

Machine Learning Methods in Finance: Recent Applications and Prospects

Daniel Hoang, Kevin Wiegatz

Long Zhen

Contents

- Introduction to ML
- Taxonomy of ML
- Summary

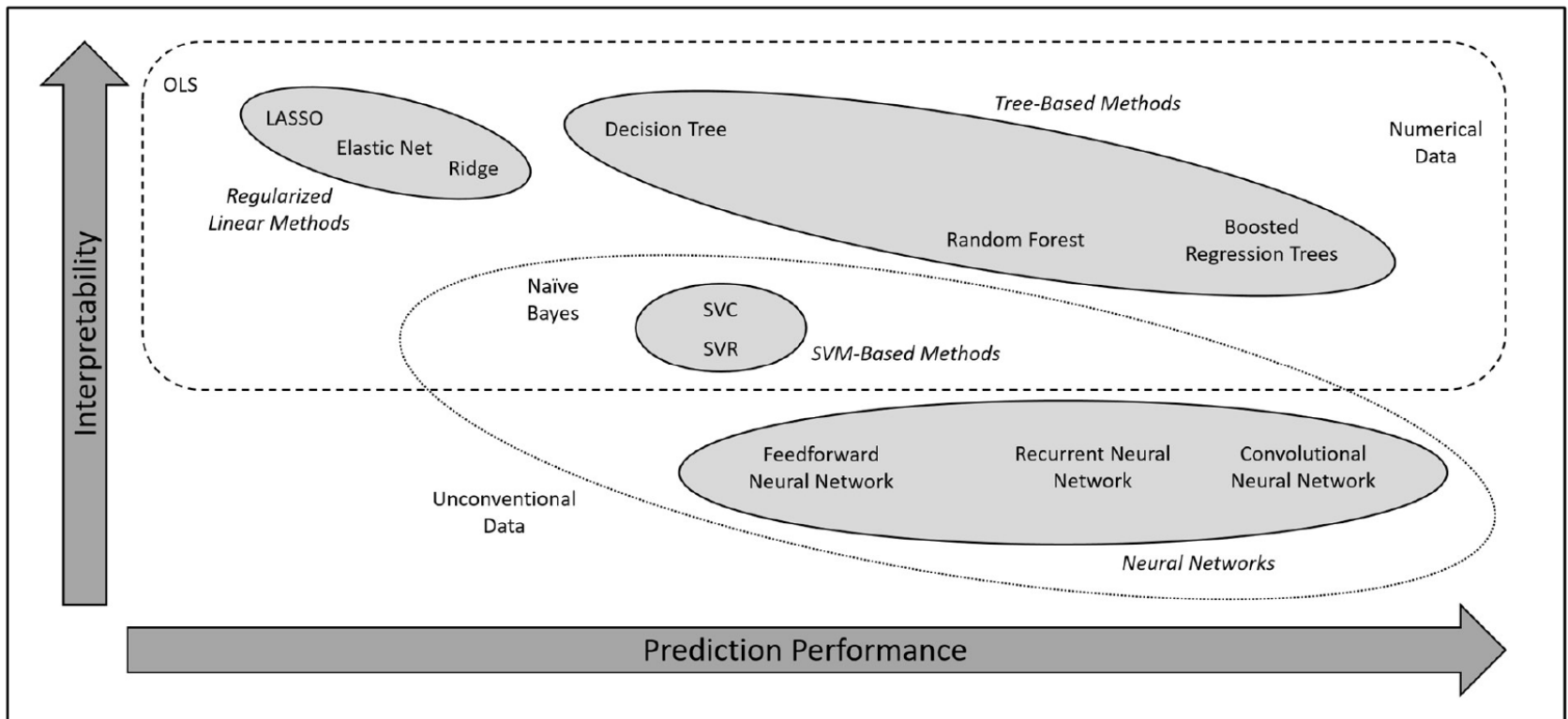
1. Introduction to ML

- Instead of causal explanations, ML serves for prediction (supervised) or data structure inference(unsupervised).

Approach	Data	Method	Results		Usage	Purpose
Traditional Econometrics	Labeled Data $(X_i, y_i)_i$	Linear Regression (OLS)	Explanatory Model	+	Statistical Significance	Explanation “ β ”
Supervised Learning	Labeled Data $(X_i, y_i)_i$	Supervised ML Method	Prediction Model	+	Prediction Performance	Prediction “ \hat{y} ”
Unsupervised Learning	Unlabeled Data $(X_i)_i$	Unsupervised ML Method	Data Structure Model	+	Data Structure Characteristics	Data Structure Inference “ \hat{X} ”

1.1 Supervised learning

- Purpose: prediction
- Training data & test data
- Prediction performance vs. interpretation



- **OLS: BLUE**; best interpretation, weak performance
 - Nonlinear transformation and interactions?
- → **LASSO/Ridge/Elastic Net**: introduce bias by adding penalty term
- **Regression tree**
 - → random forest: bootstrap samples and build separate trees
 - → GBDT: build trees iteratively

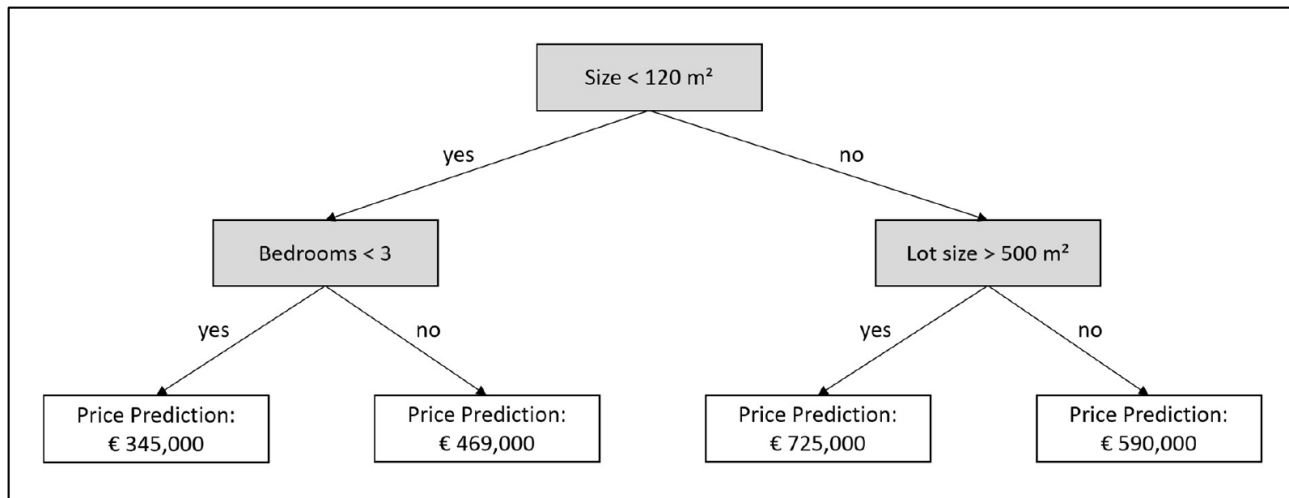
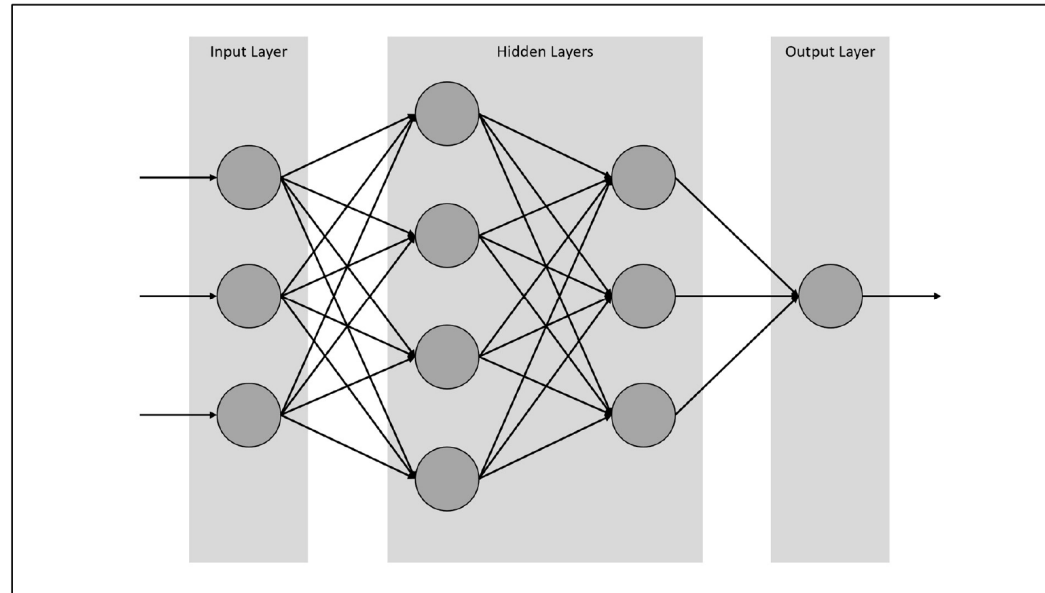


Figure 4. Illustrative depiction of a decision tree trained for house price prediction

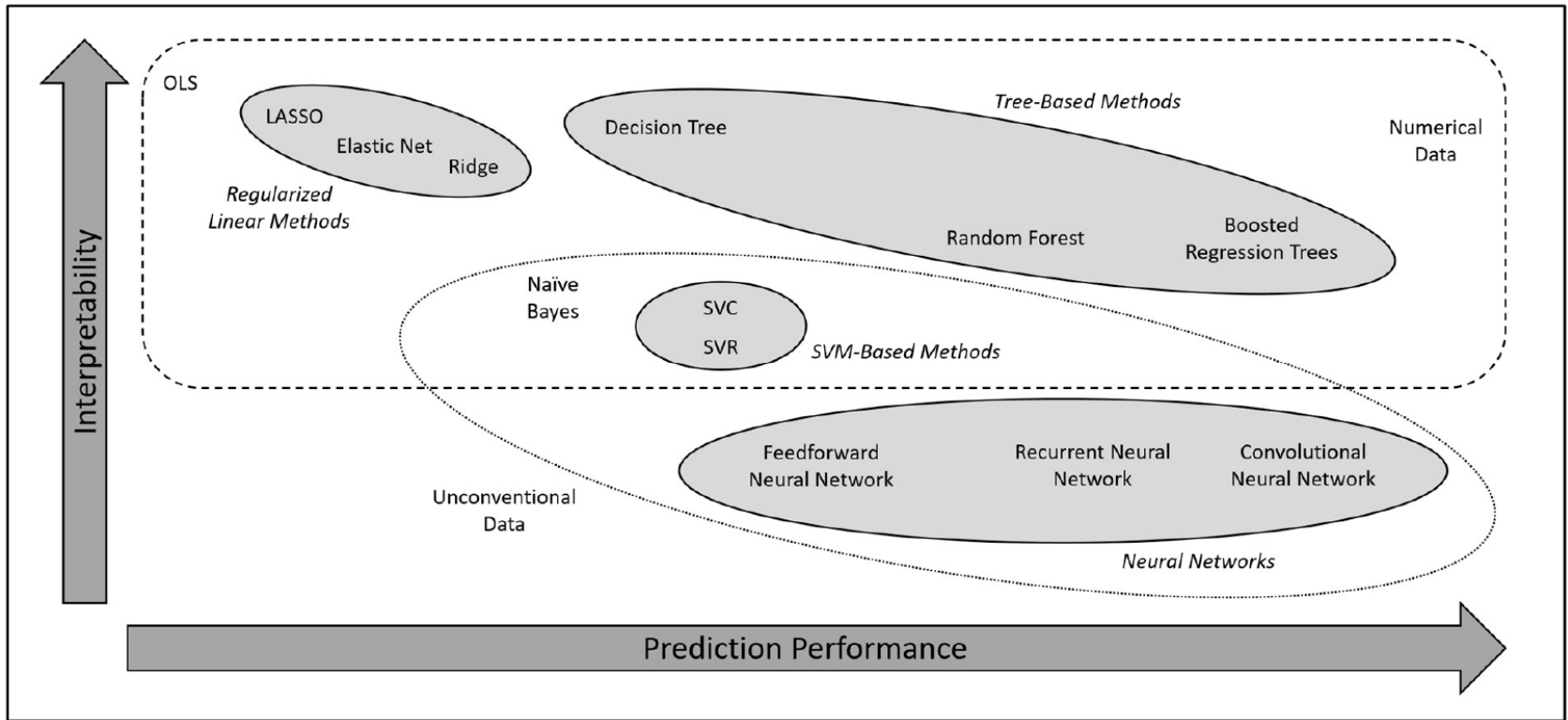
Hard to interpret:

- Feed-forward NN



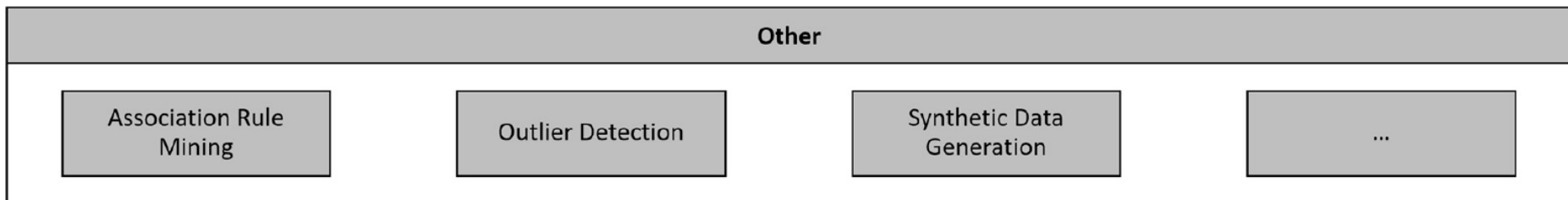
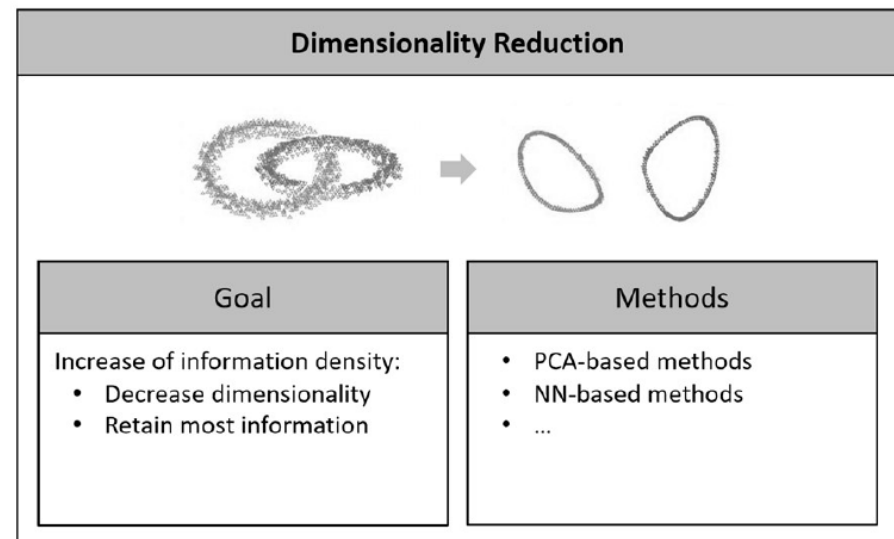
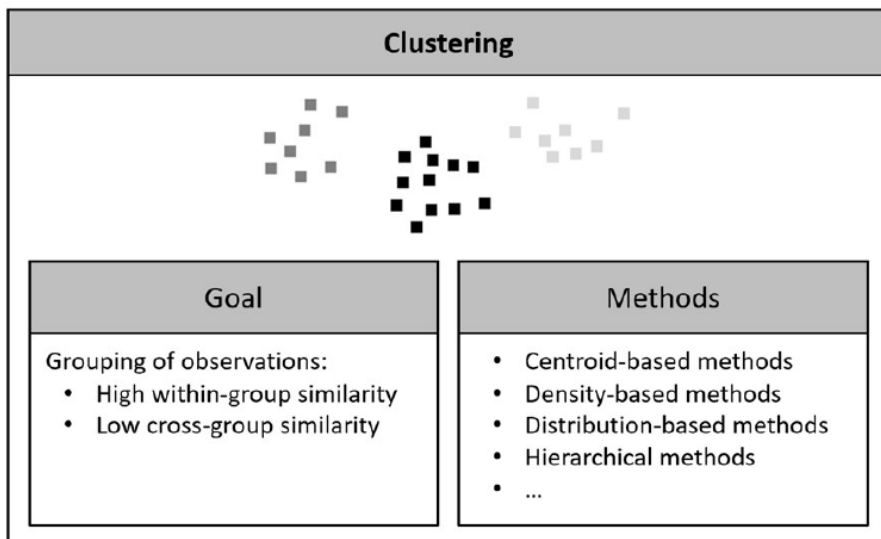
- → recurrent NN: for sequential data
 - Eg. Gated recurrent units/LSTM
- → convolutional NN: for visual data

- Older methods:
 - Naïve Bayes
 - SVM/SVR



1.2 Unsupervised learning

- Purpose: data structure inference



Clustering

- Centroid-based methods: ^{重心法} eg. K-means
 - After the initial positioning of the centroids, they iteratively update their position to arrive at suitable clusters
- Density-based methods: eg. DBSCAN
 - Group observations with many similar observations nearby into clusters
- Distribution-based methods: eg. Gaussian mixture models
 - Based on whether observations belong to the same statistical distribution
- Hierarchical methods: eg. BIRCH
 - Iteratively combine smaller clusters into larger clusters

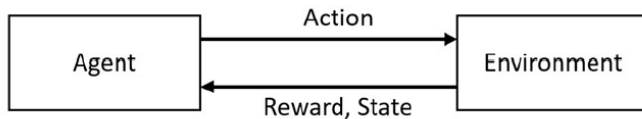
Dimensionality reduction

- Decrease the dimensionality while retain most of the inherent information
- PCA-based: cover as much of the data's variance as possible
- NN-based: autoencoder
 - Consist of an encoder network that creates a condensed representation of the input data and a subsequent decoder network that reconstructs the original data. A special bottleneck layer connect the encoder and decoder to train

- Others:
- Association rule mining: identify relations between variables
- Outlier detection methods

1.3 Reinforcement learning and other types

Reinforcement Learning



Description

- Markov decision process model with environment and agent whose actions bring rewards and change the environment
- Goal: policy that maximizes the expected total reward
- Data: description of the Markov decision process
- Idea: the algorithm trades short-term vs. long-term rewards

Semi-Supervised Learning



Description

- Combination of supervised and unsupervised learning
- Goal: prediction (as in supervised learning)
- Data: few labeled examples and many unlabeled examples
- Idea: the additional unlabeled data adds information about the probability distribution of the input data

Other

演绎学习
Deductive Learning

联邦学习
Federated Learning

遗传算法
Genetic Algorithms

...

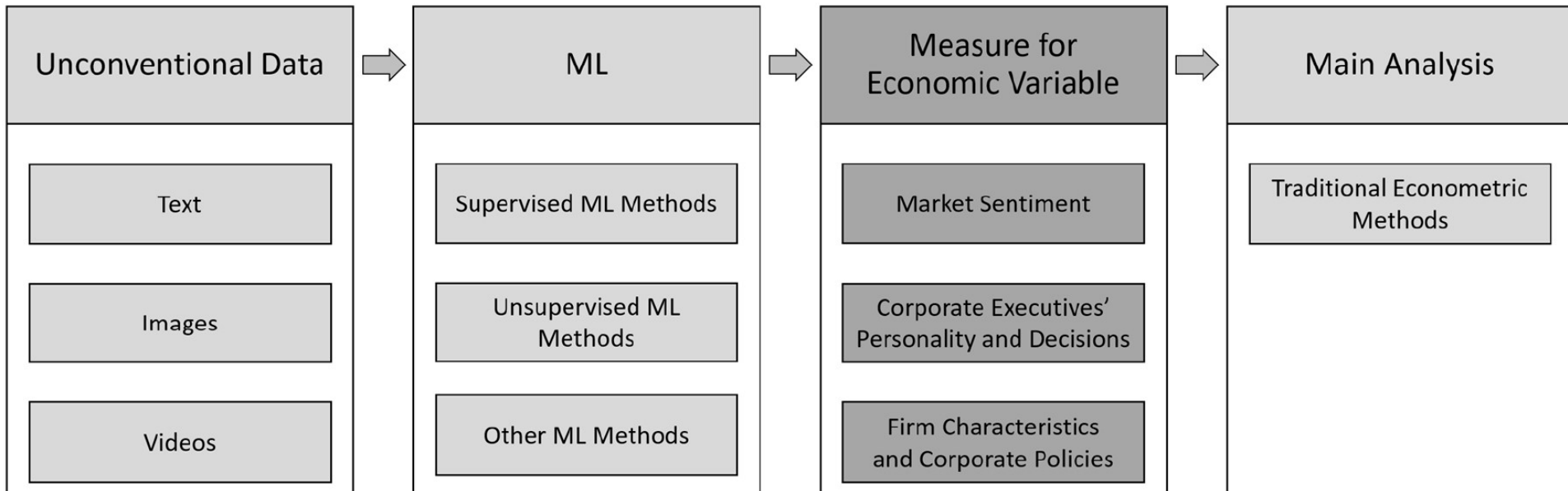
- Semi-supervised: eg. Active learning (closely related variant of semi-supervised learning): train the model with labeled examples → calculate importance score for the unlabeled → label and reduce costs

2. Taxonomy of ML Applications in Financial Economics

- ML solves different problems compared to OLS
- Three typical applications:

1	Construction of Superior and Novel Measures	$y = \beta X + \varepsilon$
2	Reduction of Prediction Error in Economic Prediction Problems	$\hat{y} = f(X)$
3	Extension of the Existing Econometric Toolset	$y = \beta X + \varepsilon$ & ML

2.1 Construction of superior and novel measures



Measures of Market Sentiment

Stock Market Sentiment

Methods	Data
Antweiler and Frank (2004)	Naïve Bayes/SVM
Renault (2017)	Yahoo Finance message board
Bartov et al. (2017)	StockTwits
Barbon et al. (2019)	Twitter
Ke et al. (2019)	Naïve Bayes
Huang et al. (2014)	Firm-specific news
Manela and Moreira (2017)	Customized ML
Vamosy (2020)	Dow Jones Newswire articles
Liew and Wang (2016)	Naïve Bayes
	Analyst reports
	Wall Street Journal front-page articles
	Deep learning
	StockTwits posts
	Commercial ML
	Twitter posts

Product Market Sentiment

Methods	Data
Tang (2018)	Commercial ML
	Twitter posts

Measures of Corporate Executives' Personality and Decisions

Executives' Personality	Methods	Data	Measures
Gow et al. (2016)	Naïve Bayes/SVM	Yahoo Finance message board	Naïve Bayes/SVM
Hrazdil et al. (2020)	IBM service	Conference calls	CEO and CFO personality
Hsieh et al. (2020)	Face detection	Executives' business headshot images	trustworthiness
Du et al. (2019)	Textual analysis	Mutual fund managers' letters to shareholders	Managers' level of confidence
Executives' Decisions	Methods	Data	Measures
Bandiera et al. (2020)			Whether perform low/high-level tasks
Barth et al. (2020)		Earnings conference call	How they withhold information
Hu and Ma (2020)		Youtube videos	How founders speak? What they say? How they present visually

Measures of Firm Characteristics and Corporate Policies

Corporates

Methods	Data	Measures
	Conference call transcripts	Corporate cultures
	Annual report	Financial constraints
	SEC letters	Regulatory IPO concern

Li et al. (2020)

Buehlmaier and
Whited (2018)

Lowry et al. (2020)

Financials

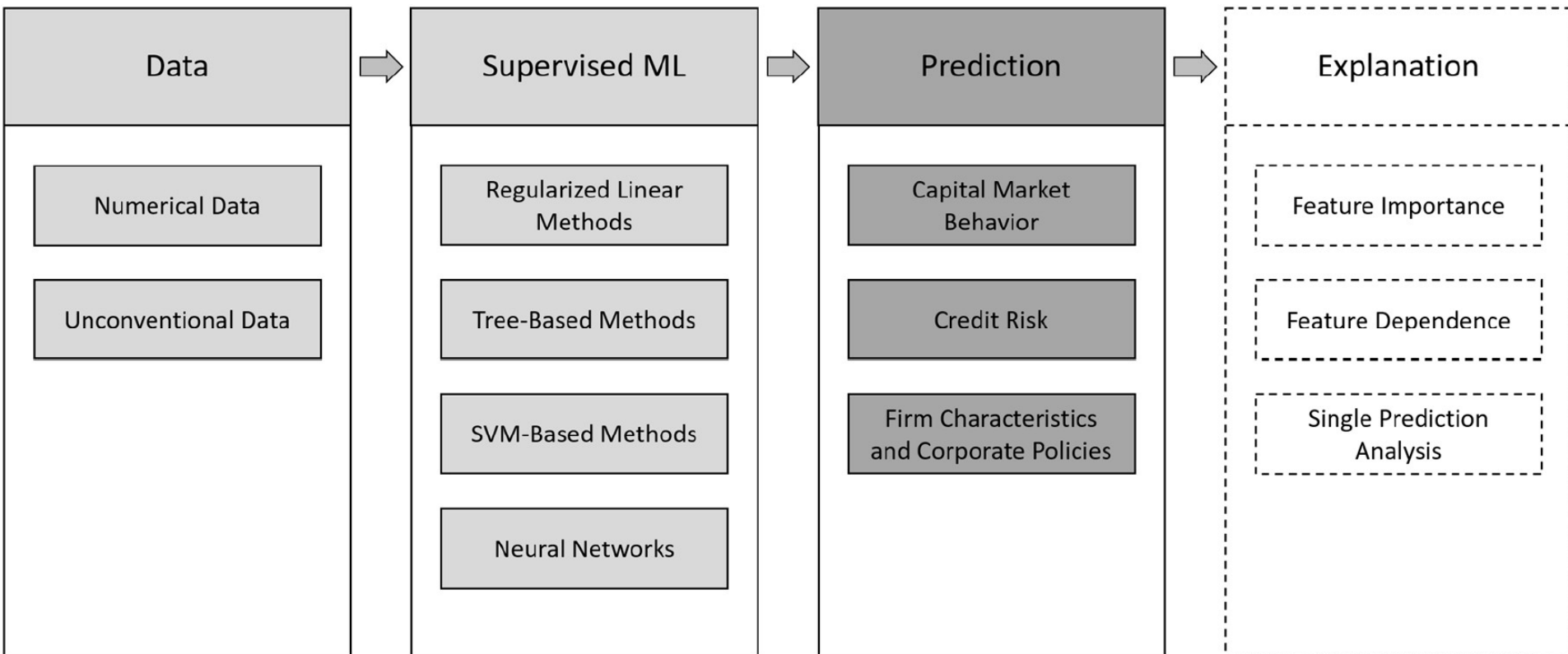
Methods	Data	Measures
Commercial ML	Banks' annual report	Risk exposure
Cluster		Venture capital relatedness
	Customer complaint texts	Bank misconduct

Hanley and Hoberg
(2019)

Bubna et al. (2020)

Bertsch et al. (2020)

2.2 Reduction of Prediction Error in Economic Prediction



- Explanation

- Feature importance: importance score
 - eg. Permutation importance (permute X_i and calculate the score)
- Feature dependence: relations between predictor and target
 - Partial dependence plots
- Single prediction analysis: disentangle the contribution of every predictor to a specific prediction value
 - Sharpley Additive Explanation(SHAP)

Prediction of Capital Market Behavior

Individual Stock Returns

Equity Risk Premium

Stochastic Discount Factor

Option Prices

Other

Rasekschaffe and Jones (2019)

Kelly et al. (2019)

Jacobsen et al. (2019)

Chen et al. (2019)

Hutchinson et al. (1994)

US treasury bonds
Bianchi et al. (2020)

Martin and Nagel (2019)

Freyberger et al. (2020)

Routledge (2019)

Kozak et al. (2018)

Yao et al. (2000)

bond liquidity
Reichenbacher et al. (2020)

Gu et al. (2019, 2020)

Grammig et al. (2020)

Adämmer and Schüssler (2020)

direction of changes in exchange rates
Spiegeleer et al. (2018)

Colombo et al. (2019)

Rossi (2018)

Chinco et al. (2019)

future stock volatility
Kogan et al. (2009)

Moritz and Zimmermann (2016)

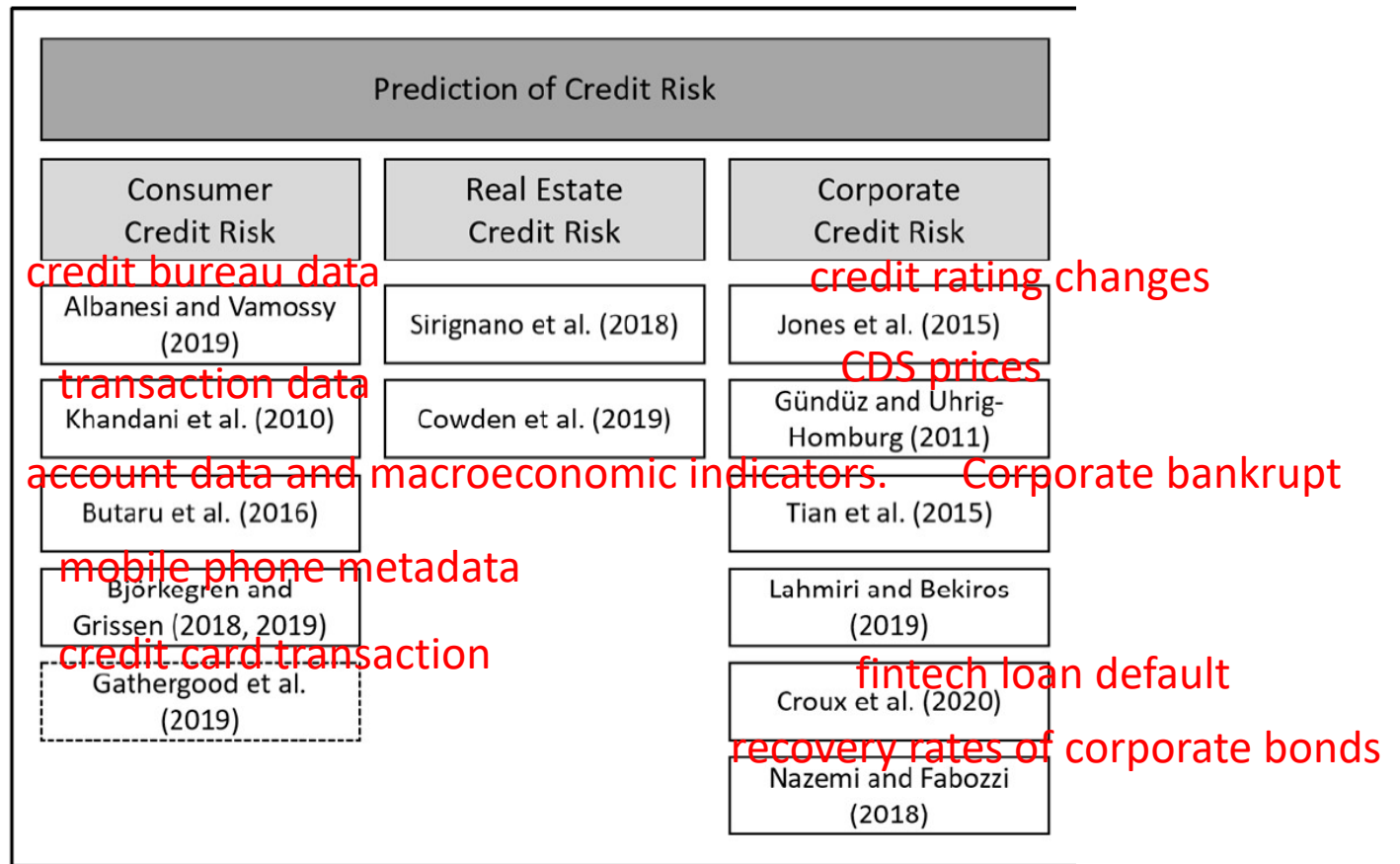
Amel-Zadeh et al. (2020)

VIX
Osterrieder et al. (2020)

lifespan of orders
McInish et al. (2019)

investors' portfolio allocation

Rossi and Utkus (2020)



Prediction of Firm Characteristics and Corporate Policies		
Firm Fundamentals	Accounting Fraud	Startups' Success
<div>leverage</div> Amini et al. (2019)	Bao et al. (2020)	Xiang et al. (2012)
<div>earnings</div> Van Binsbergen et al. (2020)	Brown et al. (2020)	Ang et al. (2020)

2.3 Extension of the Existing Econometric Toolset

- ML-enhanced instrumental variable regression
- Causal trees and forests

Causal ML		
Instrumental Variable Regression	Causal Trees and Forests	Other Causal ML
Belloni et al. (2012)	Athey and Imbens (2016)	Lee et al. (2010)
Carrasco (2012)	Wager and Athey (2019)	Mullainathan and Spiess (2017)
Hansen and Kozbur (2014)	Gulen et al. (2020)	Chernozhukov et al. (2017, 2018)
Hartford et al. (2017)	O'Malley (2018)	Athey et al. (2019)
Angrist and Frandsen (2019)		

- Recall 2SLS:

$$y_i = \alpha + \beta x_i + \epsilon_i (i = 1, \dots, n)$$
$$\text{cov}(x_i, \epsilon_i) \neq 0$$

- IV: $\text{cov}(z_i, \epsilon_i) = 0$; $\text{cov}(z_i, x_i) \neq 0$

- First stage:

$$x_i = \gamma + \delta z_i + u_i$$
$$\hat{x}_i = \hat{\gamma} + \hat{\delta} z_i$$

- Second stage:

$$y_i = \alpha + \beta \hat{x}_i + (\epsilon_i + \beta \hat{u}_i)$$

- Causal ML & IV

- Better predictions for the IV results in more precise estimates in the second stage.

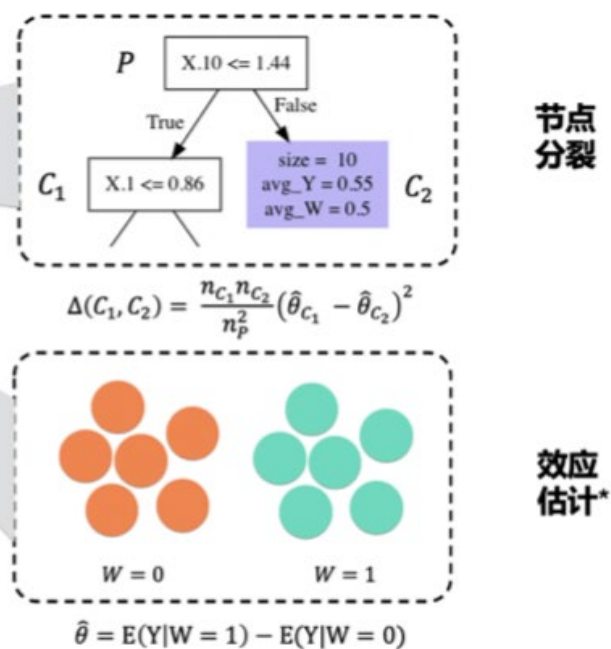
Extension of the Existing Econometric Toolset

- ML-enhanced instrumental variable regression
- Causal trees and forests

Causal ML			
	Instrumental Variable Regression	Causal Trees and Forests	Other Causal ML
LASSO	Belloni et al. (2012)	Athey and Imbens (2016)	propensity score Lee et al. (2010)
ridge	Carrasco (2012)	Wager and Athey (2019)	Data balance Mullainathan and Spiess (2017)
ridge	Hansen and Kozbur (2014)	Gulen et al. (2020)	calculate treatment effects Chernozhukov et al. (2017, 2018)
neural networks	Hartford et al. (2017)	O'Malley (2018)	Athey et al. (2019)
	Angrist and Frandsen (2019)		

- Causal trees and forests

- 从协变量中，找到一个最优分裂节点，最大化子节点间处理效应差异



* 适当假设下, 认为节点内数据同质无混淆

Summary

- ML methods:
 - Supervised learning
 - Unsupervised learning
 - Others (reinforcement learning...)
- Taxonomy of ML
 - Measure construction
 - Prediction
 - Econometric extension
- Summary