# Maching Learning Basics

Jae Yeon Kim

17 April, 2019

# DGP

- ▶ Statistics starts think of the data as being generated by a black box.
    - ▶ y <- nature <- x
- ▶ Data analysis
    - ▶ Prediction (algorithms; e.g., mean)
    - ▶ Information (inference; e.g., confidence intervals)

# Two cultures (Breiman 2001)

▶ Statistical inference

    ▶ y <- some probability models (e.g., linear regression, logistic regression, Cox model) <- x

▶ Machine learning

    ▶ y <- unkown <- x

    ▶ y <-> decision trees, neutral nets <-> x

    ▶ "The problem is to find an algorithm f(x) such that for future x in a test set, f(x) will be a good predictor of y."

# Data Modeling vs. Algorithmic Modeling

"There are **two cultures** in the use of statistical modeling to reach conclusions from data. One assumes that the data are generated by a **given stochastic data model**. The other uses **algorithmic models** and treats the data mechanism as **unknown**."

"Algorithmic models, both in theory and practice, has developed rapidly in fields of outside statistics. It can be used on large complex data sets and as a more accurate and informative alternative to data modeling on smaller data sets."

# Statistical model

▶ The following discussions come from Athey and Imbens's review paper (2019)

▶ Specifying a target (i.e., an estimand; a joint distribution of data)

▶ Fitting a model to data using an objective function (e.g., the sum of squared errors)

▶ Reporting point estimates and standard errors

▶ Validation by yes-no using goodness-of-fit tests and residual examination

# ML

- ▶ Developing algorithms

- ▶ Prediction power not structural/causal parameters

- ▶ Using out-of-sample comparisons (cross-validation) not in-sample goodness-of-fit measures

- ▶ The major problem is to avoid "the curse of dimensionality"

  - ▶ You can reduce dimensionality by limiting covariates (features) to essential ones

  - ▶ Or you can add many functions of the predictor variables (e.g., Support Vector Machines)

# Underdeveloped in the ML literature

▶ No interventions $\rightarrow$ no causal arguments (Judea Pearl's criticism)

▶ Need to exploit the structure of the problems

    ▶ The causal nature of the estimands

# ML - Terminology

- Sample to estimate parameters = Training sample
- Estimating the model = Being trained
- Regressors, covariates, or predictors = Features
- Regression parameters = weights
- Prediction problems = Supervised + Unsupervised

# ML and Causal inference

▶ Causal inference

  ▶ Designs: randomized experiments, uncofoundedness, instrumental variables, regression discontinuity, panel data, difference-in-differences

  ▶ Questions: ATE, LATE, CATE, Optimal Treatment Assignment Policies, Structural Parameter Estimation, Multiple Testing (Heterogeneous Treatment Effects)

ML: has a missing data problem (How can ML handle what if questions?)