

Machine Learning for Stock Selection

Keywan Rasekhschaffe, PhD

Robert C. Jones

Agenda

- Issues of overfitting and maximizing the signal to noise ratio
- Data prep and feature engineering: often overlooked
- Evaluating your algorithm choice: what do you want to achieve?
- Pitfalls

[Rasekhschaffe, Jones \(2019\)](#) provide a concrete example

What is Machine Learning

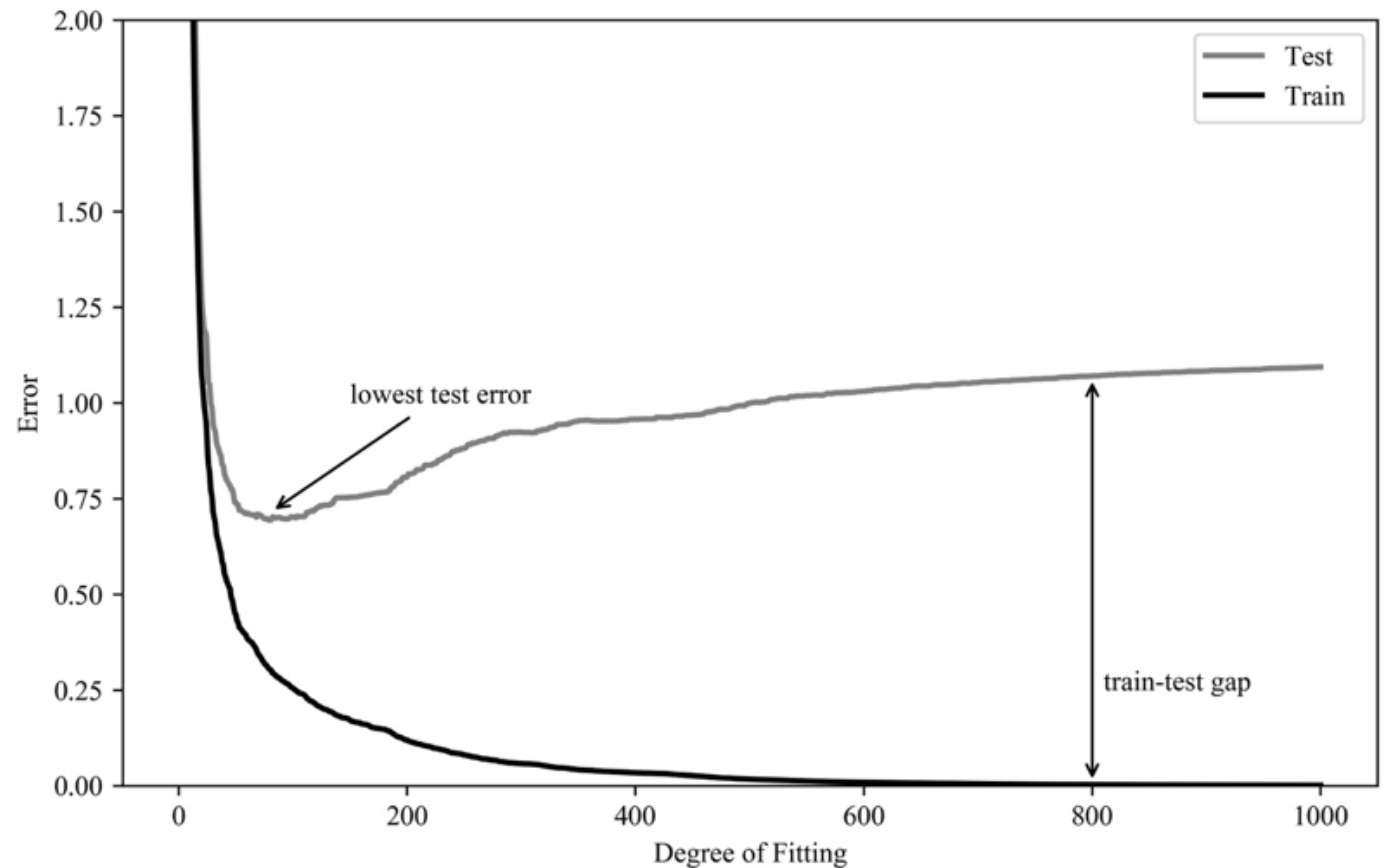
- **A branch of computer science concerned with the design of algorithms that learn from examples to make predictions**
- **Examples:**
 - Neural Networks, Deep Learning, CART, Random Forests, Gradient Boosting, Support Vector Machines
- **Many potential benefits:**
 - Can recognize complex, often non-linear patterns in large data sets
 - Many techniques work well with correlated features
 - Self correcting with additional training
 - Can detect dynamic relationships and is robust to factor decay
 - Potentially Less crowded (smart alpha)
- **One BIG problem:**
 - Overfitting

Two Ways to Overcome Overfitting

- **Feature engineering**
 - Use domain knowledge to increase the signal-to-noise ratio
 - Requires expertise in finance, statistics and data
- **Forecast ensembles**
 - The wisdom of crowds
 - Different algorithms
 - Different training sets
 - Different forecast horizons
 - Use averaging to cancel offsetting biases
 - Result: more signal, less noise

Overfitting and Underfitting

- More complexity initially leads to improved out-of-sample performance, but then results deteriorate
- Not optimal to overfit or underfit
- Train-test gap widens with increasing complexity: never judge results in-sample



Feature Engineering

- What problems do we ask the algo to solve?
- Incorporates domain knowledge
- Seldom discussed, but at least as important as prediction algorithms
- What type of data do we allow it to consider?
- Success requires expertise in investing, data and statistics

Examples of Feature Engineering

Investment Expertise

- How to classify stocks (region, industry, style)
- How to define success (returns, excess returns, alphas)
- Which factors to consider

Data Expertise

- Factor calculations
- Adjustments for errors and missing data
- Proper lags and adjustments to avoid look-ahead bias

Statistical Expertise

- Clustering
- Normalization
- Dimensionality reduction

Ensembles: Combing Diverse ML Forecasts

Diversification (I)

- Forecast combinations from **different algorithms**
- Different algos detect different attributes in the data

Diversification (II)

- Forecast combinations from **different forecast horizons**
- Different factors are important for different forecast horizons
- Technical factors are more important for shorter horizons
- Fundamental factors are more important for longer horizons

Diversification (III)

- Forecast combinations from **different training windows**
- Different factors are important in different market environments
- Training windows based on time, seasonality, or economic / market conditions

Sample Algorithms

Gradient Boosted Classification and Regression Trees

- Ensemble technique that combines weak learners (algorithms) into a strong learner (ensemble)

Random Forests

- Ensemble technique that overcomes overfitting by equal weighting forecasts from many trees

Deep Learning

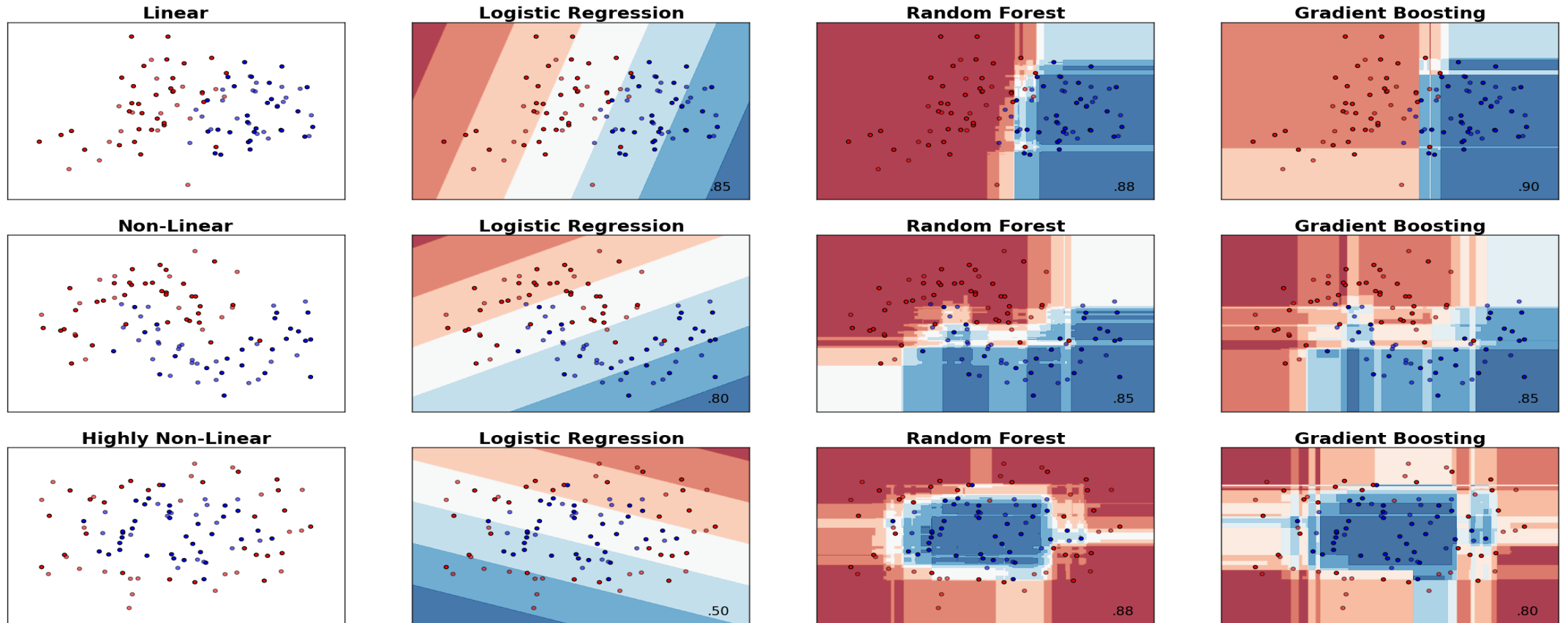
- Recent evolution of neural network (back propagation) algorithms
- Less likely to overfit

Support Vector Machines

- Learns by analogy
- Effective in high dimensional spaces

Why Use Multiple Algorithms?

Different algorithms detect different attributes in the data



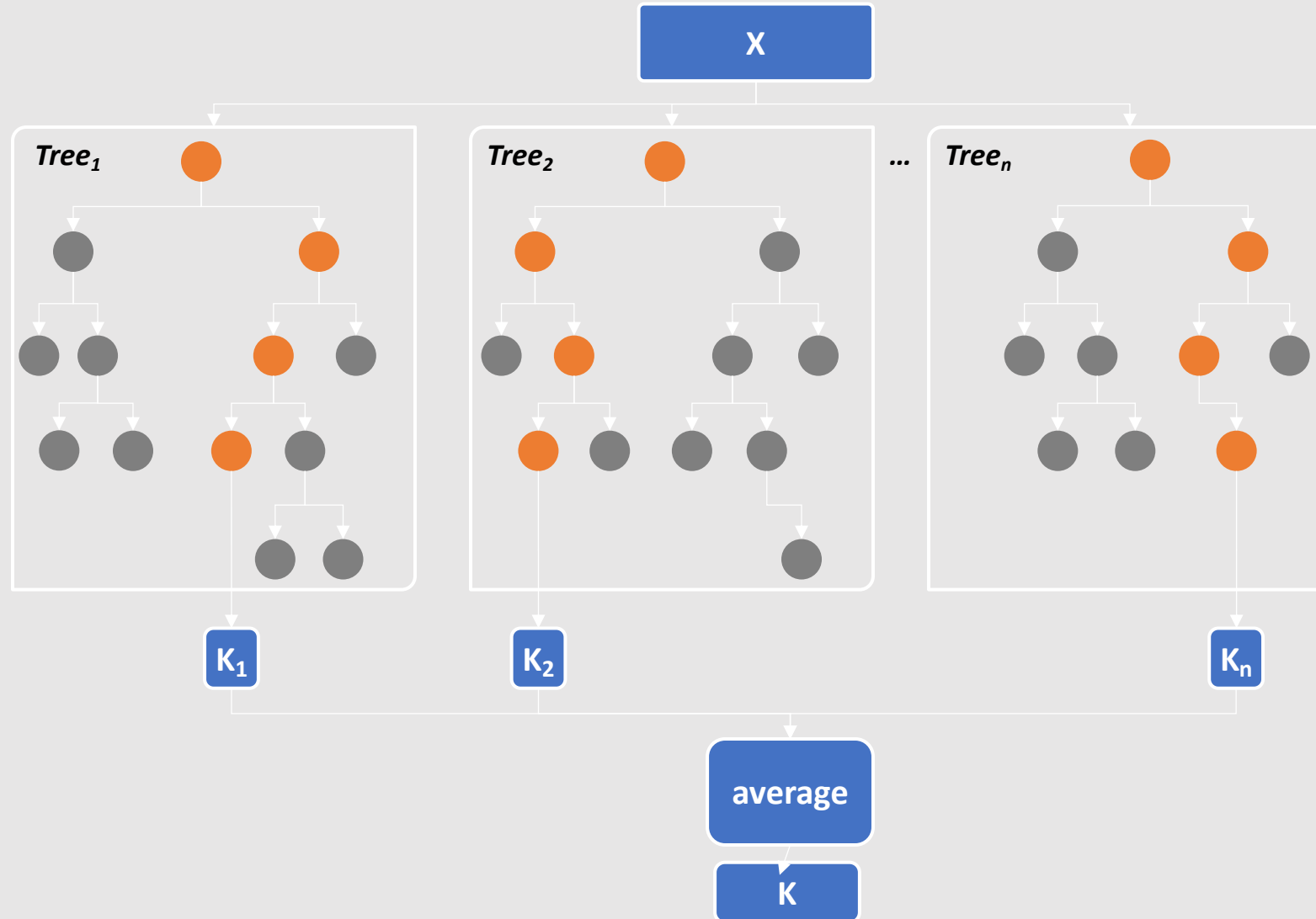
CART Models

- **CART** = **C**lassification **A**nd **R**egression **T**ree
- Classifies data into binary branches via hierarchical splits
- Easy to overfit
- “Pruning” can reduce overfitting
- Important building block for ensemble algos
- Although still popular in finance, single trees are no longer widely used in serious ML applications

Random Forests

- Increase number of trees and construct a forest
- Randomly selects subsets of features
- Ensemble “bagging” (bootstrap aggregating) algorithm
- Uses majority decision or average of base estimators
- Model variance decreases

Random Forest Diagram



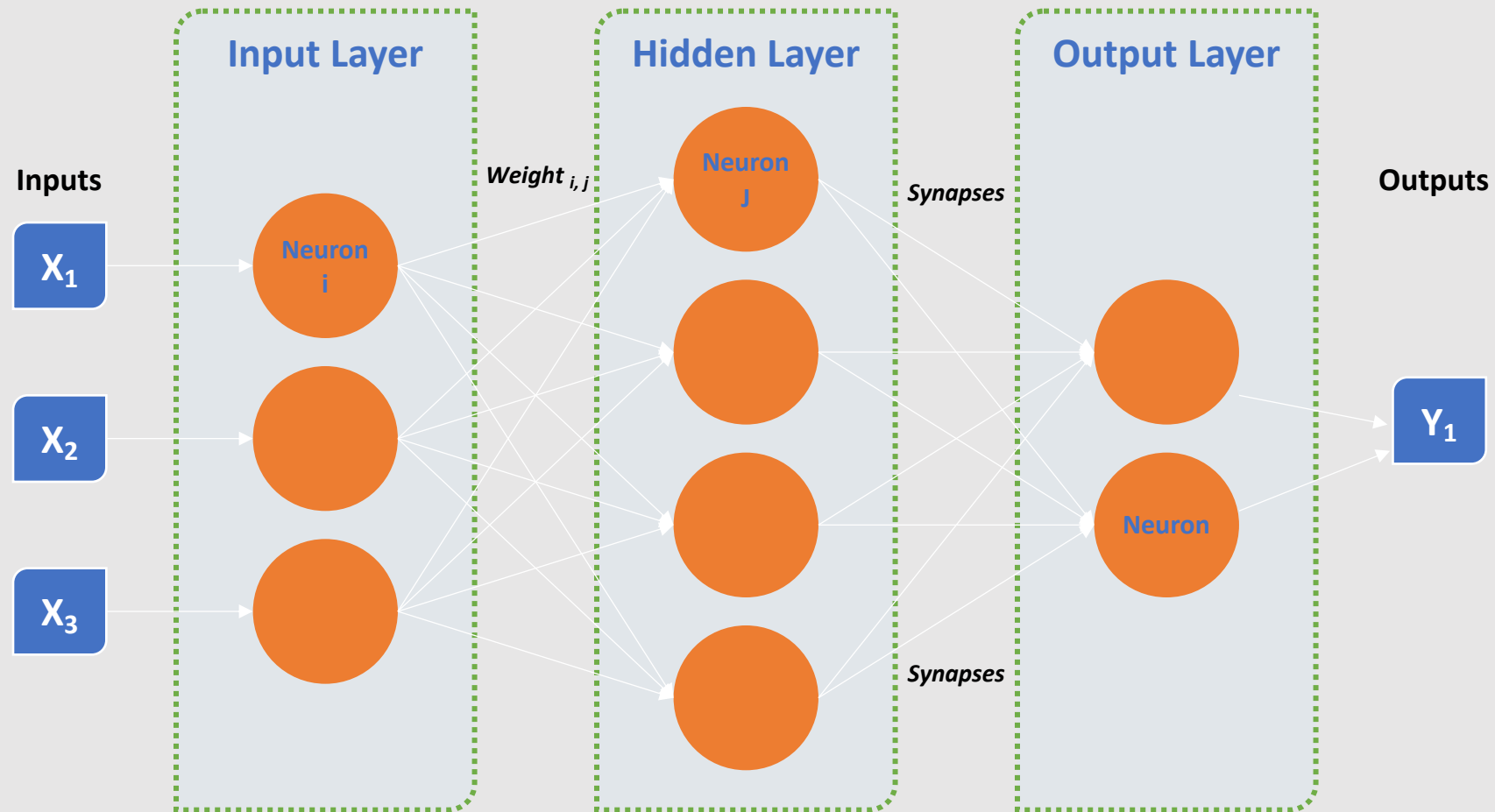
Boosting

- Combines weak learners (factors or simple models) into a strong learner (ensemble model)
- Ensemble technique like Random Forests
- Iteratively re-weight observations
- Put more emphasis on misclassified observations
- Each subsequently chosen weak learner gets assigned a decreasing weight
- Base learners can be trees, decision stumps, etc.
- Examples are AdaBoost and **Gradient Boosted Regression Trees**

Artificial Neural Networks

- Interconnected group of “neurons”
- Inspired by biology
- Input layer <---> hidden layer <---> output layer is the most common architecture
- Shallow architectures work quite well for forecasting returns

Artificial Neural Network Example



Deep Learning

- We can avoid overfitting using regularization and randomization techniques
- Dropout: randomly disable neurons in the learning process
- Forces the network to generalize by learning multiple representations of the same data
- Gold standard for computer vision / audio
- The algorithm (sometimes) takes care of feature engineering
- Requires carefully tuned architecture / parameters
- Can be computationally expensive

Support Vector Machines

- Finds hyperplane with largest distance between classes in the training data
- Can use linear or non-linear kernel
- Linear kernel related to logistic regression, but more effective in high-dimensional spaces
- “Soft margin” often used in practice when classes aren’t perfectly separable
- Flexible and powerful, but often slow

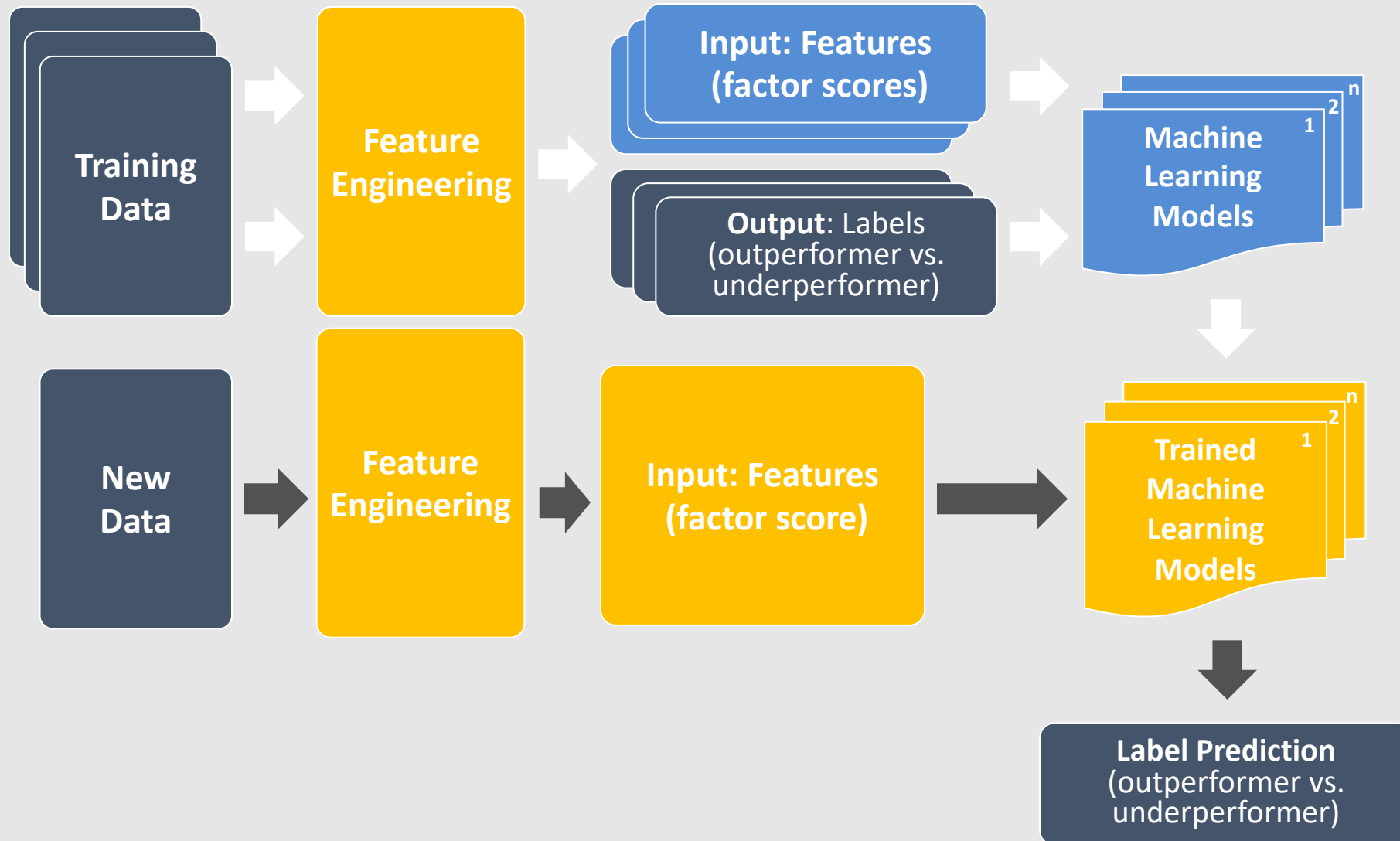
Forecast Combinations and Diversification

- Algos
- Training windows
- Forecast horizons
- Factor library (or libraries)

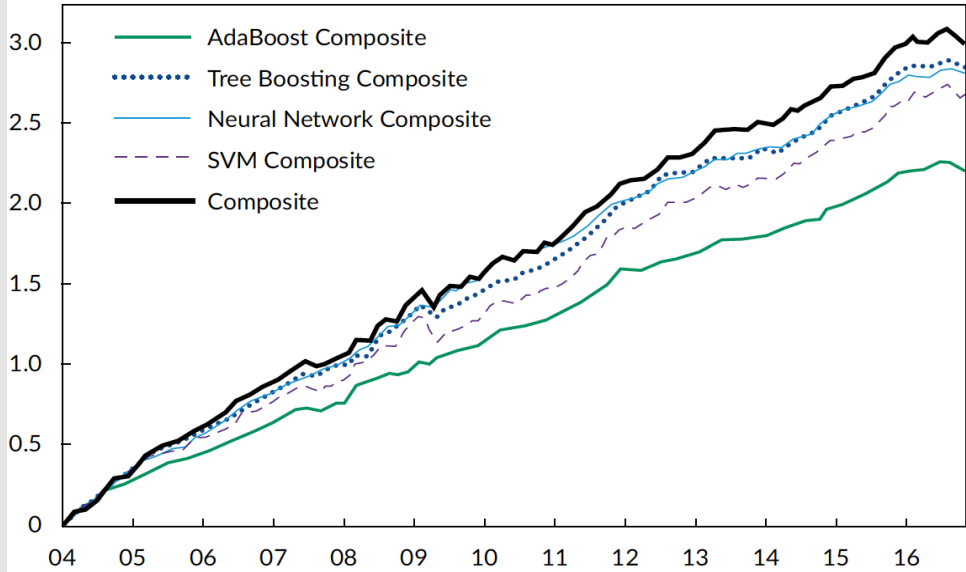
Different Training Windows

- Best to train algos using data from periods that are “similar” to the forecast period
- **Recency**: More recent data likely to be more relevant
- **Seasonality**: Data from same season (month) likely to be more relevant
- **Context**: Data from periods with similar macro conditions likely to be more relevant
 - ✓ Similar risk environment (VIX, credit spreads)
 - ✓ Similar monetary conditions (level of short rates, slope of the yield curve)
 - ✓ Similar valuation levels (dividend yield, E/P)
- **Hedging**: Data from periods where other forecasts are noisy are likely to improve overall accuracy

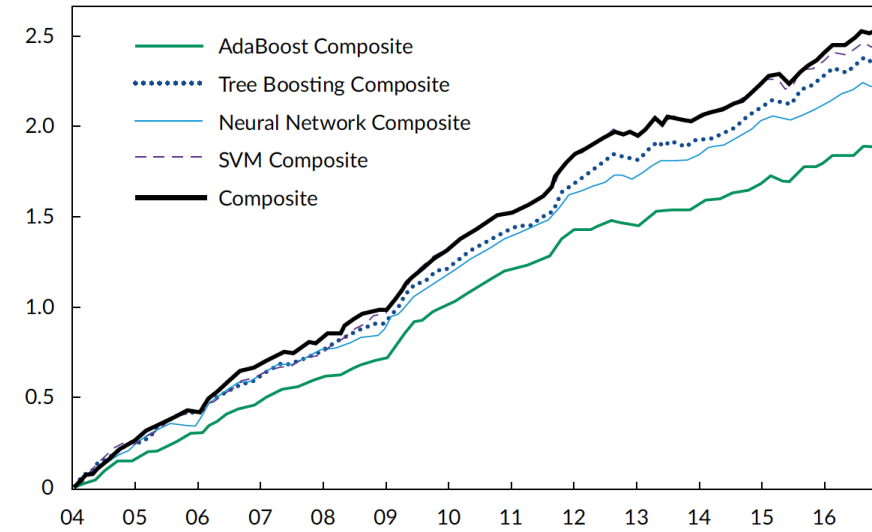
Overview of Our Approach



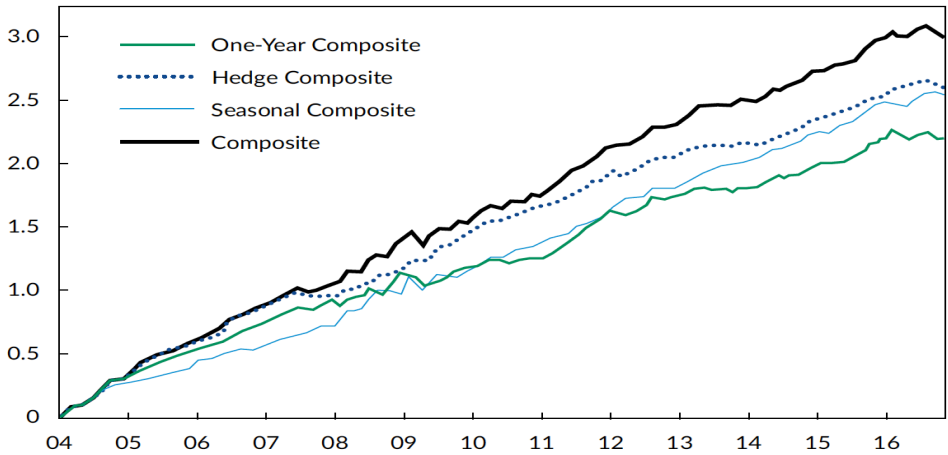
A. US Algorithms



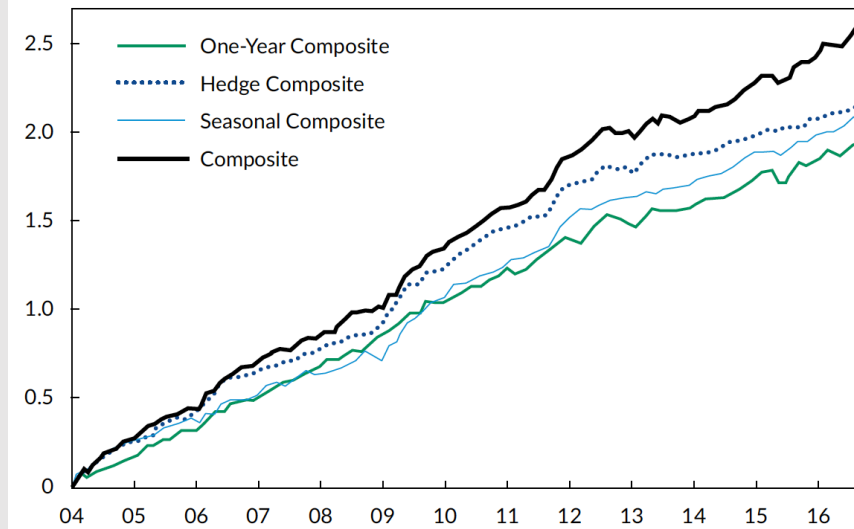
B. ROW Algorithms



C. US Training Windows



D. ROW Training Windows



Results I

Table 2. Information Coefficients for Models, Composite, and Benchmarks

Model	US Zone		ROW	
	Rank IC	t-Statistic	Rank IC	t-Statistic
ML composite	6.48%	15.87	6.43%	16.37
12-Month period				
AdaBoost	3.19	9.07	3.54	10.98
SVM	4.61	8.00	5.08	10.35
Tree boosting	4.66	9.59	4.79	11.12
Neural network	3.99	12.94	3.17	11.97
Seasonal				
AdaBoost	3.17	11.53	3.25	11.75
SVM	5.01	9.65	4.79	10.56
Tree boosting	5.00	12.04	4.53	11.89
Neural network	5.20	15.39	4.41	11.92
Hedge				
AdaBoost	3.92	14.06	3.53	13.52
SVM	4.94	13.21	5.23	14.36
Tree boosting	5.00	14.14	4.78	16.16
Neural network	4.58	14.24	4.47	14.61
Top 10 factors benchmark	2.81	6.33	4.49	9.44
OLS benchmark	3.36	5.78	2.71	4.72

Correlations

Table 3. Correlations between Machine Learning Portfolios and the Market Factor

Portfolio	US Equal Weighted	US Risk Weighted	ROW Equal Weighted	ROW Risk Weighted	MKT
US equal weighted	1.00	0.81	0.32	0.29	-0.57
US risk weighted		1.00	0.28	0.30	-0.20
ROW equal weighted			1.00	0.96	-0.16
ROW risk weighted				1.00	-0.05
MKT					1.00

Note: MKT refers to the Fama-French (1992) market risk factor.

Table 4. Regression Results vs. Fama–French–Carhart Factors (t-statistics in parentheses)

A. Regression of MLA decile spreads on Fama–French–Carhart factors

	US Zone				ROW			
	Equal Weighted	Risk Weighted	Top 10 Benchmark	OLS Benchmark	Equal Weighted	Risk Weighted	Top 10 Benchmark	OLS Benchmark
Excess return	1.60 (6.04)	1.95 (11.32)	1.23 (6.32)	1.10 (3.56)	1.50 (9.32)	1.64 (11.61)	1.06 (5.91)	0.79 (5.23)
Alpha	1.84 (9.48)	2.01 (11.93)	1.17 (7.05)	1.03 (3.79)	1.53 (9.48)	1.65 (11.53)	1.11 (6.19)	0.85 (6.18)
MKT	-0.38 (-6.93)	-0.09 (-1.91)	0.07 (1.53)	0.02 (0.28)	-0.04 (-0.80)	0.01 (0.14)	-0.07 (-1.44)	-0.12 (-3.03)
SMB	-0.23 (-2.44)	-0.07 (-0.84)	0.03 (0.38)	-0.08 (-0.63)	-0.07 (-0.93)	-0.05 (-0.67)	-0.1 (-1.22)	0.02 (0.27)
HML	-0.06 (-0.67)	0.00 (0.05)	0.37 (4.97)	0.05 (0.34)	-0.09 (-1.25)	-0.11 (-1.70)	0.1 (1.18)	0.01 (0.19)
MOM	0.18 (3.76)	0.08 (1.90)	-0.11 (2.71)	0.22 (3.43)	-0.06 (-1.44)	-0.04 (-1.15)	0.05 (1.16)	0.14 (4.10)
Adjusted R^2	0.47	0.07	0.29	0.32	0.01	0.00	0.03	0.20
No. of observations	154	154	154	154	154	154	154	154

(continued)

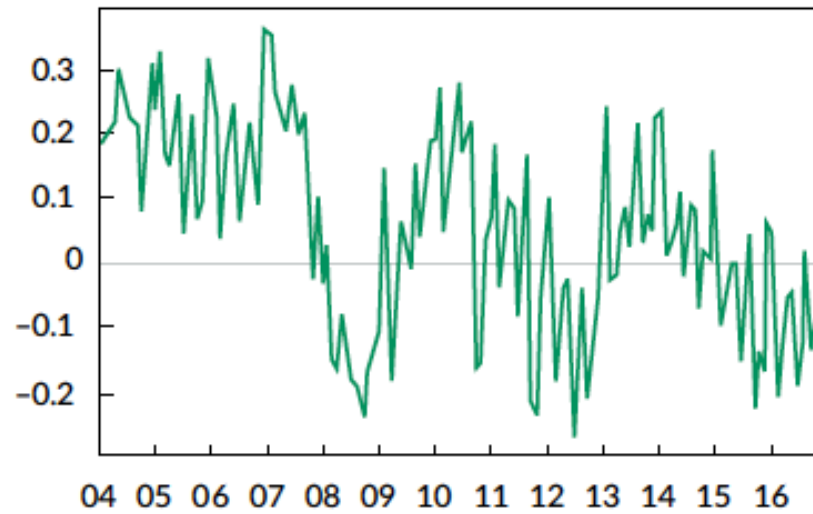
Table 5. Fama–MacBeth Regressions (t-statistics in parentheses)

Forecast/Characteristic	US Zone		ROW	
	Equal Weighted	Risk Weighted	Equal Weighted	Risk Weighted
ML composite	1.71 (10.04)	1.77 (12.05)	1.00 (9.22)	1.12 (10.44)
Earnings revision	0.21 (2.17)	0.15 (1.57)	0.33 (4.46)	0.32 (4.25)
Dividend yield	0.11 (0.94)	0.09 (0.91)	0.16 (1.38)	0.19 (1.71)
Return on equity	0.28 (1.81)	0.10 (0.80)	-0.00 (-0.03)	-0.04 (-0.48)
Book-to-price ratio	0.29 (1.47)	0.19 (1.32)	0.44 (3.17)	0.37 (2.78)
Momentum	-0.44 (-1.21)	-0.32 (-1.29)	0.27 (1.08)	0.33 (1.45)
Growth in earnings per share	0.12 (1.08)	0.05 (0.64)	-0.02 (-0.35)	-0.04 (-0.65)
1-Month reversal	0.21 (0.87)	0.12 (0.65)	0.16 (0.83)	0.12 (0.67)
Low volatility	-0.10 (-0.30)	0.16 (0.80)	-0.13 (-0.58)	0.05 (0.28)
Earnings yield	-0.38 (-1.94)	-0.17 (-1.34)	0.09 (0.86)	0.25 (2.52)
Accounting accruals	0.21 (2.24)	0.15 (1.93)	0.29 (3.77)	0.28 (3.39)
Intercept	0.06 (0.09)	0.01 (0.02)	-0.31 (-0.66)	-0.35 (-0.80)

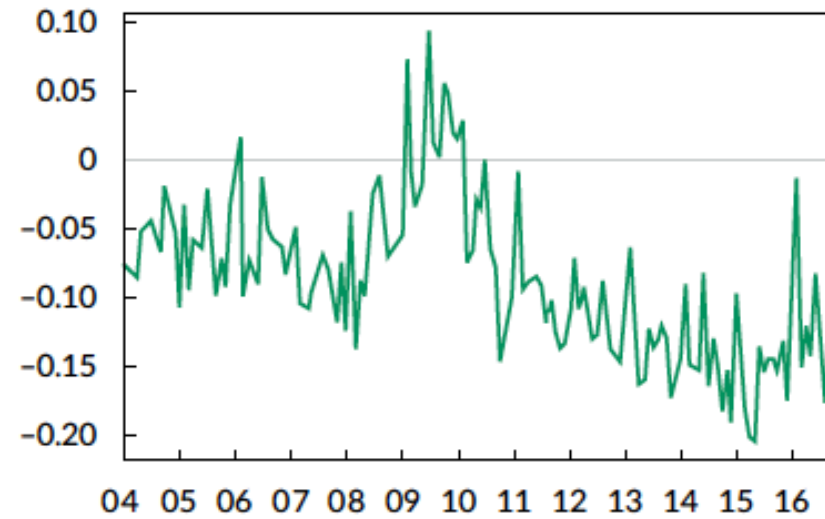
Table 6. Long and Short Portfolios (t-statistics in parentheses)

Measure	US Zone		ROW	
	Long	Short	Long	Short
Excess return	1.90 (3.54)	-0.03 (-0.05)	1.50 (4.18)	-0.13 (-0.33)
Fama-French-Carhart four-factor alpha	1.13 (4.43)	-0.96 (-3.28)	0.95 (4.72)	-0.69 (-2.97)

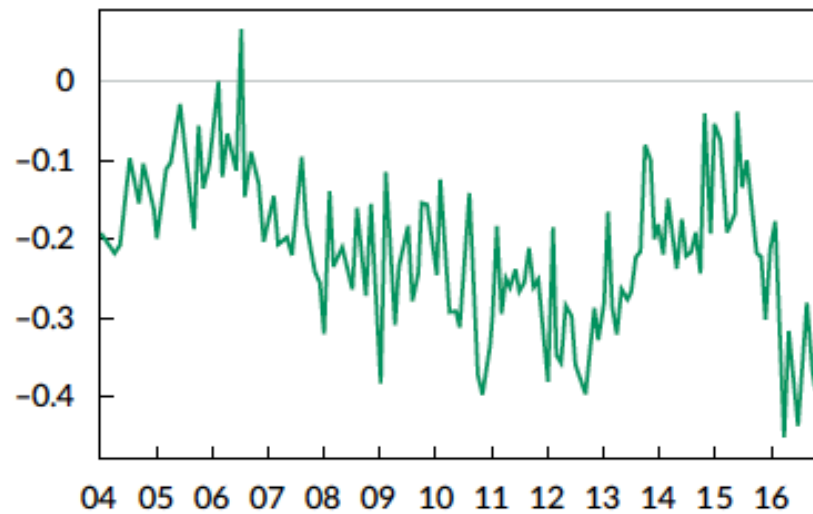
A. Correlation with HML



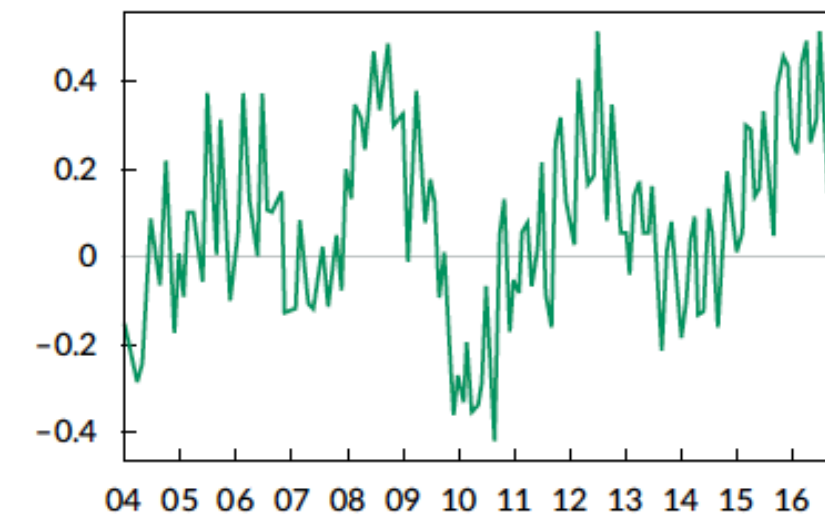
B. Correlation with SMB



C. Correlation with MKT



D. Correlation with MOM



Pitfalls

- The algos only see the data they are presented with: no solution for overfit inputs
- If there is look-ahead bias any algo worth its salt will suggest to trade on that!
- How is internal and external data checked for errors?
- Does your data team apply the same process in production and on historical data?
- Are costs accounted for? Cross-sectional variation in hidden costs? Asymmetric between long and short?

Conclusion

- Feature Engineering can Increase the signal-to-noise ratio before any estimation takes place
- Forecast combinations can lead to more signal and less noise because a portion of errors cancels out
- ML only effective if the past is informative about the future: Data needs to be point-in-time
- Inputs ideally should have low and comparable selection bias

Q&A



BAYSIANS
AGAINST
DISCRIMINATION

SUPPORT
VECTOR
MACHINES

REPEAL
POWER
LAWS

END
DUALITY
GAP

FREE
VARIABLES!

BAN
GENETIC
ALGORITHMS

Map Reduce
Map Reuse
Map Recycle