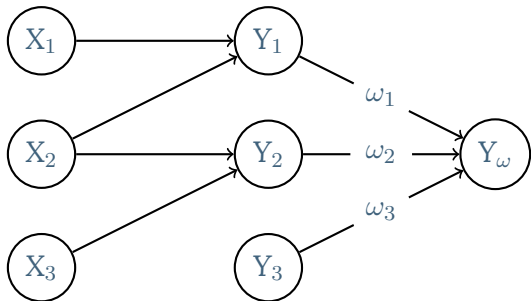
# Motivation

Predicting  $Y_3$  adds noise, wastes type-1 error.



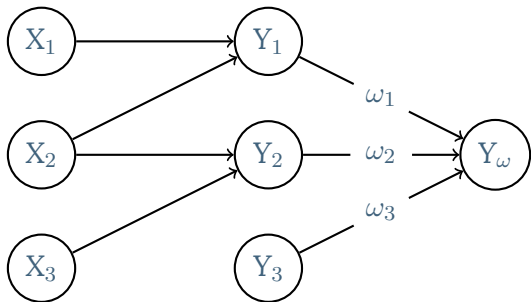
# Motivation

Predicting  $Y_\omega$  uses all of  $X$ .



# Motivation

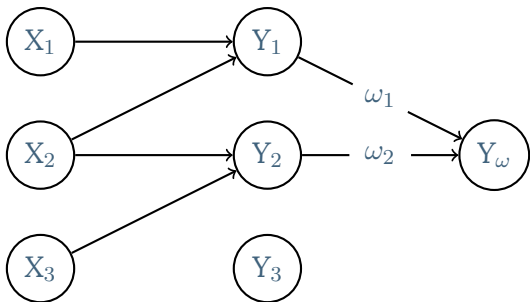
We could be clever in **choosing weights** if we knew the DAG.





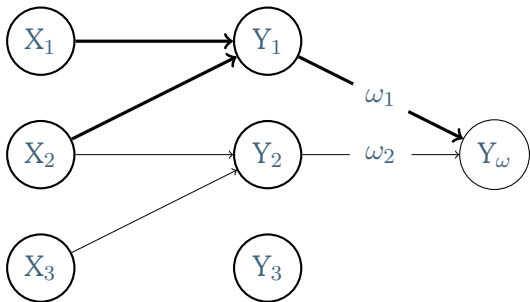
# Motivation

Outcome  $Y_3$  gets **no weight**.



# Motivation

Outcomes  $Y_1$  and  $Y_2$  get weight based on **accuracy** of predictions.



# Outline

Motivation

Defining target parameters

- ▶ Measuring accuracy of predictions
- ▶ Optimal predictor and weights

Estimation

- ▶ Super learning
- ▶ Estimating weights

Evaluating predictions

- ▶ Estimation and inference
- ▶ Variable importance

Simulation

Data Analysis

Conclusions and future directions

# Measuring accuracy

Suppose we are given  $\omega$  and  $\psi_\omega : \mathcal{X} \rightarrow \mathbb{R}$  to predict  $Y_\omega$ .

A **measure of accuracy** of  $\psi_\omega$  is MSE:

$$\mathcal{E}_0(\psi_\omega) = \mathbb{E}_0 [\{Y_\omega - \psi_\omega(\mathbf{X})\}^2] .$$

However, MSE depends on scale and variability of  $Y_\omega$ .

- ▶ Hard to compare across studies, outcomes.

# Measuring accuracy

To obtain a **scale-free measure**, we use a nonparametric version of  $R^2$ .

Let  $\mu_{0,\omega} = E_0(Y_\omega)$  be predictions made ignoring  $X$ .

$$R_{0,\omega}^2(\psi_\omega) = \frac{E_0 [\{Y_\omega - \mu_{0,\omega}\}^2] - E_0 [\{Y_\omega - \psi_\omega(X)\}^2]}{E_0 [\{Y_\omega - \mu_{0,\omega}\}^2]} .$$

# Measuring accuracy

To obtain a **scale-free measure**, we use a nonparametric version of  $R^2$ .

Let  $\mu_{0,\omega} = E_0(Y_\omega)$  be predictions made ignoring  $X$ .

$$R_{0,\omega}^2(\psi_\omega) = \frac{\overbrace{E_0 [\{Y_\omega - \mu_{0,\omega}\}^2]}^{\text{MSE of } \mu_{0,\omega}} - E_0 [\{Y_\omega - \psi_\omega(X)\}^2]}{E_0 [\{Y_\omega - \mu_{0,\omega}\}^2]} .$$

# Measuring accuracy

To obtain a **scale-free measure**, we use a nonparametric version of  $R^2$ .

Let  $\mu_{0,\omega} = E_0(Y_\omega)$  be predictions made ignoring  $X$ .

$$R_{0,\omega}^2(\psi_\omega) = \frac{\overbrace{E_0 [\{Y_\omega - \mu_{0,\omega}\}^2]}^{\text{MSE of } \mu_{0,\omega}} - \overbrace{E_0 [\{Y_\omega - \psi_\omega(X)\}^2]}^{\text{MSE of } \psi_\omega}}{E_0 [\{Y_\omega - \mu_{0,\omega}\}^2]} .$$

# Measuring accuracy

To obtain a **scale-free measure**, we use a nonparametric version of  $R^2$ .

Let  $\mu_{0,\omega} = E_0(Y_\omega)$  be predictions made ignoring  $X$ .

$$R_{0,\omega}^2(\psi_\omega) = \frac{\overbrace{E_0 [\{Y_\omega - \mu_{0,\omega}\}^2]}^{\text{MSE of } \mu_{0,\omega}} - \overbrace{E_0 [\{Y_\omega - \psi_\omega(X)\}^2]}^{\text{MSE of } \psi_\omega}}{\underbrace{E_0 [\{Y_\omega - \mu_{0,\omega}\}^2]}_{\text{MSE of } \mu_{0,\omega}}} .$$



# Predicting combined outcome

The criteria  $R_{0,\omega}^2$  provides a way to **compare** prediction functions for  $Y_\omega$ .

- ▶ Large  $R_{0,\omega}^2$  indicates **accurate** predictions.
- ▶  $R_{0,\omega}^2 = 1$  indicates **perfect** predictions.
- ▶  $R_{0,\omega}^2 < 0$  means predictions worse than  $\mu_{0,\omega}$ !

The **maximizer** over all  $\mathcal{X} \rightarrow \mathbb{R}$  is

$$\psi_{0,\omega}(\mathbf{X}) = E_0(Y_\omega | \mathbf{X}) .$$

This fact plays a key role in how we will construct a **prediction function**.

# Optimal weights

For any  $\omega$ , the function  $\psi_{0,\omega}$  gives the **most accurate** predictions of  $Y_\omega$ .

We also want **weights** that lead to most accurate predictions of combined outcome.

Formally, we define

$$\omega_0 = \operatorname{argmax}_\omega R_{0,\omega}^2(\psi_{0,\omega}) .$$

The **statistical goal** is to estimate  $\psi_{0,\omega_0}$  and  $\omega_0$ .

# Caveats

The sense in which this combination is **optimal** is strictly related to prediction.

- ▶ Not how well combined outcome measures “latent” outcome.

The optimal weights could give **zero weight** to some outcomes.

- ▶ These outcomes may still be important!

The procedure is best viewed as an **exploratory** analysis.

- ▶ However, it can be fully pre-specified!

# Outline

Motivation

Defining target parameters

- ▶ Measuring accuracy of predictions
- ▶ Optimal predictor and weights

Estimation

- ▶ Super learning
- ▶ Estimating weights

Evaluating predictions

- ▶ Estimation and inference
- ▶ Variable importance

Simulation

Data Analysis

Conclusions and future directions

# Estimation

The statistical goal is to estimate  $\psi_{0,\omega_0}$  and  $\omega_0$ .

At first glance, it looks like a **difficult optimization** problem,

$$\begin{aligned}\psi_{0,\omega_0} &= \operatorname{argmax}_{\psi} R_{0,\omega_0}^2(\psi) \\ \omega_0 &= \operatorname{argmax}_{\omega} R_{0,\omega}^2(\psi_{0,\omega})\end{aligned}$$

The problem is made easier by recognizing

$$\psi_{0,\omega}(X) = E_0(Y_\omega | X) = E_0\left(\sum_{j=1}^J \omega_j Y_j \mid X\right) = \sum_{j=1}^J \omega_j E_0(Y_j | X) .$$

# Estimation

For any  $\omega$ ,  $\psi_{0,\omega}$  is weighted sum of **conditional means**.

This allows the **optimization** to be **split up**:

1. Estimate  $E_0(Y_j | X)$  for  $j = 1, \dots, J$ .
2. Combine estimates  $\sum_{j=1}^J \omega_j \hat{E}(Y_j | X)$ .
3. Optimize over weights using estimated prediction function.

However, we must take care to **avoid overfitting!**

# Outline

Motivation

Defining target parameters

- ▶ Measuring accuracy of predictions
- ▶ Optimal predictor and weights

Estimation

- ▶ Super learning
- ▶ Estimating weights

Evaluating predictions

- ▶ Estimation and inference
- ▶ Variable importance

Simulation

Data Analysis

Conclusions and future directions

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?



# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines (with different degrees?)

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines (with different degrees?)
- ▶ Penalized linear regression

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines (with different degrees?)
- ▶ Penalized linear regression, with different penalties?

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines (with different degrees?)
- ▶ Penalized linear regression, with different penalties?
- ▶ Random forests



# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines (with different degrees?)
- ▶ Penalized linear regression, with different penalties?
- ▶ Random forests, with different tuning parameters?

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines (with different degrees?)
- ▶ Penalized linear regression, with different penalties?
- ▶ Random forests, with different tuning parameters?
- ▶ Gradient boosting?

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines (with different degrees?)
- ▶ Penalized linear regression, with different penalties?
- ▶ Random forests, with different tuning parameters?
- ▶ Gradient boosting? Support vector machines?

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines (with different degrees?)
- ▶ Penalized linear regression, with different penalties?
- ▶ Random forests, with different tuning parameters?
- ▶ Gradient boosting? Support vector machines? Deep neural networks?

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines (with different degrees?)
- ▶ Penalized linear regression, with different penalties?
- ▶ Random forests, with different tuning parameters?
- ▶ Gradient boosting? Support vector machines? Deep neural networks?
- ▶ **Highly adaptive lasso?**

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines (with different degrees?)
- ▶ Penalized linear regression, with different penalties?
- ▶ Random forests, with different tuning parameters?
- ▶ Gradient boosting? Support vector machines? Deep neural networks?
- ▶ **Highly adaptive lasso?**
- ▶ Variable selection?

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines (with different degrees?)
- ▶ Penalized linear regression, with different penalties?
- ▶ Random forests, with different tuning parameters?
- ▶ Gradient boosting? Support vector machines? Deep neural networks?
- ▶ **Highly adaptive lasso?**
- ▶ Variable selection?
- ▶ Ad infinitum...

# Predicting outcomes

How should we estimate  $E_0(Y_j | X)$  for a given  $j$ ?

- ▶ Linear regression, with interactions, and nonlinear terms, or splines (with different degrees?)
- ▶ Penalized linear regression, with different penalties?
- ▶ Random forests, with different tuning parameters?
- ▶ Gradient boosting? Support vector machines? Deep neural networks?
- ▶ **Highly adaptive lasso?**
- ▶ Variable selection?
- ▶ Ad infinitum...

The **best algorithm** for estimating depends on the (unknown) truth!

- ▶ Might be different for different outcomes.



# Predicting outcomes

We use  $\psi : \mathcal{S} \rightarrow \Psi$  to denote an **algorithm**.

- ▶  $\mathcal{S}$  is all subsets of  $\{1, \dots, n\}$ .
- ▶  $\Psi$  is set of  $\mathcal{X} \rightarrow \mathbb{R}$ .

Given a data set, an algorithm:

1. takes as input a **subset of observations**;
2. uses observations to create a **prediction function**;
3. returns prediction function.

We refer to this process as **training** an algorithm.

# Predicting each outcome

Say we have  $M$  algorithms that could be used to estimate  $E_0(Y_j | X)$ .

How can we **evaluate** these algorithms?

# Predicting each outcome

Say we have  $M$  algorithms that could be used to estimate  $E_0(Y_j | X)$ .

How can we **evaluate** these algorithms?

- ▶ Train algorithms on full data, see which has largest empirical  $R^2$ .

# Predicting each outcome

Say we have  $M$  algorithms that could be used to estimate  $E_0(Y_j | X)$ .

How can we **evaluate** these algorithms?

- ▶ Train algorithms on full data, see which has largest empirical  $R^2$ .
  - ▶ **Overfit!**

# Predicting each outcome

Say we have  $M$  algorithms that could be used to estimate  $E_0(Y_j | X)$ .

How can we **evaluate** these algorithms?

- ▶ Train algorithms on full data, see which has largest empirical  $R^2$ .
  - ▶ **Overfit!**
- ▶ Train algorithms on full data, collect more data to evaluate  $R^2$ .

# Predicting each outcome

Say we have  $M$  algorithms that could be used to estimate  $E_0(Y_j | X)$ .

How can we **evaluate** these algorithms?

- ▶ Train algorithms on full data, see which has largest empirical  $R^2$ .
  - ▶ **Overfit!**
- ▶ Train algorithms on full data, collect more data to evaluate  $R^2$ .
  - ▶ **Expensive!**

# Predicting each outcome

Say we have  $M$  algorithms that could be used to estimate  $E_0(Y_j | X)$ .

How can we **evaluate** these algorithms?

- ▶ Train algorithms on full data, see which has largest empirical  $R^2$ .
  - ▶ **Overfit!**
- ▶ Train algorithms on full data, collect more data to evaluate  $R^2$ .
  - ▶ **Expensive!**
- ▶ Cross validation!

# Cross validation

Consider randomly splitting the data into  $K$  different pieces.

$S_1$
$S_2$
$S_3$
$S_4$
$S_5$



# Cross validation

Define first training,  $T_1$ , and validation,  $V_1$ , sample.

$V_1$
$T_1$
$T_1$
$T_1$
$T_1$

# Cross validation

Train  $M$  algorithms using  $T_1$ .

- ▶  $\psi_{j,m}(T_1)$ , for  $m = 1, \dots, M$

Withhold validation sample  $V_1$  from training process.

- ▶ As though we did another experiment of size  $|V_1|$ !

Use validation sample to **estimate MSE** of each algorithm

$$\hat{\mathcal{E}}_{j,m,1}(\psi_{j,m}) = \frac{1}{|V_1|} \sum_{i \in V_1} \{Y_{j,i} - \psi_{j,m}(T_1)(X_i)\}^2 .$$

# Cross validation

Define second training,  $T_2$ , and validation,  $V_2$ , sample.

$T_2$
$V_2$
$T_2$
$T_2$
$T_2$

# Cross validation

Train  $M$  algorithms using  $T_2$ .

- ▶  $\psi_{j,m}(T_2)$ , for  $m = 1, \dots, M$

Withhold validation sample  $V_2$  from training process.

- ▶ As though we did another experiment of size  $|V_2|$ !

Use validation sample to **estimate MSE** of each algorithm

$$\hat{\mathcal{E}}_{j,m,2}(\psi_{j,m}) = \frac{1}{|V_2|} \sum_{i \in V_2} \{Y_{j,i} - \psi_{j,m}(T_2)(X_i)\}^2 .$$

# Cross validation

Continue until each split has been **validation** once.

$T_3$
$T_3$
$V_3$
$T_3$
$T_3$

# Cross validation

Continue until each split has been **validation** once.

$T_4$
$T_4$
$T_4$
$V_4$
$T_4$

# Cross validation

Continue until each split has been **validation** once.

$T_5$
$T_5$
$T_5$
$T_5$
$V_5$

# Cross-validation selector

The **cross-validated MSE** of algorithm  $m$  is

$$\hat{\mathcal{E}}_{j,m}(\psi_{j,m}) = \frac{1}{K} \sum_{k=1}^K \hat{\mathcal{E}}_{j,m,k}(\psi_{j,m}) .$$

We call the algorithm  $m^*$  with the lowest MSE the **cross-validation selector**.

We might use  $\psi_{j,m^*}(F_n)$  as estimate of  $E_0(Y_j | X)$ , where  $F_n = \{1, \dots, n\}$ .



# Ensemble estimator

What if  $\psi_{j,1}$  and  $\psi_{j,2}$  capture **different features**?

Using  $\psi_{j,SL} = 0.5\psi_{j,1} + 0.5\psi_{j,2}$  might be better than  $\psi_{j,1}$  and  $\psi_{j,2}$  alone.

More generally, consider an **ensemble prediction** function

$$\psi_{j,SL} = \sum_{m=1}^M \alpha_{j,m} \psi_{j,m} , \alpha_{j,m} \geq 0 \text{ for all } m , \text{ and } \sum_{m=1}^M \alpha_{j,m} = 1 .$$

Easy to find  $\alpha_j$  that minimizes **cross-validated MSE**.

# Super learner

Stacked regression originally proposed in Breiman, 1996.

Referred to as a **super learner** due to **oracle inequality** results (van der Laan and Dudoit, 2003).

- ▶ MSE of the super learner is asymptotically equivalent to the oracle estimator.
- ▶ Even when one considers many estimators.
- ▶ Often seen to have good finite-sample performance (van der Laan et al, 2007).

Proving new **oracle results** is an open area of research!

- ▶ Recent work on Big Data oracle inequalities.

# Outline

Motivation

Defining target parameters

- ▶ Measuring accuracy of predictions
- ▶ Optimal predictor and weights

Estimation

- ▶ Super learning
- ▶ Estimating weights

Evaluating predictions

- ▶ Estimation and inference
- ▶ Variable importance

Simulation

Data Analysis

Conclusions and future directions

# Statistical goal

Difficult **optimization** problem,

$$\begin{aligned}\psi_{0,\omega_0} &= \operatorname{argmax}_{\psi} R_{0,\omega_0}^2(\psi) \\ \omega_0 &= \operatorname{argmax}_{\omega} R_{0,\omega}^2(\psi_{0,\omega}) ,\end{aligned}$$

made easier because

$$\psi_{0,\omega} = E_0\left(\sum_{j=1}^J \omega_j Y_j \mid X\right) = \sum_{j=1}^J \omega_j E_0(Y_j \mid X) .$$

For any  $\omega$ , **combine super learners** to estimate  $\psi_{0,\omega}$ ,

$$\psi_{n,\omega} = \sum_{j=1}^J \omega_j \psi_{j,SL}(F_n) .$$

# Estimating optimal weights

How do we estimate  $\omega_0$ ?

# Estimating optimal weights

How do we estimate  $\omega_0$ ?

- ▶ Maximize empirical  $R^2$  over weights.

$$\omega_n = \operatorname{argmax}_{\omega} \left( 1 - \frac{\frac{1}{n} \sum_{i=1}^n \{Y_{\omega,i} - \psi_{n,\omega}(X_i)\}^2}{\frac{1}{n} \sum_{i=1}^n \{Y_{\omega,i} - \bar{Y}_{\omega}\}^2} \right)$$

# Estimating optimal weights

How do we estimate  $\omega_0$ ?

- ▶ Maximize empirical  $R^2$  over weights.

$$\omega_n = \operatorname{argmax}_{\omega} \left( 1 - \frac{\frac{1}{n} \sum_{i=1}^n \{Y_{\omega,i} - \psi_{n,\omega}(X_i)\}^2}{\frac{1}{n} \sum_{i=1}^n \{Y_{\omega,i} - \bar{Y}_{\omega}\}^2} \right)$$

- ▶ **overfit!**

# Estimating optimal weights

How do we estimate  $\omega_0$ ?

- ▶ Maximize empirical  $R^2$  over weights.

$$\omega_n = \operatorname{argmax}_{\omega} \left( 1 - \frac{\frac{1}{n} \sum_{i=1}^n \{Y_{\omega,i} - \psi_{n,\omega}(X_i)\}^2}{\frac{1}{n} \sum_{i=1}^n \{Y_{\omega,i} - \bar{Y}_{\omega}\}^2} \right)$$

- ▶ **overfit!**
- ▶ Collect more data, maximize  $R^2$  on new data.



# Estimating optimal weights

How do we estimate  $\omega_0$ ?

- ▶ Maximize empirical  $R^2$  over weights.

$$\omega_n = \operatorname{argmax}_{\omega} \left( 1 - \frac{\frac{1}{n} \sum_{i=1}^n \{Y_{\omega,i} - \psi_{n,\omega}(X_i)\}^2}{\frac{1}{n} \sum_{i=1}^n \{Y_{\omega,i} - \bar{Y}_{\omega}\}^2} \right)$$

- ▶ **overfit!**
- ▶ Collect more data, maximize  $R^2$  on new data.
  - ▶ **expensive!**

# Estimating optimal weights

How do we estimate  $\omega_0$ ?

- ▶ Maximize empirical  $R^2$  over weights.

$$\omega_n = \operatorname{argmax}_{\omega} \left( 1 - \frac{\frac{1}{n} \sum_{i=1}^n \{Y_{\omega,i} - \psi_{n,\omega}(X_i)\}^2}{\frac{1}{n} \sum_{i=1}^n \{Y_{\omega,i} - \bar{Y}_{\omega}\}^2} \right)$$

- ▶ **overfit!**
- ▶ Collect more data, maximize  $R^2$  on new data.
  - ▶ **expensive!**
- ▶ Cross validation!

# Nested cross validation

Define first training,  $T_1$ , and validation,  $V_1$ , sample.

$V_1$
$T_1$
$T_1$
$T_1$
$T_1$

# Cross validation

Train  $J$  super learners using  $T_1$ .

- ▶  $\psi_{j,SL}(T_1)$ , for  $j = 1, \dots, J$

For any  $\omega$ , we can construct combined super learner

$$\psi_{\omega,SL}(T_1) = \sum_{j=1}^J \omega_j \psi_{j,SL}(T_1) .$$

Withhold validation sample  $V_1$  from super learner fitting.

- ▶ As though we did another experiment of size  $|V_1|!$

For any  $\omega$ , use validation sample to estimate MSE

$$\hat{\mathcal{E}}_{\omega,1}(\psi_{\omega,SL}) = \frac{1}{|V_1|} \sum_{i \in V_1} \{Y_{\omega,i} - \psi_{\omega,SL}(T_1)(X_i)\}^2 .$$

# Nested cross validation

Define second training,  $T_2$ , and validation,  $V_2$ , sample.

$T_2$
$V_2$
$T_2$
$T_2$
$T_2$

# Nested cross validation

Train  $J$  super learners using  $T_2$ .

- ▶  $\psi_{j,SL}(T_2)$ , for  $j = 1, \dots, J$

For any  $\omega$ , we can construct combined super learner

$$\psi_{\omega,SL}(T_2) = \sum_{j=1}^J \omega_j \psi_{j,SL}(T_2) .$$

Withhold validation sample  $V_2$  from super learner fitting.

- ▶ As though we did another experiment of size  $|V_2|$ !

For any  $\omega$ , use validation sample to estimate MSE

$$\hat{\mathcal{E}}_{\omega,2}(\psi_{\omega,SL}) = \frac{1}{|V_2|} \sum_{i \in V_2} \{Y_{\omega,i} - \psi_{\omega,SL}(T_2)(X_i)\}^2 .$$

# Nested cross validation

Continue until each split has been used once.

$T_3$
$T_3$
$V_3$
$T_3$
$T_3$

# Nested cross validation

Continue until each split has been used once.

$T_4$
$T_4$
$T_4$
$V_4$
$T_4$



# Nested cross validation

Continue until each split has been used once.

$T_5$
$T_5$
$T_5$
$T_5$
$V_5$

# Estimating weights

For any  $\omega$ , we have cross-validated estimate of  $R^2$ ,

$$R_{n,\omega}^2(\psi_{\omega,SL}) = 1 - \frac{\frac{1}{K} \sum_{k=1}^K \hat{\mathcal{E}}_{\omega,k}(\psi_{\omega,SL})}{\frac{1}{n} \sum_{i=1}^n \{Y_{\omega,i} - \bar{Y}_{\omega}\}^2} .$$

Estimate of **optimal weights** is

$$\omega_n = \operatorname{argmax}_{\omega} R_{n,\omega}^2(\psi_{\omega,SL}) .$$

# Outline

Motivation

Defining target parameters

- ▶ Measuring accuracy of predictions
- ▶ Optimal predictor and weights

Estimation

- ▶ Super learning
- ▶ Estimating weights

Evaluating predictions

- ▶ Estimation and inference
- ▶ Variable importance

Simulation

Data Analysis

Conclusions and future directions

# Estimating predictive performance

How accurate is  $\psi_{n,\omega_n}$  in predicting  $Y_{\omega_n}$ ?

# Estimating predictive performance

How accurate is  $\psi_{n,\omega_n}$  in predicting  $Y_{\omega_n}$ ?

- ▶ Report  $R_{n,\omega_n}^2(\psi_{\omega_n,SL})$ , call it a day.

# Estimating predictive performance

How accurate is  $\psi_{n,\omega_n}$  in predicting  $Y_{\omega_n}$ ?

- ▶ Report  $R_{n,\omega_n}^2(\psi_{\omega_n,SL})$ , call it a day.
  - ▶ **overfit!**

# Estimating predictive performance

How accurate is  $\psi_{n,\omega_n}$  in predicting  $Y_{\omega_n}$ ?

- ▶ Report  $R_{n,\omega_n}^2(\psi_{\omega_n,SL})$ , call it a day.
  - ▶ **overfit!**
- ▶ Collect more data to evaluate predictions.

# Estimating predictive performance

How accurate is  $\psi_{n,\omega_n}$  in predicting  $Y_{\omega_n}$ ?

- ▶ Report  $R_{n,\omega_n}^2(\psi_{\omega_n,SL})$ , call it a day.
  - ▶ **overfit!**
- ▶ Collect more data to evaluate predictions.
  - ▶ **expensive!**



# Estimating predictive performance

How accurate is  $\psi_{n,\omega_n}$  in predicting  $Y_{\omega_n}$ ?

- ▶ Report  $R_{n,\omega_n}^2(\psi_{\omega_n,SL})$ , call it a day.
  - ▶ **overfit!**
- ▶ Collect more data to evaluate predictions.
  - ▶ **expensive!**
- ▶ More cross validation!!!

# Doubly nested cross validation

Pictures omitted for the sanity of audience.

# Estimating predictive performance

Cross-validate the **entire procedure** to estimate performance.

- ▶ Compute  $\omega_n$  and  $\psi_{n,\omega_n}$  in training sample
- ▶ Estimate  $R^2$  in validation sample
- ▶ Average over splits

Formally, we define  $\omega : \mathcal{S} \rightarrow \Omega$  and  $\psi_{\omega,SL} : \mathcal{S} \rightarrow \Psi$ .

The cross-validated  $R^2$  estimate is

$$R_n^2(\omega_n, \psi_{n,\omega_n}) = 1 - \frac{\sum_{k=1}^K \frac{1}{|V_k|} \sum_{i \in V_k} \{Y_{\omega(T_k),i} - \psi_{\omega(T_k),SL}(X_i)\}^2}{\sum_{k=1}^K \frac{1}{|V_k|} \sum_{i \in V_k} \{Y_{\omega(T_k),i} - \bar{Y}_{\omega(T_k)}\}^2}.$$

# Inference for predictive performance

In spite of the highly adaptive estimation procedure,

$$n^{1/2} \{R_n^2(\omega_n, \psi_{n,\omega_n}) - R_0^2(\omega_n, \psi_{n,\omega_n})\} \rightarrow \text{Normal}(0, \sigma^2).$$

A sufficient condition is that  $\omega_0$  is **unique**.

- ▶ Possible to relax this condition (Luedtke et al, 2016).

Variance derived via delta method for **influence functions**.

- ▶ Consistently estimated with closed-form estimator.

Variance estimates used to construct **closed-form confidence intervals** and **hypotheses tests**.

- ▶ Machine learning with inference!

# Variable importance

Define  $R_0^2(\omega_n^d, \psi_{n,\omega_n}^d)$  as the true  $R^2$  of the estimated super learner that leaves out  $X_d$ ,  $d = 1, \dots, D$ .

The “importance” of  $X_d$  could be quantified by

$$\Delta_{0n}^d = R_0^2(\omega_n, \psi_{n,\omega_n}) - R_0^2(\omega_n^d, \psi_{n,\omega_n}^d).$$

How much did  $X_d$  improve predictions of combined outcome?

- ▶ Similar to random forest variable importance measures

# Variable importance

Variable importance can be estimated as

$$\Delta_n^d = R_n^2(\omega_n, \psi_{n, \omega_n}) - R_n^2(\omega_n^d, \psi_{n, \omega_n}^d) .$$

We can still establish

$$n^{1/2}(\Delta_n^d - \Delta_{0n}^d) \rightarrow \text{Normal}(0, \sigma_d^2) .$$

Variance can again be derived using delta method for **influence functions**.

Did  $X_d$  **significantly improve** predictions of combined outcome?

# Outline

Motivation

Defining target parameters

- ▶ Measuring accuracy of predictions
- ▶ Optimal predictor and weights

Estimation

- ▶ Super learning
- ▶ Estimating weights

Evaluating predictions

- ▶ Estimation and inference
- ▶ Variable importance

Simulation

Data Analysis

Conclusions and future directions

# Simulation

## Covariates:

$X_1, \dots, X_6 \sim \text{Uniform}(0, 4)$ ,  $X_7, \dots, X_9 \sim \text{Bernoulli}(0.5)$

## Outcomes:

$$Y_1 = X_1 + 2X_2 + 4X_3 + X_7 + 2X_8 + 4X_9 + 2X_4 + \epsilon_1 ,$$

$$Y_2 = X_1 + 2X_2 + 4X_3 + X_7 + 2X_8 + 4X_9 + 2X_5 + \epsilon_2 , \text{ and}$$

$$Y_3 = X_1 + 2X_2 + 4X_3 + X_7 + 2X_8 + 4X_9 + 2X_6 + \epsilon_3 ,$$

with  $\epsilon_j \sim \text{Normal}(0, 5^2)$ ,  $j = 1, 2, 3$ .

## True parameters:

$$\omega_0 = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right), R_{0, \omega_0}^2 = 0.80, \Delta_0^2 = 0.12.$$

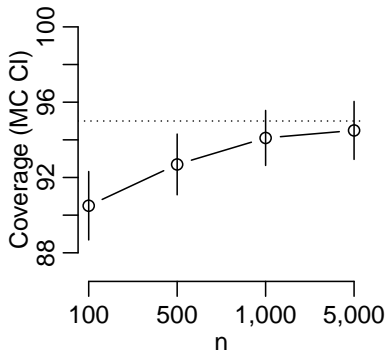
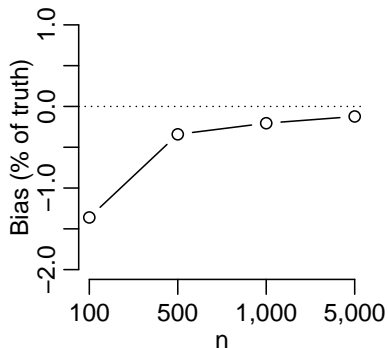
## Super learner:

intercept only, main terms, and stepwise linear model.



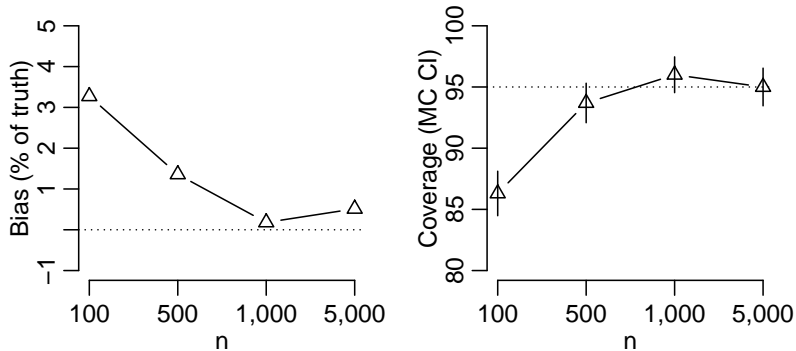
# Simulation

Bias and coverage for  $R_{n,\omega_n}^2$  (1000 replications).



# Simulation

Bias and coverage for  $\Delta_n^2$  (1000 replications).



# Outline

Motivation

Defining target parameters

- ▶ Measuring accuracy of predictions
- ▶ Optimal predictor and weights

Estimation

- ▶ Super learning
- ▶ Estimating weights

Evaluating predictions

- ▶ Estimation and inference
- ▶ Variable importance

Simulation

Data Analysis

Conclusions and future directions

# Data analysis

Can we use neonatal information to **predict neurocognitive outcomes** later in life?

- ▶ Early identification of at-risk children.

What covariates are important for making predictions?

- ▶ Informs what information to collect to screen children.

# Motivation: PROBIT

The **Promotion of Breastfeeding Intervention Trial** enrolled pregnant mothers in 1996–97.

<b>Group name</b>	<b>Variables</b>
Breastfeeding	breastfeeding encouraged
Socioeconomic	household animals
Parental	age, height, weight, education, siblings employment status
Birth	gestational age, Apgar score
Growth	WAZ, HAZ, HCAZ (0,1,2,3,6,9,12 months)
Other	mother smoked during pregnancy, mother drank during pregnancy
WASI score (age 6)	Matrix, Block, Vocabulary, Similarities

# Future directions

Effect estimation for discovery in **high dimensions**.

- ▶ Maximize effect instead of  $R^2$ ?
- ▶ Cross-validated TMLE (Zheng and van der Laan, 2010)

The method could be extended to **binary outcomes**, other performance metrics, and dependent data.

**Nonlinear combinations** of outcomes are also of interest.

- ▶ Alternating conditional expectations (Breiman and Friedman, 1985)

R packages:

**r2weight**

- ▶ Available on GitHub:

<https://github.com/benkeser/r2weight>

**SuperLearner** (Polley et al, 2016)

- ▶ Demonstration – <http://benkeser.github.io/sllecture/>

# References I

- [1] Michael S Kramer, Beverley Chalmers, Ellen D Hodnett, Zinaida Sevkovskaya, Irina Dzikovich, Stanley Shapiro, Jean-Paul Collet, Irina Vanilovich, Irina Mezen, Thierry Ducruet, et al. Promotion of breastfeeding intervention trial (PROBIT): a randomized trial in the Republic of Belarus. *Journal of the American Medical Association*, 285(4):413–420, 2001.
- [2] AB Feranil, SA Gultiano, and LS Adair. The Cebu Longitudinal Health and Nutrition Survey: Two Decades Later. *Asia-Pacific Population Journal*, 23(3), 2008.
- [3] Leo Breiman. Stacked regressions. *Machine Learning*, 24(1):49–64, 1996.
- [4] Mark J van der Laan and Sandrine Dudoit. Unified cross-validation methodology for selection among estimators and a general cross-validated adaptive epsilon-net estimator: Finite sample oracle inequalities and examples. *UC Berkeley Division of Biostatistics Working Paper Series*, 2003.
- [5] Mark J van der Laan and Eric C Polley. Super learner. *Statistical Applications in Genetics and Molecular Biology*, 6(1):1–23, 2007.
- [6] Alexander R Luedtke, Mark J Van Der Laan, et al. Statistical inference for the mean outcome under a possibly non-unique optimal treatment strategy. *The Annals of Statistics*, 44(2):713–742, 2016.
- [7] Wenjing Zheng and Mark J van der Laan. Asymptotic theory for cross-validated targeted maximum likelihood estimation. *UC Berkeley Division of Biostatistics Working Paper Series*, 2010.
- [8] Leo Breiman and Jerome H Friedman. Estimating optimal transformations for multiple regression and correlation. *Journal of the American Statistical Association*, 80(391):580–598, 1985.



# References II

- [9] Eric Polley, Erin LeDell, Chris Kennedy, and Mark van der Laan. *SuperLearner: Super Learner Prediction*, 2016. R package version 2.0-21.

# Oracle Inequality

Let  $L : \Psi \times \mathcal{O} \rightarrow \mathbb{R}$  be a loss function for  $\psi_0$  in the sense that

$$\psi_0 = \operatorname{argmin}_{\psi \in \Psi} E_0 \{L(\psi)(\mathcal{O})\}.$$

Define  $d_0(\psi, \psi_0) := E_0 \{L(\psi)(\mathcal{O}) - L(\psi_0)(\mathcal{O})\}$  and let  $p$  be the proportion of observations in the validation sample.

Assume

1.  $\psi_{n,m} \in \Psi$  with probability tending to one for  $m = 1, \dots, M$ .
2. For some  $C_0 < \infty$ ,  $\sup_{\psi \in \Psi} d_0(\psi, \psi_0) < C_0$  almost surely.
3. For some  $C_1 < \infty$ ,  $E_0 \{L(\psi)(\mathcal{O}) - L(\psi_0)(\mathcal{O}) - d_0(\psi, \psi_0)\}^2 \leq C_1 d_0(\psi, \psi_0)$  for all  $\psi \in \Psi$ .

For every  $\lambda > 0$  and  $C(\lambda) := \frac{2}{3}(1 + \lambda)^2(C_0 + C_1)$ ,

$$\begin{aligned} & E_{B_n} \{d_0(\psi_{n,SL}, \psi_0)\} \\ & \leq (1 + 2\lambda) E_{B_n} \{d_0(\psi_{n,OR}, \psi_0)\} + 2C(\lambda) \left( \frac{1 + \log K(n)}{np} \right) \end{aligned}$$

# Asymptotics

The cross-validated  $R^2$  estimate is

$$\begin{aligned} R_n^2(\omega_n, \psi_n, \omega_n) &= 1 - \frac{\sum_{k=1}^K \frac{1}{|V_k|} \sum_{i \in V_k} \{Y_{\omega(T_k),i} - \psi_{\omega(T_k),SL}(X_i)\}^2}{\sum_{k=1}^K \frac{1}{|V_k|} \sum_{i \in V_k} \{Y_{\omega(T_k),i} - \bar{Y}_{\omega(T_k)}\}^2} \\ &= 1 - \frac{\theta_{1,n}}{\theta_{2,n}} \end{aligned}$$

For  $k = 1, \dots, K$ , define

$$\begin{aligned} D_{0n,k}(\psi_\omega)(O) &:= \{Y_{\omega(T_k)} - \psi_{\omega(T_k)}(T_k)(X)\}^2 - E_0 [\{Y_{\omega(T_k),i} - \psi_{\omega(T_k),SL}(X_i)\}^2] , \\ D_{0n,k}(\bar{Y}_\omega)(O) &:= \{Y_{\omega(T_k)} - \bar{Y}_{\omega(T_k)}\}^2 - E_0 [\{Y_{\omega(T_k),i} - \bar{Y}_{\omega(T_k)}(X_i)\}^2] . \end{aligned}$$

# Asymptotics

$n^{1/2}(\theta_{1,n} - \theta_{1,0}) \rightarrow \text{Normal}(0, \sigma_1^2)$ , with

$$\sigma_1^2 = \frac{1}{K} \sum_{k=1}^K E_0\{D_{0n,k}(\psi_{\omega_0})(O)^2\}.$$

$n^{1/2}(\theta_{2,n} - \theta_{2,0}) \rightarrow \text{Normal}(0, \sigma_2^2)$ , with

$$\sigma_2^2 = \frac{1}{K} \sum_{k=1}^K E_0\{D_{0n,k}(\bar{Y}_{\omega_0})(O)^2\}.$$

Let  $D_{0n,k}(\psi_\omega, \bar{Y}_\omega) = (D_{0n,k}(\psi_\omega), D_{0n,k}(\bar{Y}_\omega))^T$ ,  $g(\theta) = \log(\theta_1/\theta_2)$ , and  $\nabla g(\theta) = (1/\theta_1, -1/\theta_2)^T$ .

# Asymptotics

We have  $n^{1/2}\{g(\theta_n) - g(\theta_0)\} \rightarrow \text{Normal}(0, \sigma_3^2)$ , where  $\sigma_3^2$  is

$$\nabla g(\theta_0)^T \frac{1}{K} \sum_{k=1}^K E_0\{D_{0n,k}(\psi_{\omega_0}, \bar{Y}_{\omega_0})(O)D_{0n,k}(\psi_{\omega_0}, \bar{Y}_{\omega_0})(O)^T\} \nabla g(\theta_0).$$

# Canonical correlation

Let  $Y_\alpha = \sum_{j=1}^J \alpha_j Y_j$  and  $X_\beta = \sum_{d=1}^D \beta_d X_d$ .

The first-order **canonical variate** of  $X$  and  $Y$  is found by maximizing

$$\frac{E_0\{(Y_\alpha - \mu_{0,\alpha})(X_\beta - \mu_{0,\beta})\}}{E_0\{(Y_\alpha - \mu_{0,\alpha})^2\}E_0\{(X_\beta - \mu_{0,\beta})^2\}}$$

over  $\alpha$  and  $\beta$  under constraint that variances equal one.

The canonical correlation is the **correlation** between  $X_{\beta,0}$  and  $Y_{\alpha_0}$ .

# Canonical correlation

If  $\psi_{0,j} = X_{\beta}$  for all  $j = 1, \dots, J$ , the optimal  $R^2$  equals the squared first-order canonical correlation.

To illustrate difference, consider

- ▶  $X_d \sim \text{Normal}(0, 1)$ ,  $d = 1, 2$
- ▶  $Y_j = X_j^2$  for  $j = 1, 2$

Canonical correlation measures **linear association** between  $X$  and  $Y$ .

- ▶ Canonical correlation **equals zero**.

Optimal  $R^2$  measures how well we **predict**  $Y$  using  $X$ .

- ▶ Optimal  $R^2$  **equals one**.