



Machine Learning for Data Analysis

Jesús Fernández-Villaverde¹ and Galo Nuño²

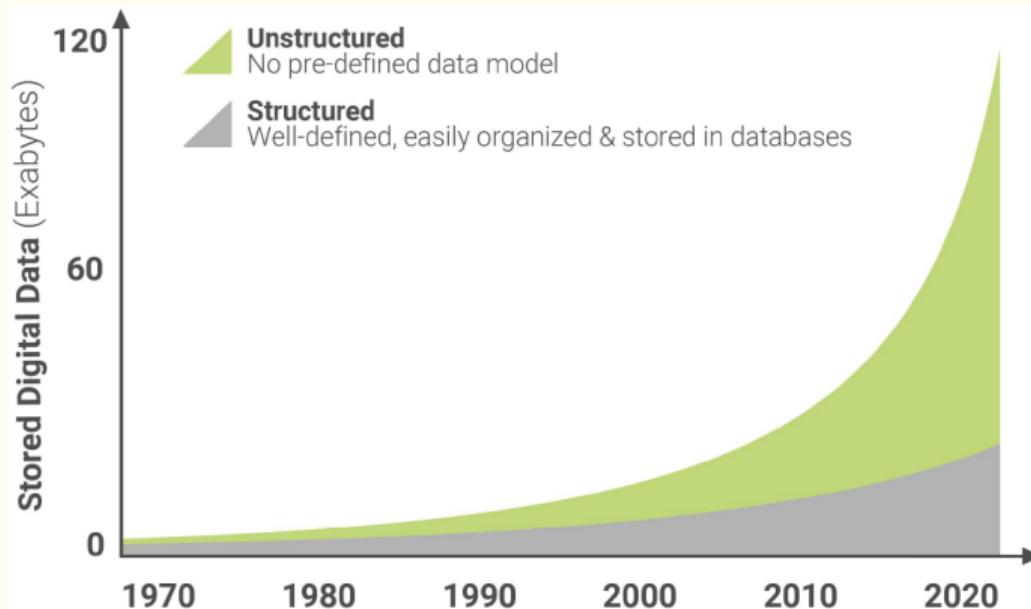
September 1, 2022

¹University of Pennsylvania

²Banco de España

New data

- Most important lesson for economists from data science: Everything is data.
- Unstructured data: Newspaper articles, business reports, congressional speeches, FOMC meetings transcripts, satellite data, photographs, audio, mobility, ...



Library data

Refine Your Search

Author

- [Galileo Galilei](#) (1208)
- [Galilei Galileo](#) (813)
- [Peiresc Nicolas C...](#) (46)
- [Meucci Ferdinando](#) (22)
- [Cicoli Andrea](#) (21)
- [Show more ...](#)

Year

- [1668](#) (117)
- [1966](#) (61)
- [1864](#) (104)
- [1899](#) (65)
- [1610](#) (62)
- [Show more ...](#)

Language

- [Italian](#) (1740)
- [Undetermined](#) (388)
- [English](#) (386)
- [Latin](#) (351)
- [French](#) (273)
- [Show more ...](#)

Content

- [Biography](#) (96)
- [Fiction](#) (10)
- [Non-Fiction](#) (3610)

Audience

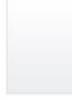
- [Juvenile](#) (8)
- [Non-Juvenile](#) (3612)

Topic

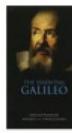
- [Physical Sciences](#) (407)
- [Language, Linguistics](#) (40)
- [Philosophy & Religion](#) (25)

Select All Clear All

Save to: Save

1. 

[Galileo on the world systems : a new abridged translation and guide](#)
by Galileo Galilei; Maurice A Finocchiaro
 eBook : Document [View all formats and languages](#) 
Language: English
Publisher: Berkeley : University of California Press, ©1997.
[View all editions](#) 

2. 

[The essential Galileo](#)
by Galileo Galilei; Maurice A Finocchiaro;
 Print book [View all formats and languages](#) 
Language: English
Publisher: Indianapolis, Ind. : Hackett Pub., 2008
[View all editions](#) 

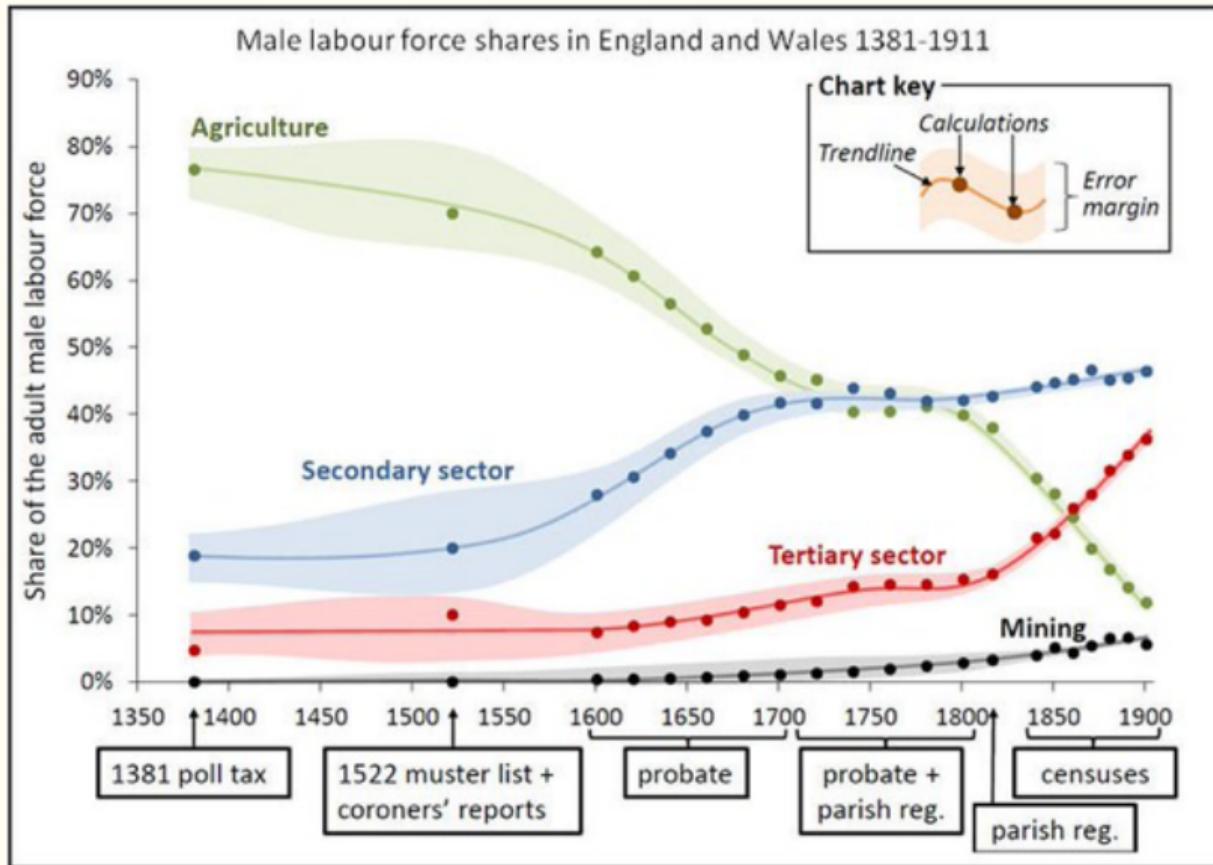
3. 

[Dialogue concerning the two chief world systems, Ptolemaic and Copernican](#)
by Galileo Galilei; Stillman Drake; Dava Sobel; Albert Einstein; Folio Society (London, England)
 Print book [View all formats and languages](#) 
Language: English
Publisher: London : Folio Society, 2013.
[View all editions](#) 

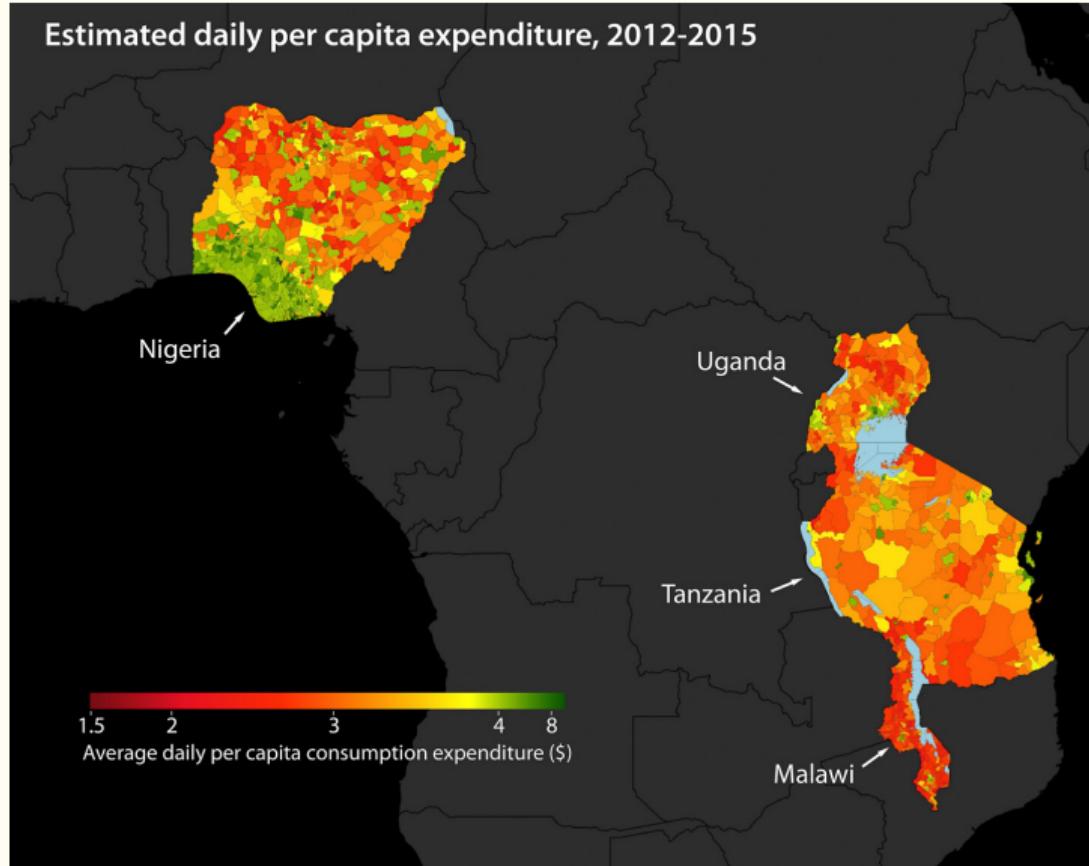
4. 

[Sidereus nuncius or The sidereal messenger](#)
by Galileo Galilei; Albert Van Helden
 Print book [View all formats and languages](#) 
Language: English
Publisher: Chicago The University of Chicago Press [2015]
[View all editions](#) 

Parish and probate data



Satellite imagery



Luminosity

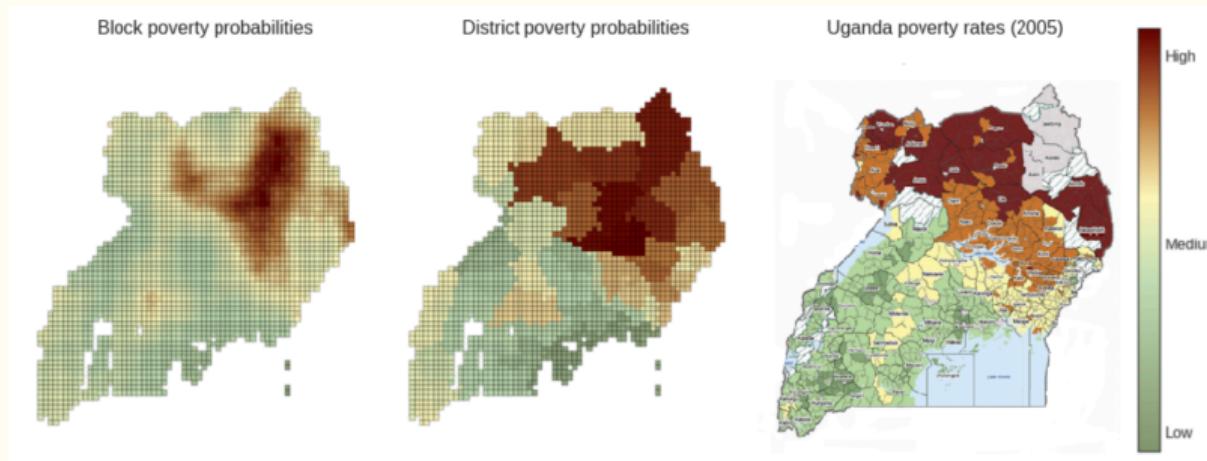


Figure 3: **Left:** Predicted poverty probabilities at a fine-grained 10km × 10km block level. **Middle:** Predicted poverty probabilities aggregated at the district-level. **Right:** 2005 survey results for comparison (World Resources Institute 2009).

Xie et. al. (2016)

Cell use

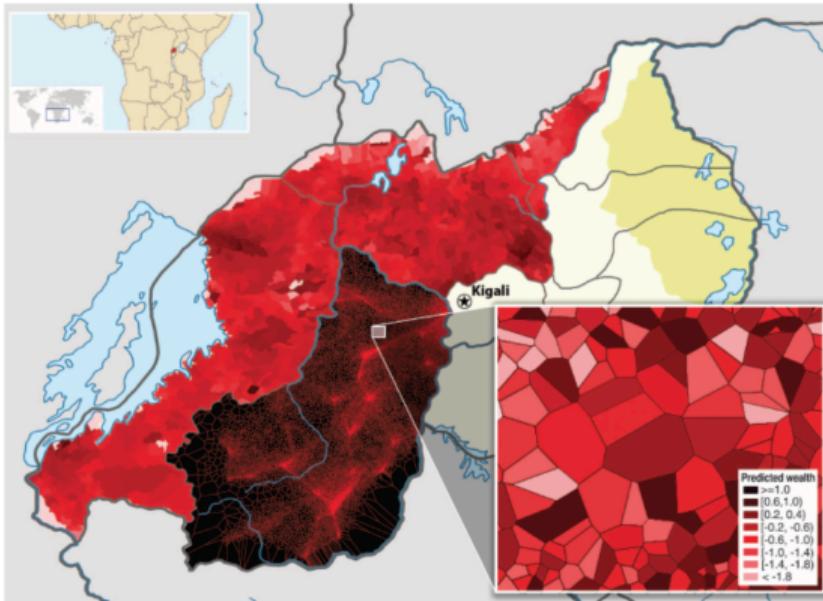


Fig. 2. Construction of high-resolution maps of poverty and wealth from call records. Information derived from the call records of 1.5 million subscribers is overlaid on a map of Rwanda. The northern and western provinces are divided into cells (the smallest administrative unit of the country), and the cell is shaded according to the average (predicted) wealth of all mobile subscribers in that cell. The southern province is overlaid with a Voronoi division that uses geographic identifiers in the call data to segment the region into several hundred thousand small partitions. (**Bottom right inset**) Enlargement of a 1-km² region near Kiyonza, with Voronoi cells shaded by the predicted wealth of small groups (5 to 15 subscribers) who live in each region.

Blumenstock et. al. (2015)

TABLE I
MOST PARTISAN PHRASES FROM THE 2005 CONGRESSIONAL RECORD^a

Panel A: Phrases Used More Often by Democrats		
<i>Two-Word Phrases</i>		
private accounts	Rosa Parks	workers rights
trade agreement	President budget	poor people
American people	Republican party	Republican leader
tax breaks	change the rules	Arctic refuge
trade deficit	minimum wage	cut funding
oil companies	budget deficit	American workers
credit card	Republican senators	living in poverty
nuclear option	privatization plan	Senate Republicans
war in Iraq	wildlife refuge	fuel efficiency
middle class	card companies	national wildlife
<i>Three-Word Phrases</i>		
veterans health care	corporation for public	cut health care
congressional black caucus	broadcasting	civil rights movement
VA health care	additional tax cuts	cuts to child support
billion in tax cuts	pay for tax cuts	drilling in the Arctic National
credit card companies	tax cuts for people	victims of gun violence
security trust fund	oil and gas companies	solvency of social security
social security trust	prescription drug bill	Voting Rights Act
privatize social security	caliber sniper rifles	war in Iraq and Afghanistan
American free trade	increase in the minimum wage	civil rights protections
central American free	system of checks and balances	credit card debt
	middle class families	

TABLE I—Continued

Panel B: Phrases Used More Often by Republicans		
<i>Two-Word Phrases</i>		
stem cell	personal accounts	retirement accounts
natural gas	Saddam Hussein	government spending
death tax	pass the bill	national forest
illegal aliens	private property	minority leader
class action	border security	urge support
war on terror	President announces	cell lines
embryonic stem	human life	cord blood
tax relief	Chief Justice	action lawsuits
illegal immigration	human embryos	economic growth
date the time	increase taxes	food program
<i>Three-Word Phrases</i>		
embryonic stem cell	Circuit Court of Appeals	Tongass national forest
hate crimes legislation	death tax repeal	pluripotent stem cells
adult stem cells	housing and urban affairs	Supreme Court of Texas
oil for food program	million jobs created	Justice Priscilla Owen
personal retirement accounts	national flood insurance	Justice Janice Rogers
energy and natural resources	oil for food scandal	American Bar Association
global war on terror	private property rights	growth and job creation
hate crimes law	temporary worker program	natural gas natural
change hearts and minds	class action reform	Grand Ole Opry
global war on terrorism	Chief Justice Rehnquist	reform social security

^aThe top 60 Democratic and Republican phrases, respectively, are shown ranked by χ^2_{pl} . The phrases are classified as two or three word after dropping common “stopwords” such as “for” and “the.” See Section 3 for details and see Appendix B (online) for a more extensive phrase list.

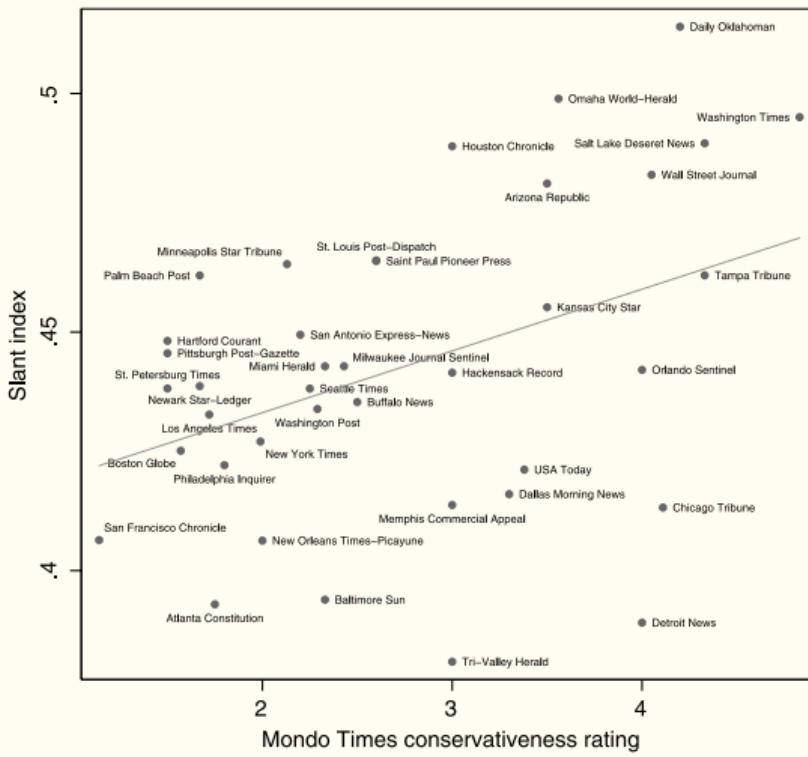
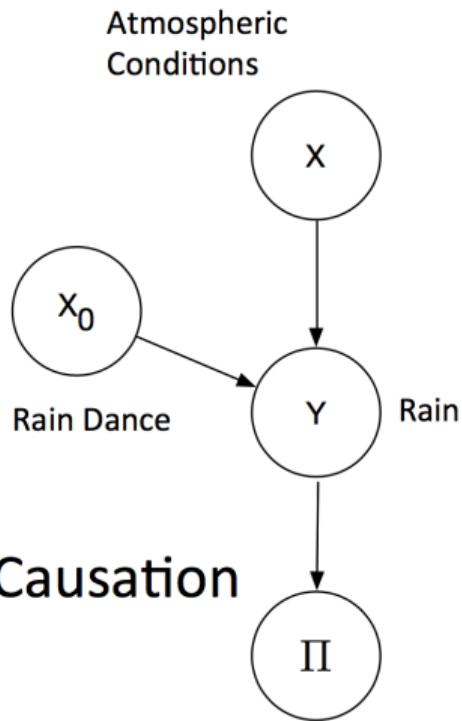


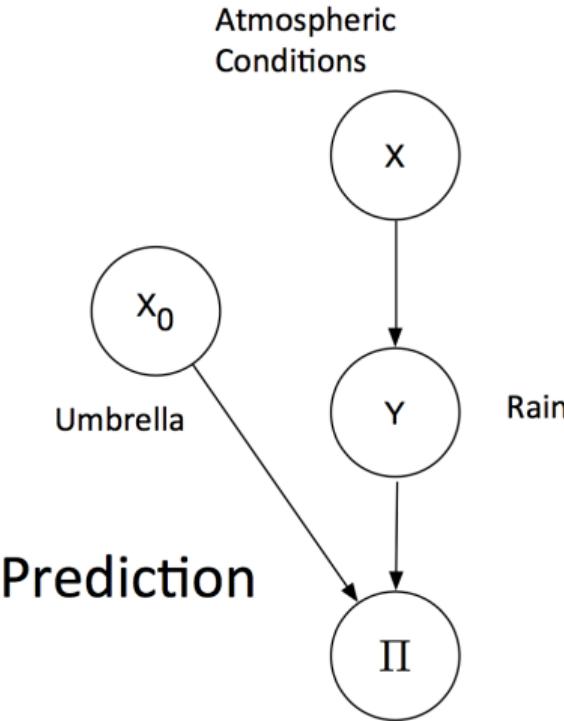
FIGURE 1.—Language-based and reader-submitted ratings of slant. The slant index (y axis) is shown against the average Mondo Times user rating of newspaper conservativeness (x axis), which ranges from 1 (liberal) to 5 (conservative). Included are all papers rated by at least two users on Mondo Times, with at least 25,000 mentions of our 1000 phrases in 2005. The line is predicted slant from an OLS regression of slant on Mondo Times rating. The correlation coefficient is 0.40 ($p = 0.0114$).

- A more general point \Rightarrow role of causality in economics:
 1. Counterfactuals.
 2. Welfare.
 3. General equilibrium effects.
 4. New changes.
 5. Less data.
- Another example by Athey (2017): hotel prices and occupancy rates. In the data, prices and occupancy rates are strongly positively correlated, but what is the expected impact of a hotel raising its prices on a given day?



Causation

Experiments



Prediction

Machine Learning

Unsupervised learning

Unsupervised learning

- Use a sample:

$$\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$$

to:

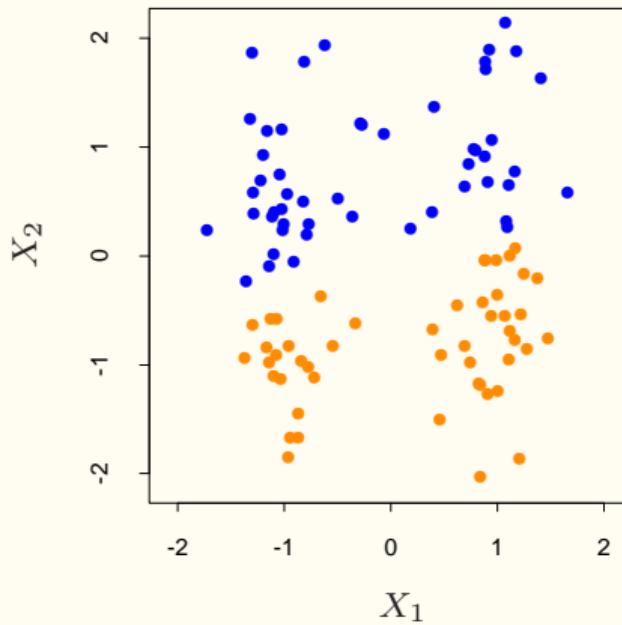
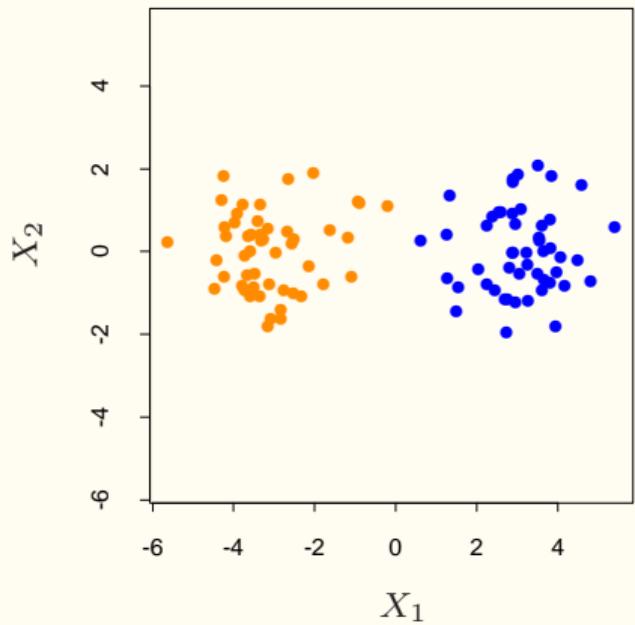
1. Group observations in interesting patterns.
2. Describe most important sources of variation in the data.
3. Dimensionality reduction.

- Example: what can we learn about the loan book of a bank without imposing too much a priori structure?

- More concretely, we search for:

$$p(\mathbf{x}_i | \theta)$$

- Clustering and association rules.



Cluster discovering

1. Select K clusters

$$K^* = \operatorname{argmax}_K p(K|\mathcal{D})$$

2. Assign each observation to a cluster

$$z_t^* = \operatorname{argmax}_k p(z_i = k | \mathbf{x}_i, \mathcal{D})$$

3. A common method to pin down K is the silhouette. For each observation i , we compute:

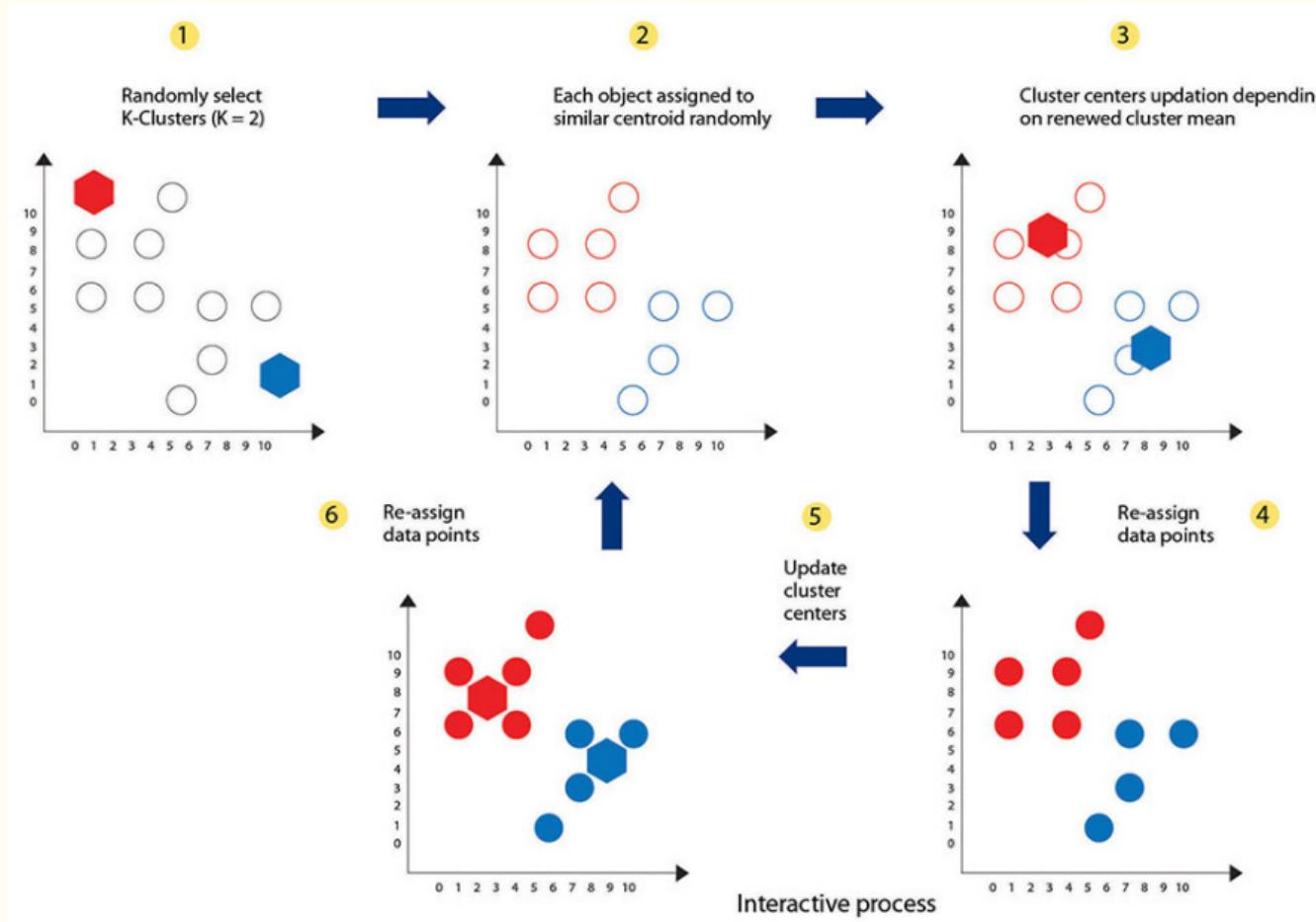
$$s_i = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

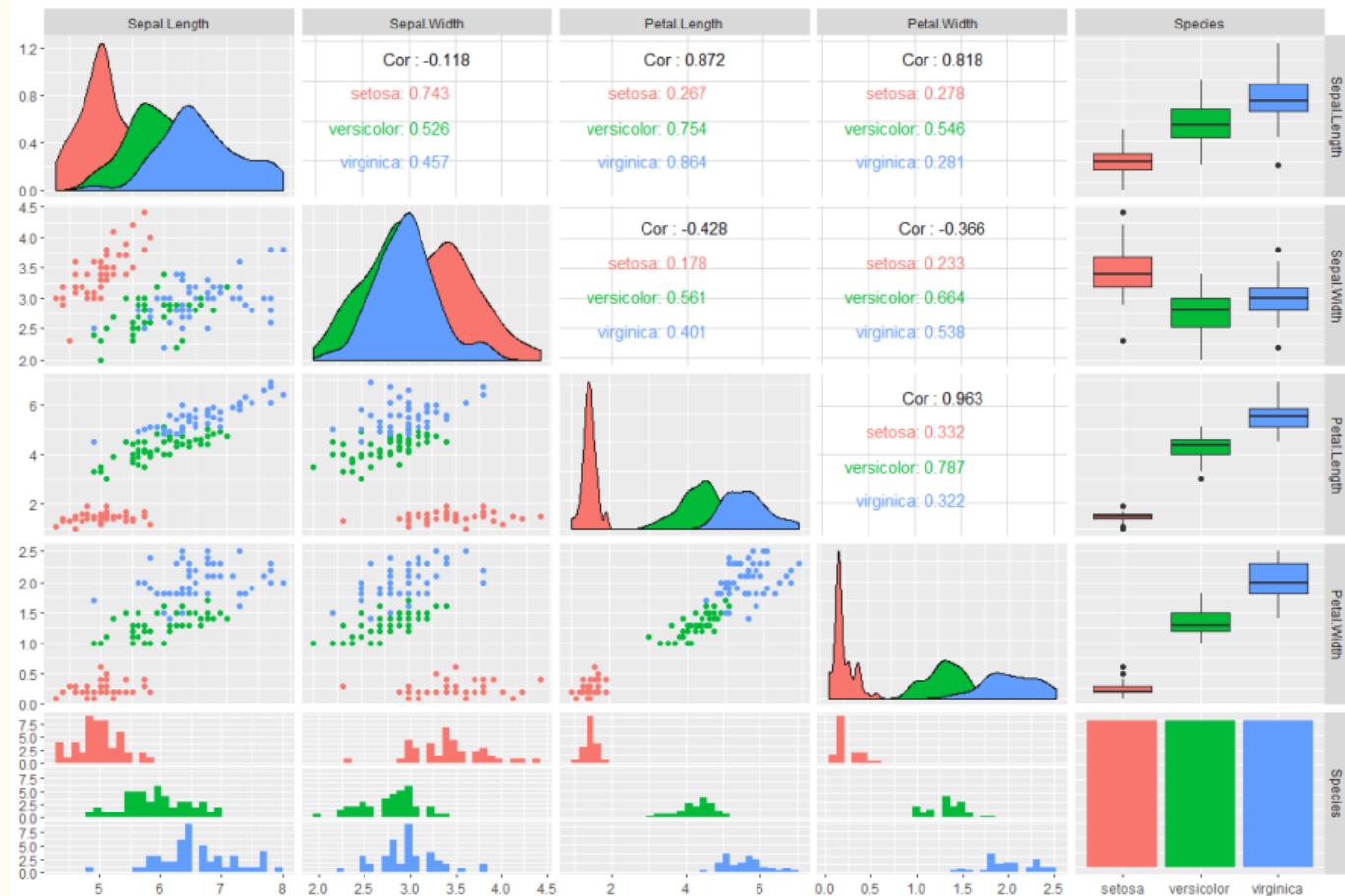
where $a(i)$ is the average distance between i and all other members of the cluster while $b(i)$ is the minimum distance between i and all other members of another cluster.

- K-means clustering by **Steinhaus (1957)**

$$\operatorname{argmax}_S \sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|^2$$

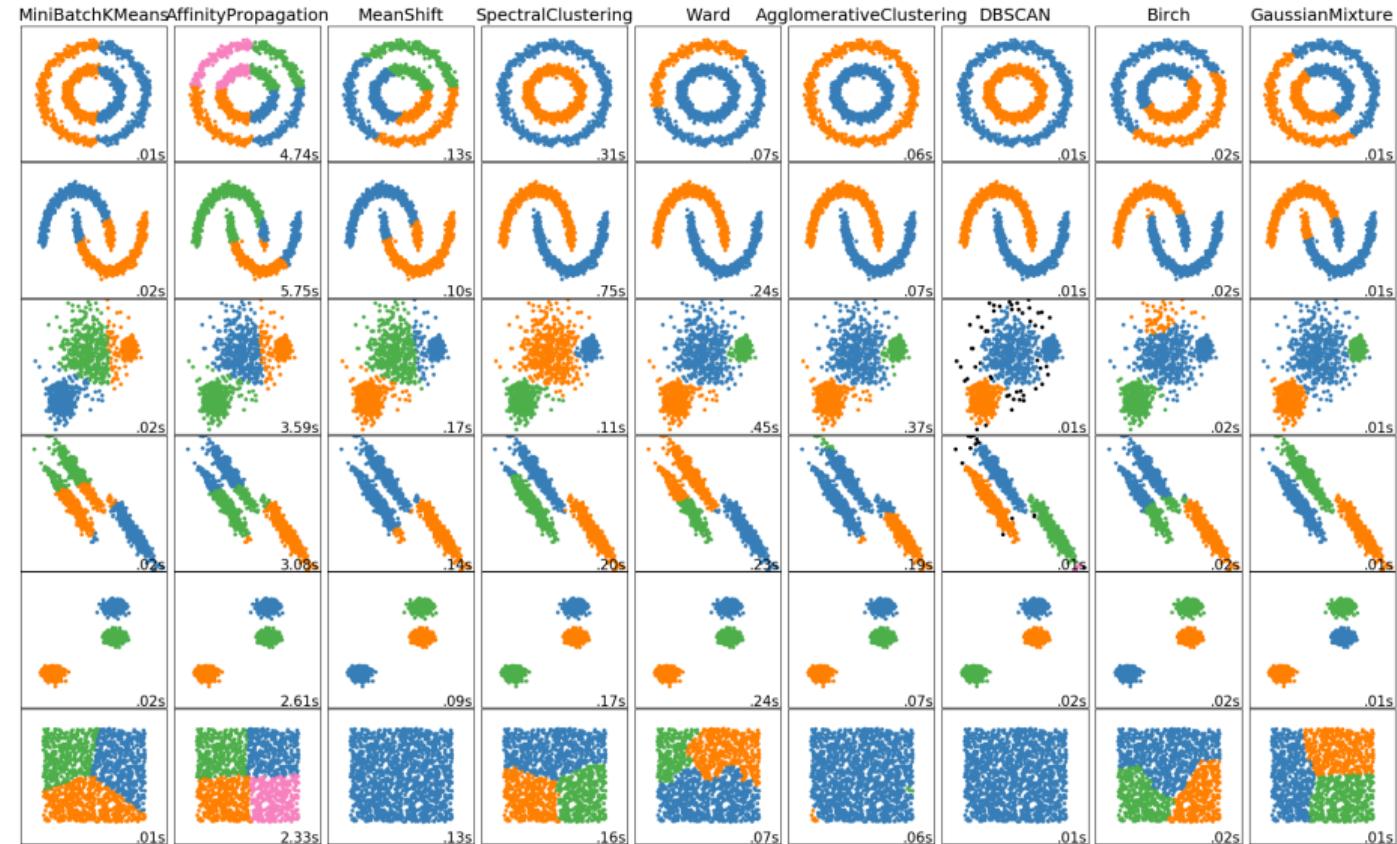
- It requires an iterative algorithm for implementation **Lloyd (1957)**.
- Related variations:
 1. k-medians \Rightarrow uses medians computed through the Taxicab geometry.
 2. k-medoids \Rightarrow minimizes a sum of pairwise dissimilarities.
 3. k-SVD.

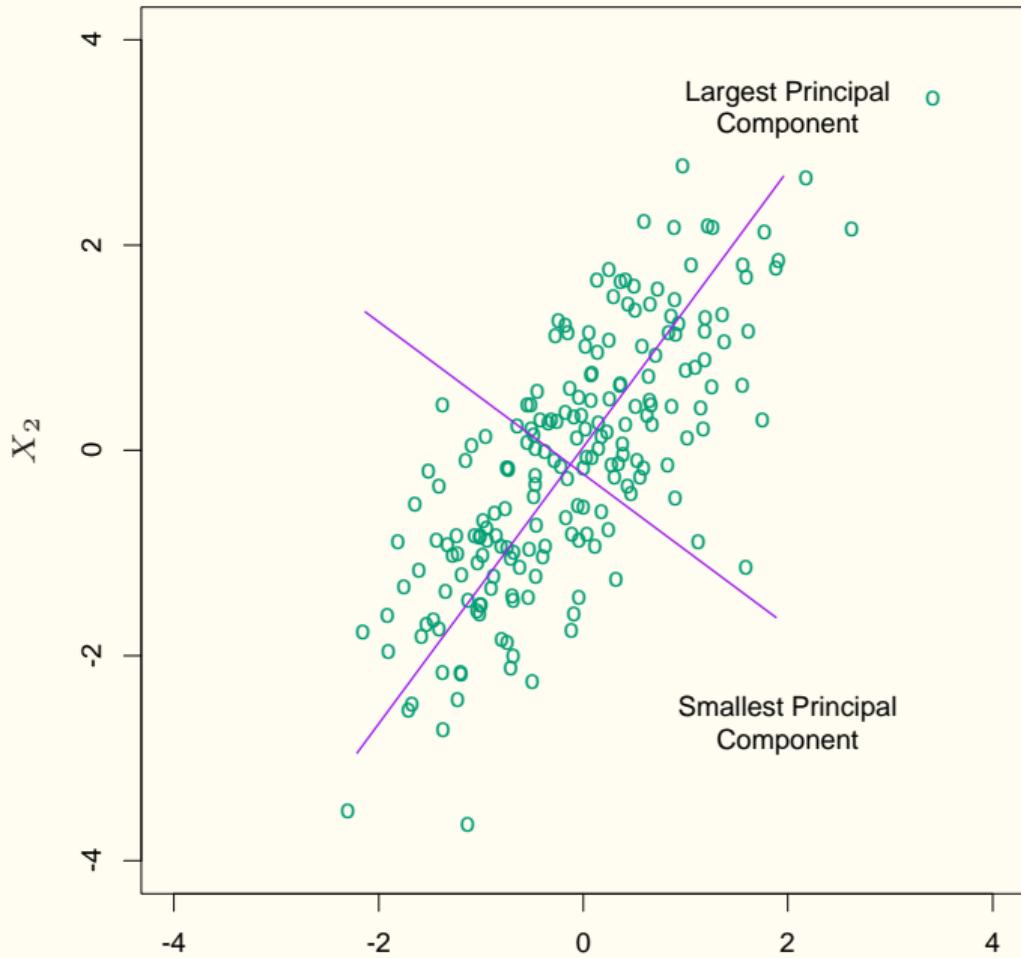




Other algorithms

- Other clustering methods:
 1. Agglomerative clustering.
 2. DBSCAN.
 3. Birch.
- Principal component analysis.
- Density estimation.
- Gaussian mixture models.
- Association rules and the Apriori algorithm ([Agrawal and Srikant, 1994](#)).





Supervised learning

Supervised learning, I

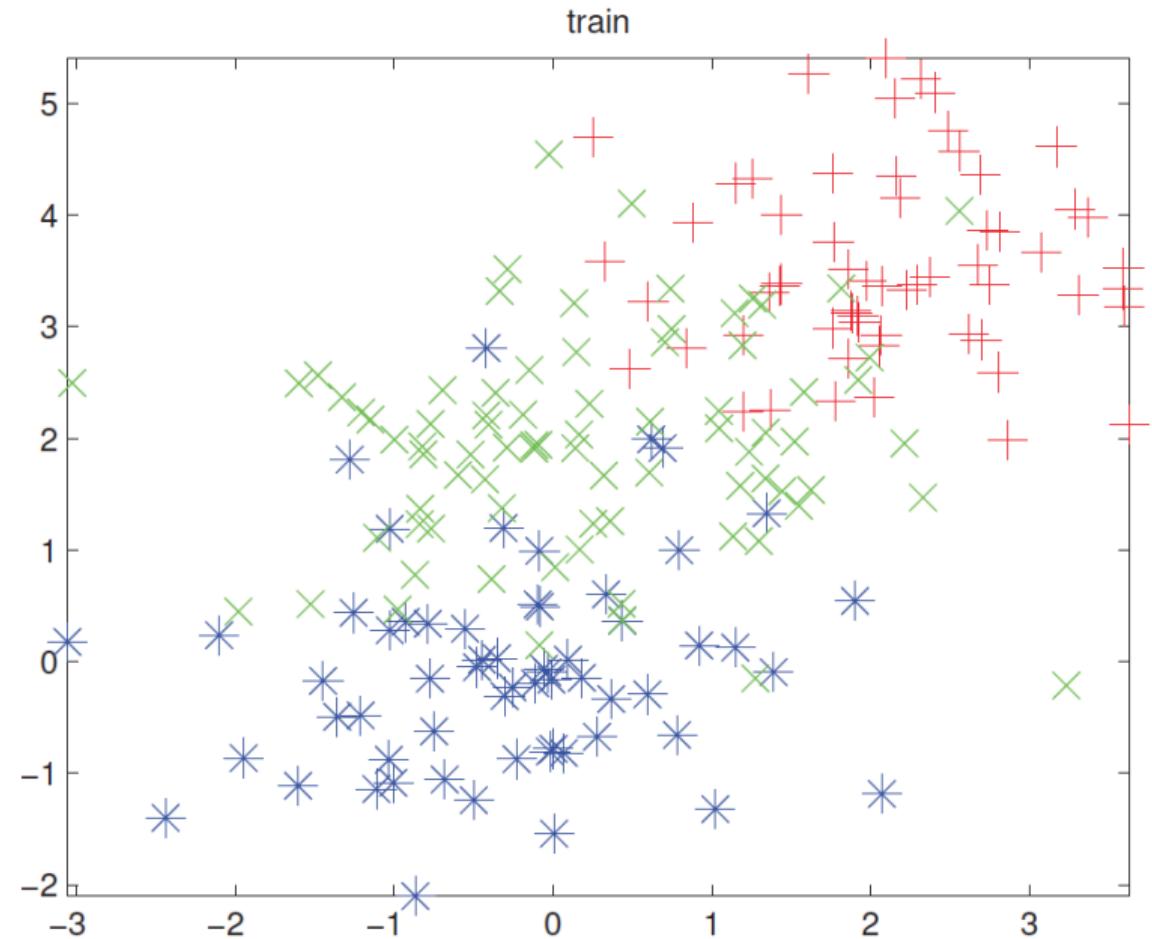
- Use a training sample

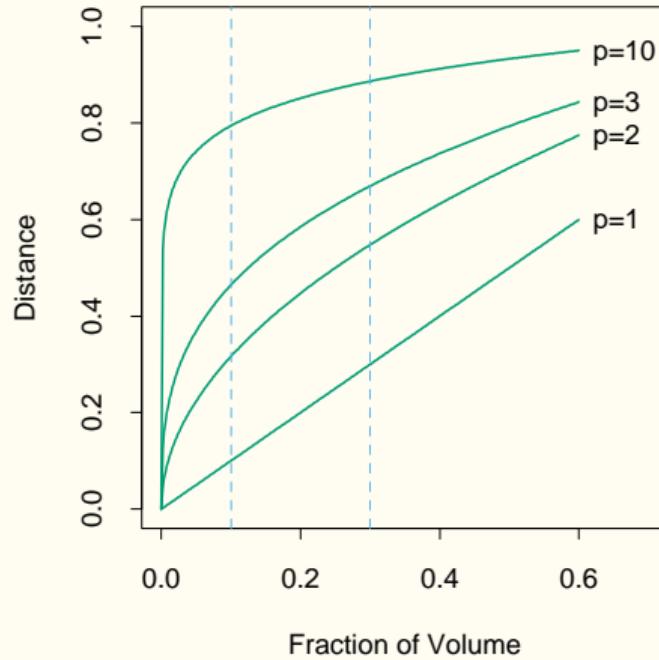
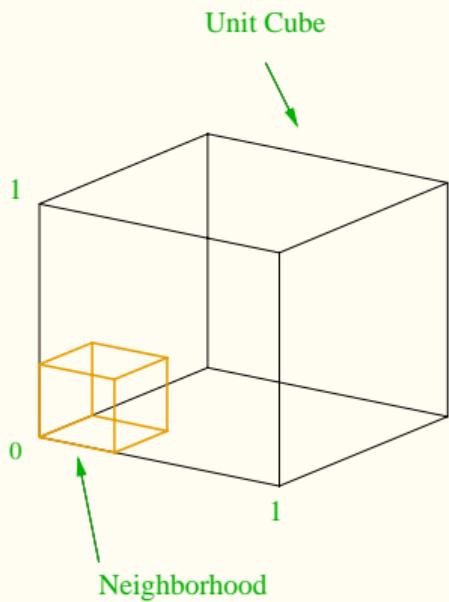
$$\mathcal{D} = \{y_i, \mathbf{x}_i\}_{i=1}^N$$

to classify new observations

$$\mathcal{F} = \{y_i, \mathbf{x}_i\}_{i=N+1}^M$$

- \mathbf{x}_i may potentially belong to a high-dimensional space.
- Standard regression in econometrics is a particular example when y_i is a continuous variable.
- But often, in machine learning, y_i is categorical (classification problem).
- Example: classify loans of a bank into good and bad risks.
- Caveat: it can be highly labor intensive!







Supervised learning, II

- Problem of function approximation:

$$y = f(\mathbf{x})$$

when we know next to nothing about the function $f(\cdot)$.

- Linked with projection and perturbation methods.
- Also, linked with traditional econometrics.
- Then:

$$\hat{y} = \hat{f}(\mathbf{x}) = \arg \max_{c \in C} p(y = c | \mathbf{x}, \mathcal{D})$$

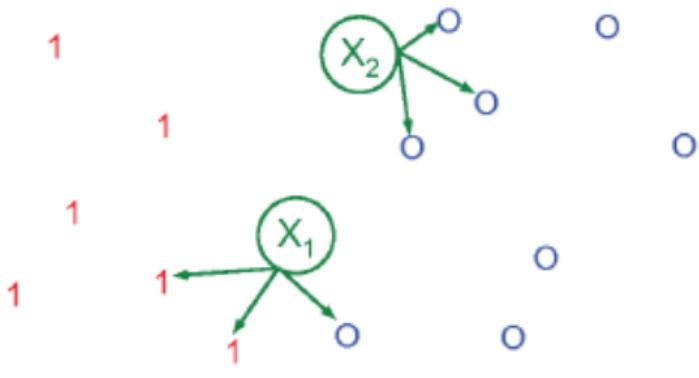
- In Bayesian, this is the mode of the posterior, but in machine learning it is called a MAP estimate (maximum a posteriori).

Example I: K-nearest algorithm

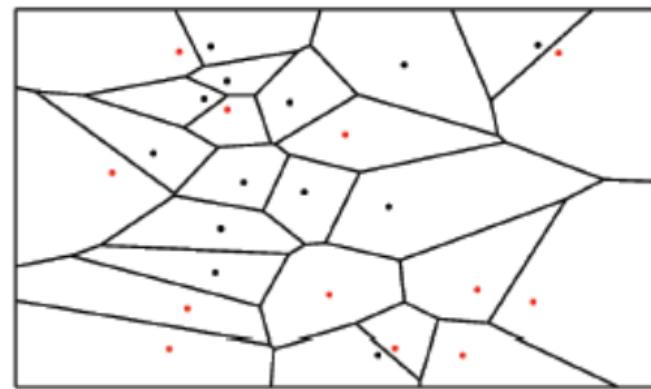
- We look at the K -nearest points to x_i in the training sample.
- Formally:

$$p(y = c | \mathbf{x}, \mathcal{D}, K) = \frac{1}{K} \sum_{i \in N_K(\mathbf{x}, K)} \mathbb{I}(y_i = c)$$

- It generates a Voronoi tessellation.
- Simple and intuitive, but it suffers from curse of dimensionality.



(a)



(b)

Bias-variance trade-off

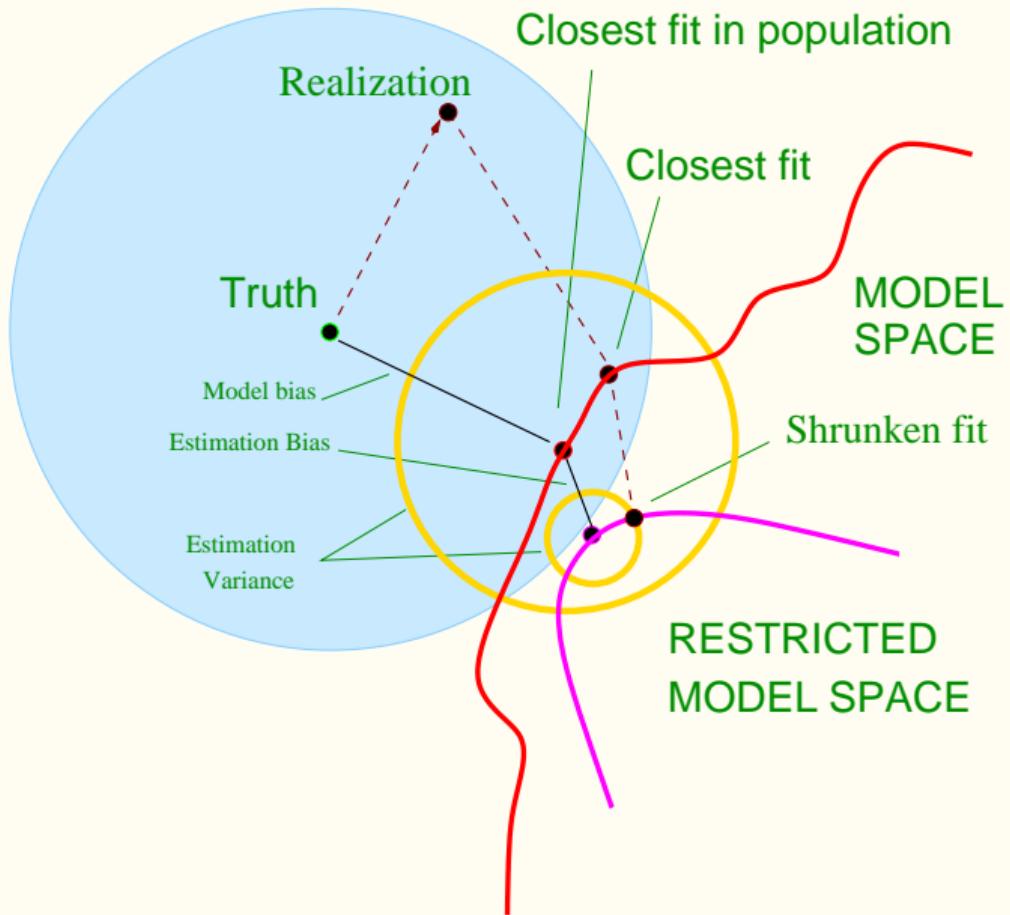
- How do we select K ? Minimize misclassification rate on a validation set.
- Bias-variance trade-off is about underfitting vs. overfitting.
- We have:

$$y = f(\mathbf{x}) = \hat{f}(\mathbf{x}) + \varepsilon, \varepsilon \sim ID(0, \sigma^2)$$

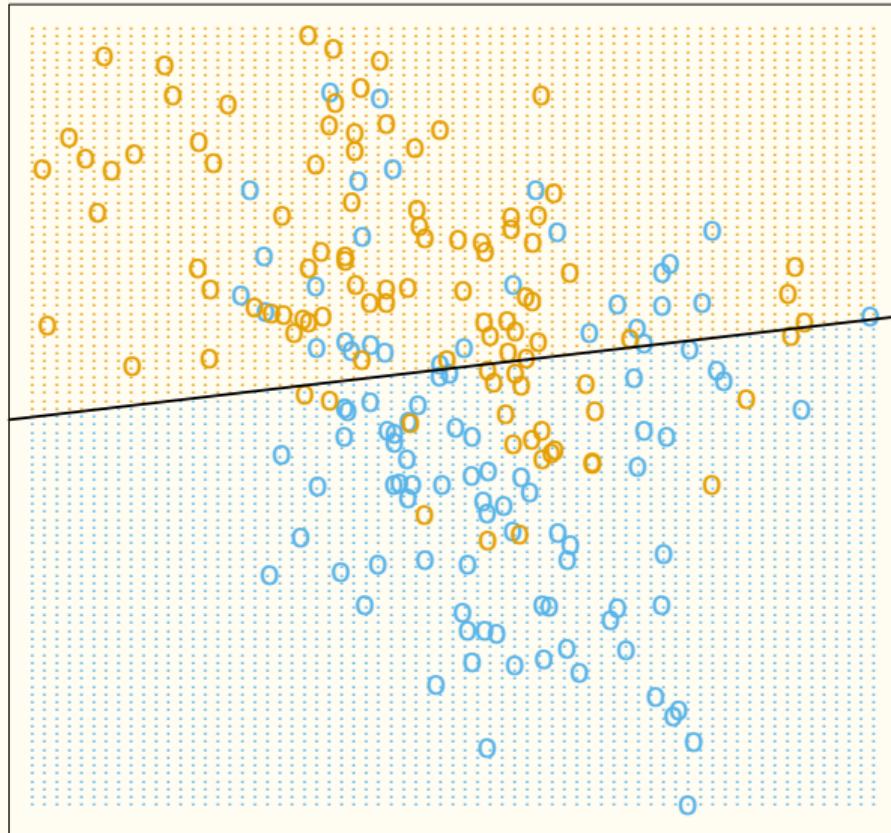
- Then:

$$\begin{aligned}\mathbb{E}(y - \hat{f}(\mathbf{x}))^2 &= \mathbb{E}(\hat{f}(\mathbf{x}) - f(\mathbf{x}))^2 \\ &= (\mathbb{E}\hat{f}(\mathbf{x}) - f(\mathbf{x}))^2 + \mathbb{E}(\hat{f}(\mathbf{x}) - \mathbb{E}\hat{f}(\mathbf{x}))^2 \\ &= Bias^2 + var\end{aligned}$$

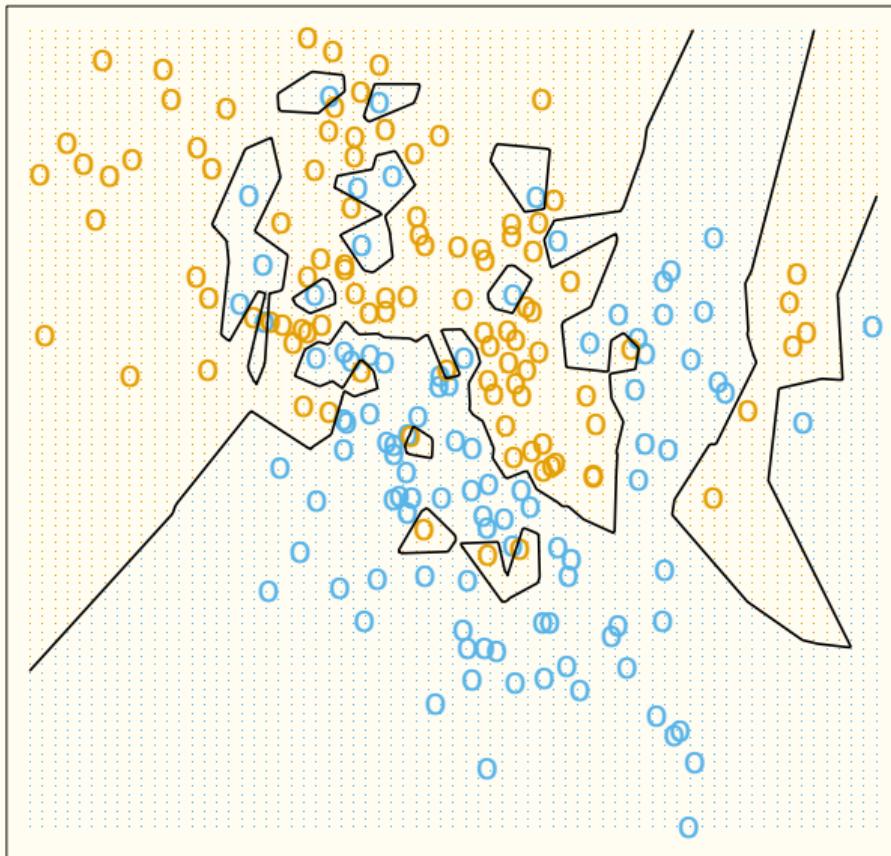
- Traditional econometrics obsession with BLUE is mysterious.



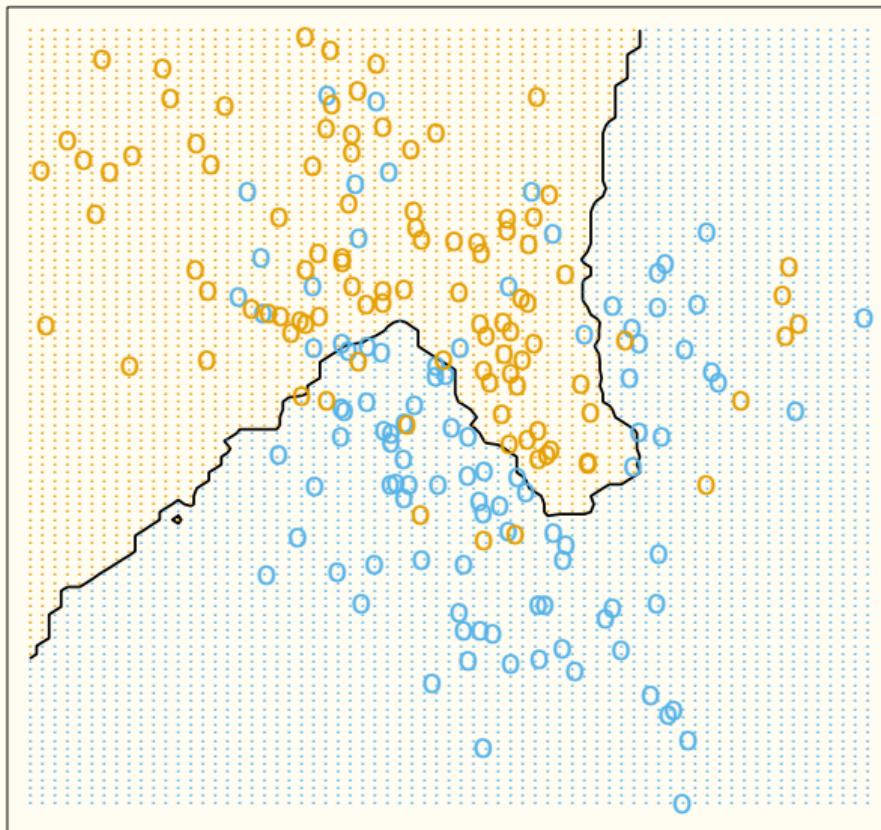
Linear Regression of 0/1 Response



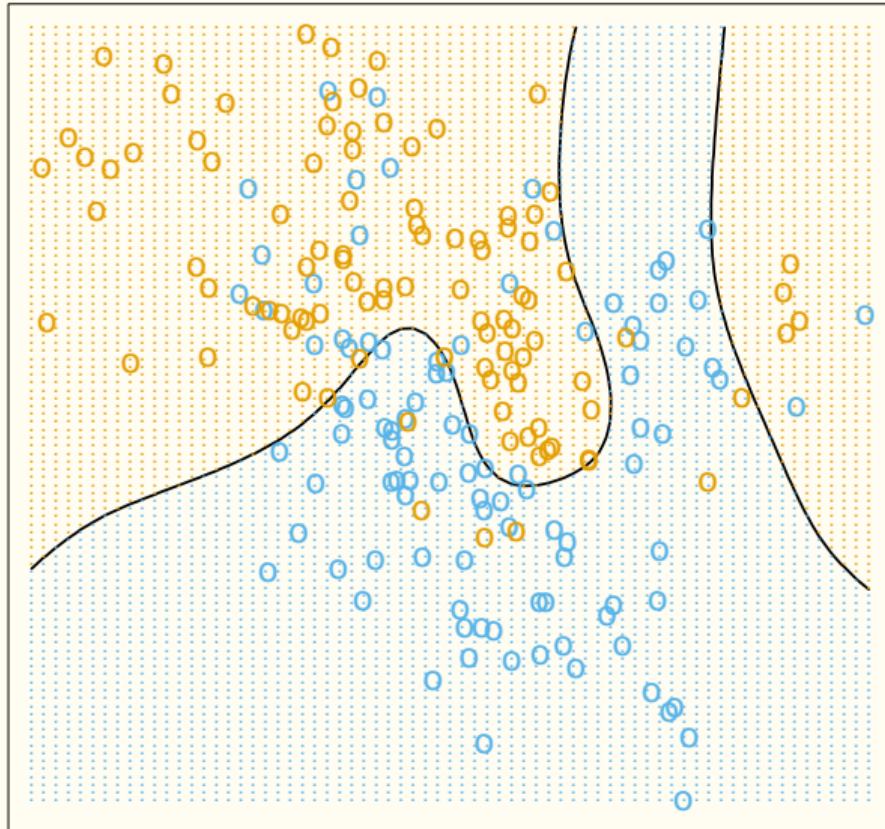
1–Nearest Neighbor Classifier

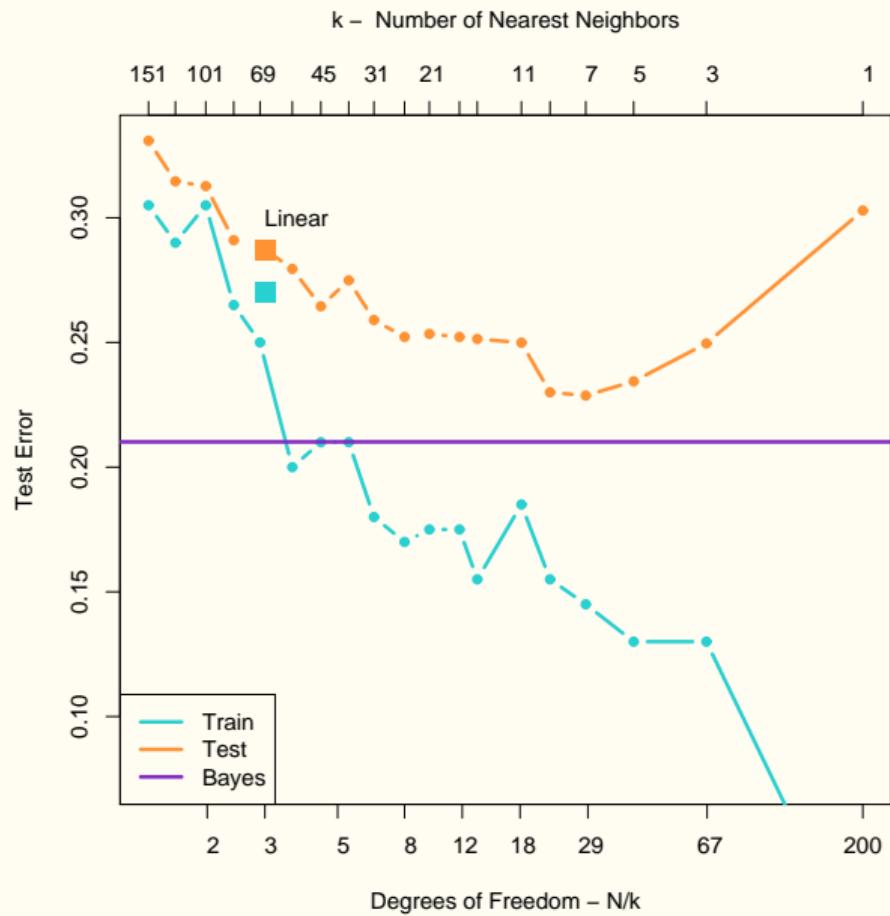


15-Nearest Neighbor Classifier

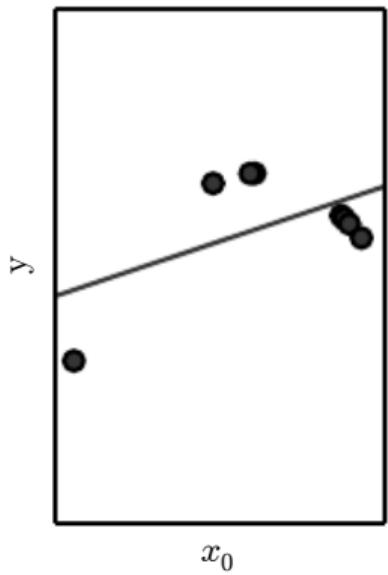


Bayes Optimal Classifier

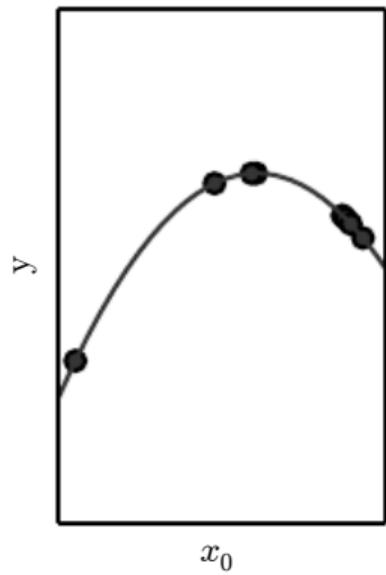




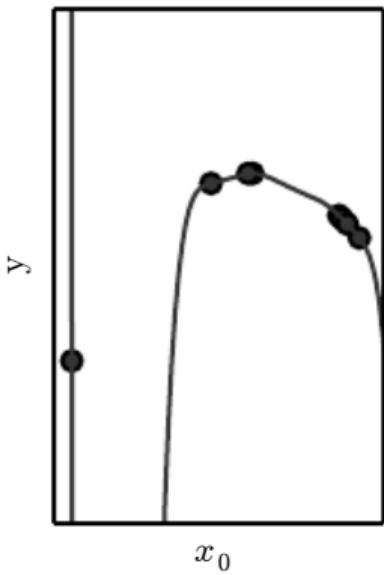
Underfitting



Appropriate capacity



Overfitting



FULL DATA



TRAIN



TEST



Example II: Regularized regression

- Standard linear regression:

$$y_i = \mathbf{x}'_i \beta + \varepsilon_t$$

- When number of regressors is 3 or more, OLS is not admissible (Stein's phenomenon).
- In particular, when the number of regressors is large, forecasting properties are poor.
- Zvi Griliches: "never trust OLS with more than five regressors."
- This problem is becoming more prominent as we get larger datasets with thousands of covariates.
- Also relevant when we have multiple potential IVs.
- Note difference between problem of learning about coefficients and forecasting!

A general formulation

- Solution: regularization (aka as penalization)
 1. Shrink estimates towards zero.
 2. Sparse representation of non-zero parameters.
- Long tradition in econometrics: James–Stein estimator.
- Estimator:

$$\beta_{reg} = \min_{\beta} \sum_{i=1}^N (y_i - \mathbf{x}'_i \beta)^2 + \lambda \|\beta\|_p^p$$

where

$$\|\beta\|_p = \left(\sum_{i=1}^K |\beta|^p \right)^{\frac{1}{p}}$$

CORE PRINCIPLES IN RESEARCH

JORGE CHAM © 2009



OCCAM'S RAZOR

"WHEN FACED WITH TWO POSSIBLE EXPLANATIONS, THE SIMPLER OF THE TWO IS THE ONE MOST LIKELY TO BE TRUE."



OCCAM'S PROFESSOR

"WHEN FACED WITH TWO POSSIBLE WAYS OF DOING SOMETHING, THE MORE COMPLICATED ONE IS THE ONE YOUR PROFESSOR WILL MOST LIKELY ASK YOU TO DO."

WWW.PHDCOMICS.COM

L_0 norm

- Estimator with $p = 0$ (i.e., the number of non-zero coefficients).
- With an appropriately selected λ , we essentially choose β according to AIC/BIC criteria, which have strong statistical foundations.
- Challenge: the optimization problem is difficult, and becomes computationally infeasible when number of coefficients is large.
- Two alternatives:
 1. *Stepwise forward regression*: start with no covariates; add whichever single covariate improves the fit most; continue until no covariate significantly improves fit.
 2. *Stepwise backward regression*: start with all covariates; drop single covariate that leads to the lowest reduction in fit; continue to drop until there is a significant reduction in fit.
- Once we select the relevant covariates, we can perform OLS, but standard errors are not correct since we do not condition on model selection.
- Thus, not very used in practice.

Ridge regression

- Estimator with $p = 2$.

- Close form:

$$\beta = (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}'\mathbf{Y}$$

- Estimators are shrink towards zero by factor $1 + \lambda$.
- It does not enforce sparsity and, thus, we do not achieve model selection.
- Bayesian interpretation $\lambda = \frac{\sigma^2}{\tau^2}$ where $\beta_{prior} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I})$ and $\varepsilon_t \sim \mathcal{N}(0, \sigma^2)$.

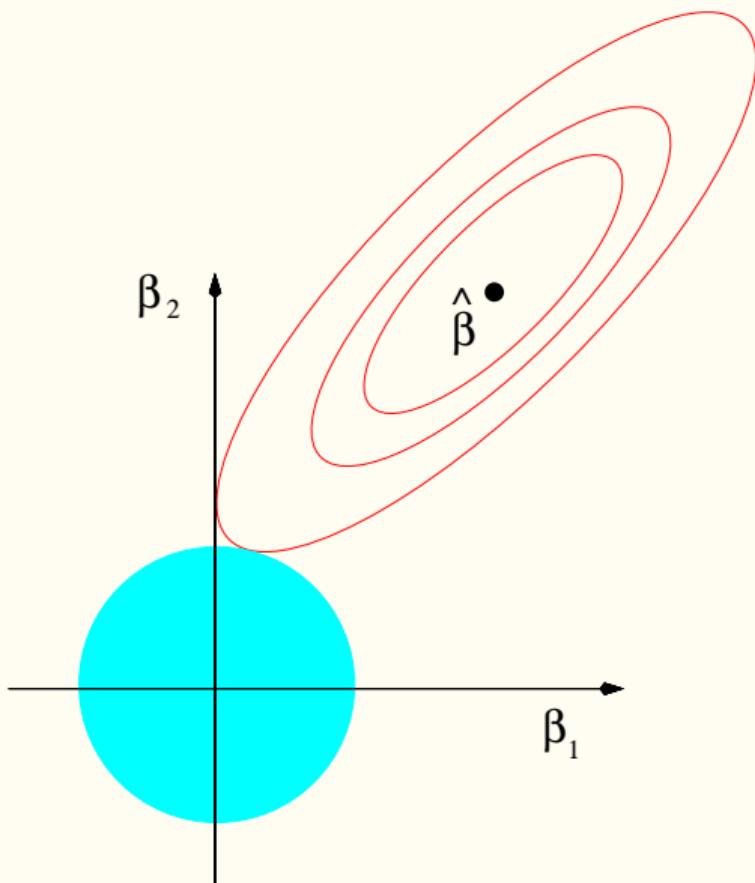
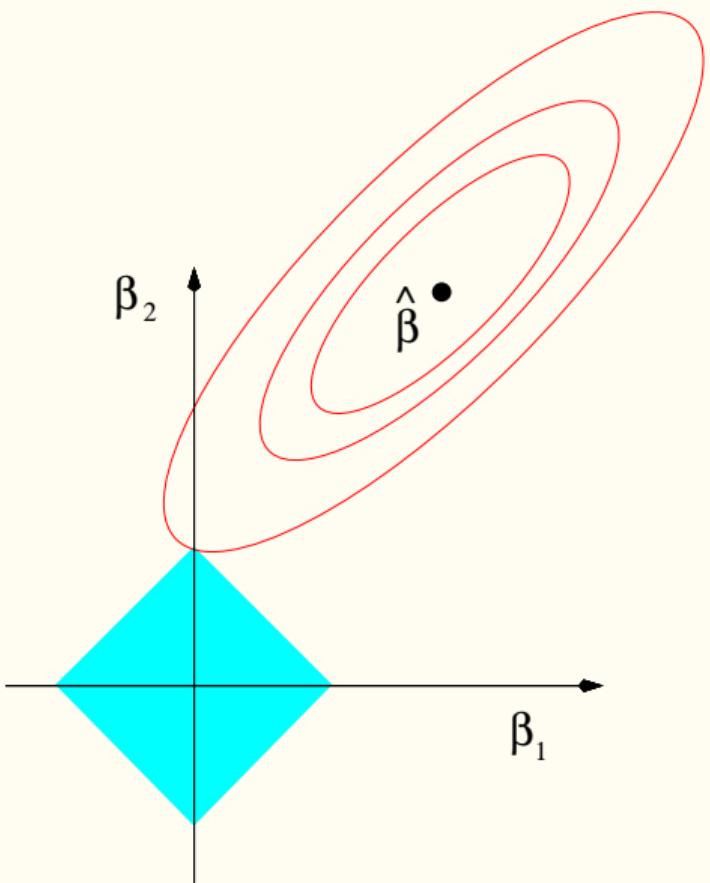
- Least absolute selection and shrinkage operator (**Tibshirani, 1996**):

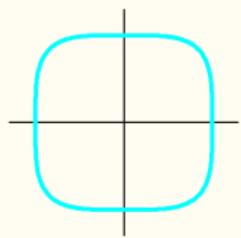
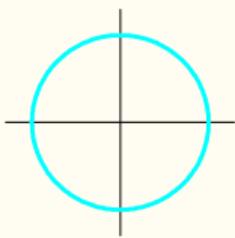
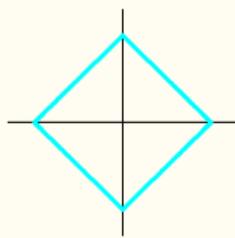
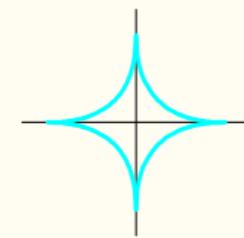
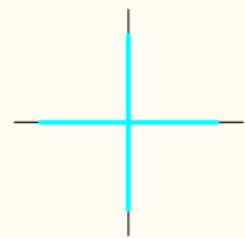
$$\min \sum_{i=1}^N (y_i - \mathbf{x}'_i \beta)^2 + \lambda \|\beta\|_1$$

- LASSO shrinks some coefficients exactly to zero and others towards zero (variation of relaxed lasso).
- Algorithms for its implementation are easy to code \Rightarrow package in R package `glmnet`.
- Bayesian interpretation: mode of posterior given a Laplace prior

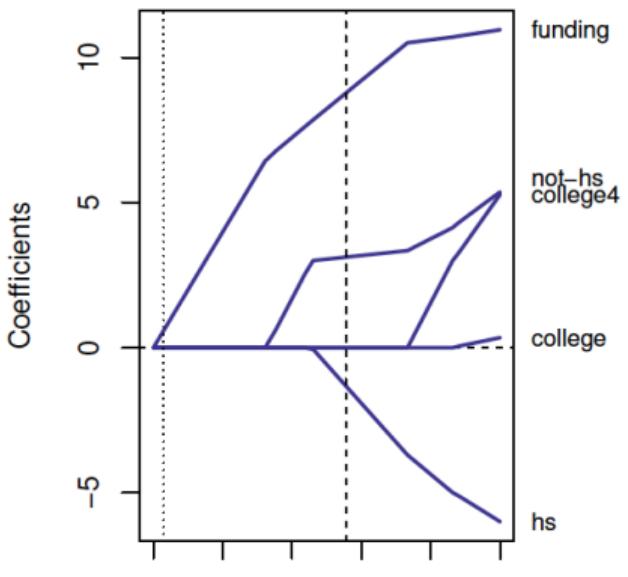
$$p(\beta) \propto \exp(-\lambda \|\beta\|_1)$$

- Bayesian alternative: spike-and-slab prior by **Mitchell and Beauchamp (1988)**.

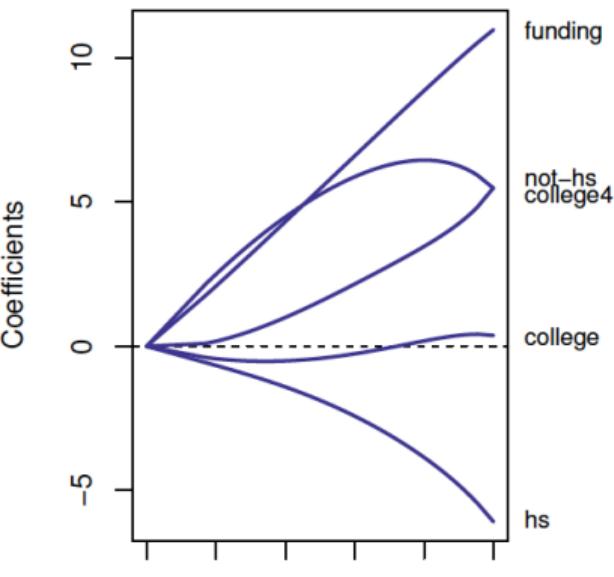


$q = 4$  $q = 2$  $q = 1$  $q = 0.5$  $q = 0.1$ 

Lasso



Ridge Regression

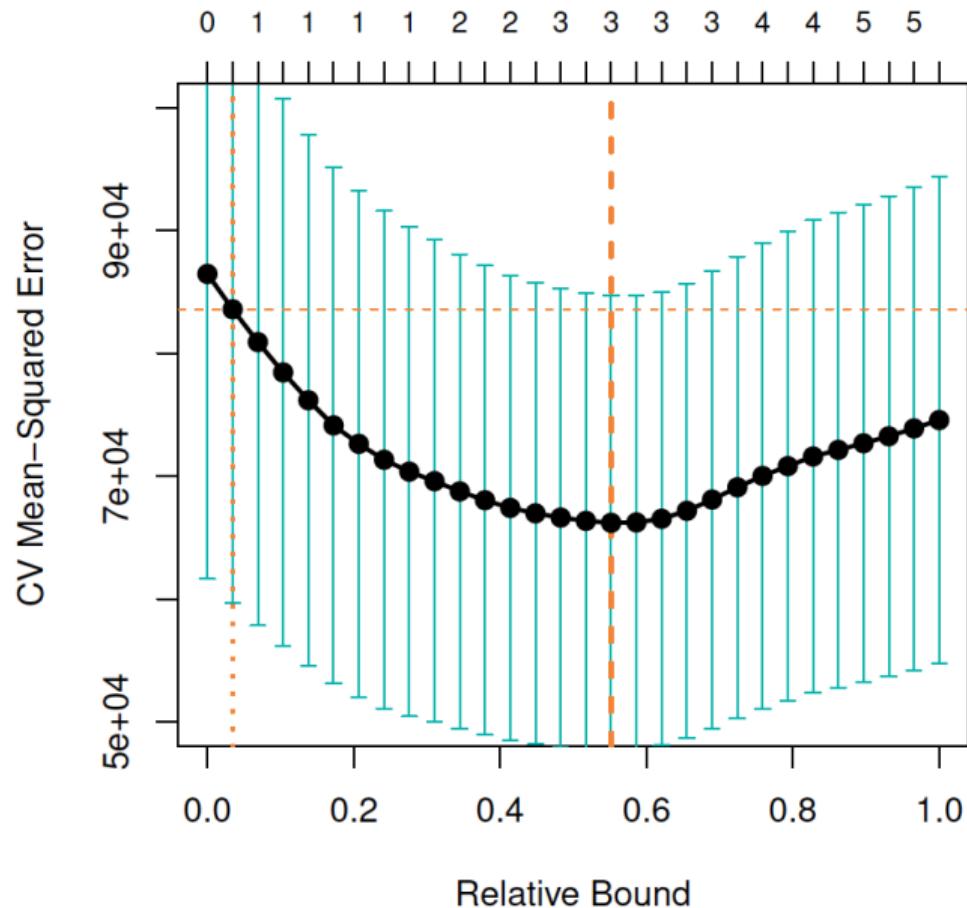


$$\|\hat{\beta}\|_1 / \|\tilde{\beta}\|_1$$

$$\|\hat{\beta}\|_2 / \|\tilde{\beta}\|_2$$

Which penalty?

- The results of LASSO are sensitive to λ .
- Usually chosen via cross-validation by eliminating some test data:
 1. Choose λ to minimize average error of held-out data.
 2. Choose largest λ such that error is within one standard deviation of minimum.
- The choice of penalty parameter also affects model selection consistency.
- **Meinshausen and Bühlmann (2006)**: the prediction-optimal choice of λ is guaranteed NOT to recover the true model asymptotically.
- When λ is chosen by cross validation, we tend to overselect variables and include false positives.
- We need stronger penalization to guarantee model selection consistency, but in finite samples the appropriate choice of λ is not clear.
- Bottom line: in applied work, the LASSO is most useful for variable screening rather than model selection.



Statistical inference

- We might also be interested in performing hypothesis testing on coefficient estimates from LASSO.
- One option is to put selected variables into an OLS regression. Highly problematic because no conditioning on model selection.
- For example, the first variable selected by LASSO will be the one with the highest partial correlation with the response.
- Suppose number of regressors is large and covariates are randomly generated. One of the covariates is likely to be highly correlated with the response; will be selected by LASSO; and have a significant *p*-value in an OLS regression.
- This *post-selection inference* problem is not fully resolved in the statistics literature.
- Possibility: sample splitting. Estimate LASSO on some data points, perform OLS regression on selected variables in the remaining data points.

Adaptive LASSO

- Proposed by Zou (2006).
- The LASSO penalty screens out irrelevant variables, but it also biases the estimates of relevant variables.
- Ideally, we want to weaken the penalty on relevant covariates and vice versa.
- Estimator:

$$\min \sum_{i=1}^N (y_i - \mathbf{x}'_i \beta)^2 + \sum_j \omega_j |\beta_j|$$

where

$$\omega_j = \frac{1}{|\hat{\beta}_j^{OLS}|^\gamma}$$

- You can still use standard software to estimate it.

Elastic net

- Advantages and disadvantages of LASSO vs. ridge: saturation in the “large number of regressors, small sample” case.
- We can combine Lasso and ridge:

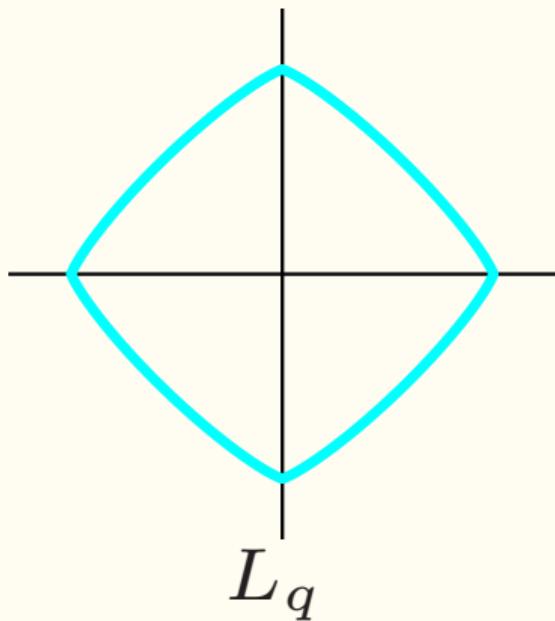
$$\min \sum_{i=1}^N (y_i - \mathbf{x}'_i \beta)^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2$$

or

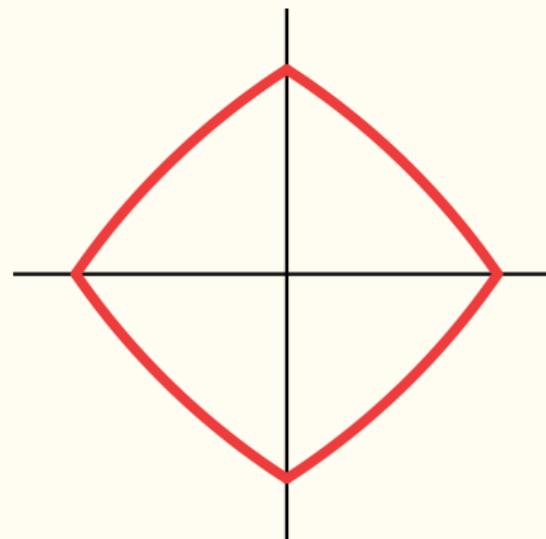
$$\min \sum_{i=1}^N (y_i - \mathbf{x}'_i \beta)^2 + \lambda (\|\beta\|_1 + \alpha \|\beta\|_2)$$

- A elastic net can be expressed as linear support vector machine.
- Particularly useful because we can use available software for support vector machines.

$$q = 1.2$$



$$\alpha = 0.2$$

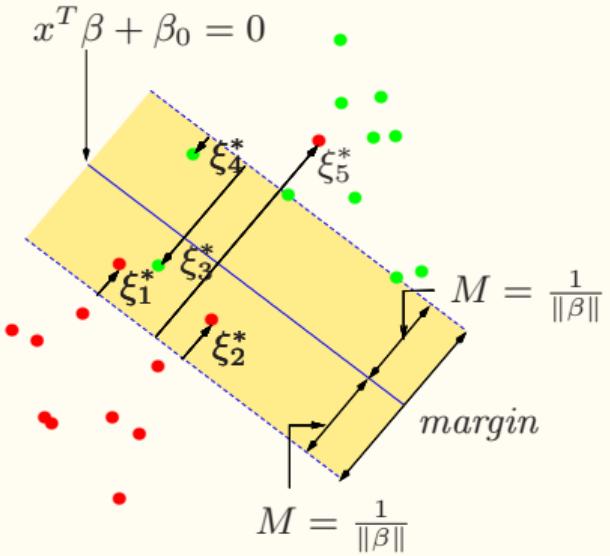
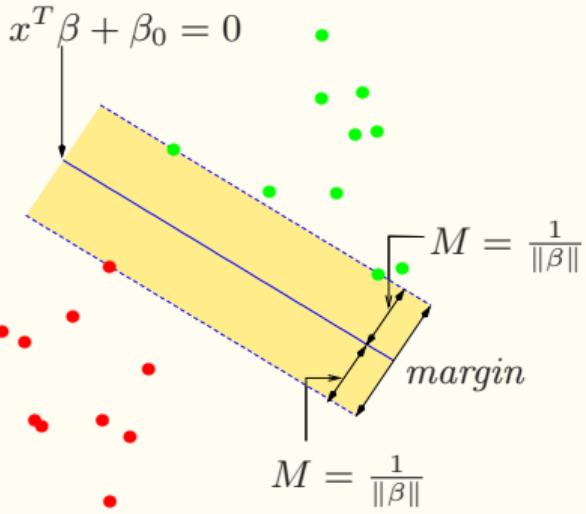


Elastic Net

- References:
 1. *Statistical Learning with Sparsity: The Lasso and Generalizations*, by Trevor Hastie, Robert Tibshirani, and Martin Wainwright (2015).
 2. "High-Dimensional Methods and Inference on Treatment and Structural Effects in Economics," by Belloni, Chernozhukov, and Hansen (2014).
 3. "Inference on Treatment Effects after Selection among High-Dimensional Controls," by Belloni, Chernozhukov, and Hansen (2014).

Example II: Support vector machines

- Developed by Vapnik and Alexey Ya. Chervonenkis (1963).
- One of the most popular algorithms in supervised learning.
- Nice example of the “kernel trick.”
- Classification problem for $\{y_i, \mathbf{x}_i\}_{i=1}^N$ with $y_i \in \{-1, 1\}$.
- Divide observations in two groups by a boundary $f(\mathbf{x}_i)$.
- Maximum-margin hyperplane principle: we search for the boundary that separates as much as possible both regions of data (hard-margin in separable case vs. soft-margin in non-separable case).
- We call observations at the margin support vectors.
- Comparison with logit or probit.



Linear boundaries

- A standard choice: linear boundary with one regressor

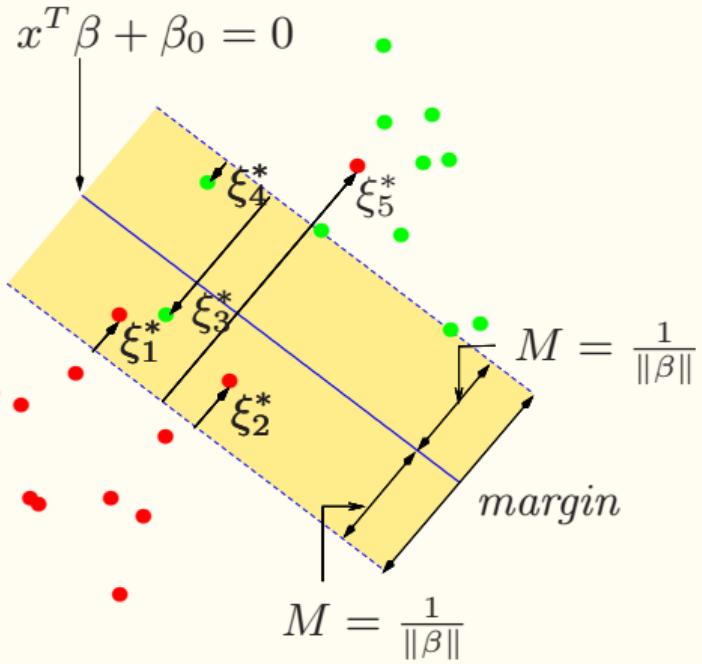
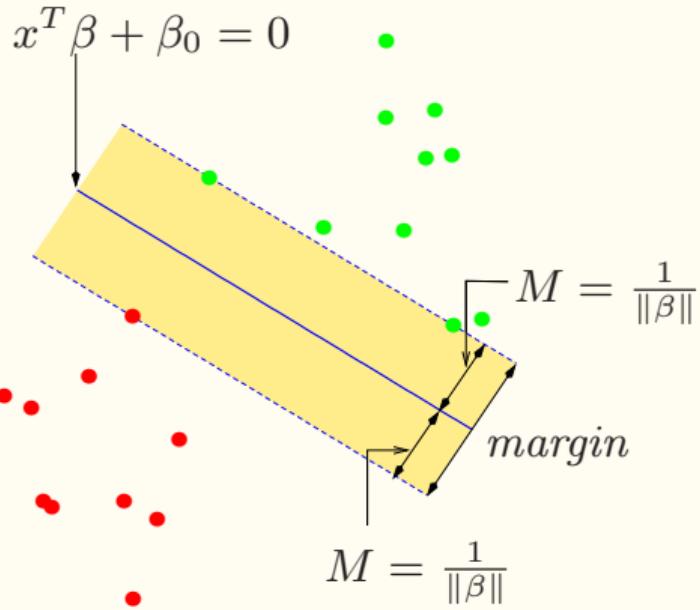
$$y_i = 1 \text{ if } \beta_0 + \beta_1 x \geq 0$$

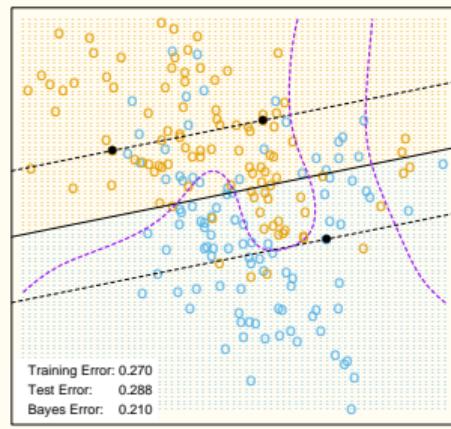
$$y_i = -1 \text{ otherwise}$$

- Find:

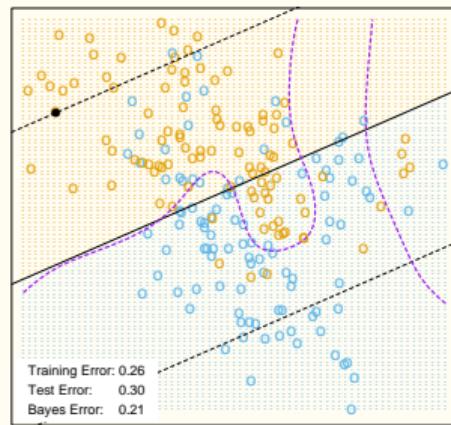
$$\min_{\beta_0, \beta_1} \frac{1}{N} \sum_{i=1}^N \max [0, (1 - y_i (\beta_0 - \beta_1) x)] + \lambda \|\beta\|_2^2$$

- The regularization term $\lambda \|\beta\|_2^2$ is known as the hinge loss.
- Note that by setting $\lambda \rightarrow 0$ in the hinge loss, we approximate a hard-margin boundary.
- You can solve the primal problem with a sub-gradient descent method or dual problem with a coordinate descent algorithm.

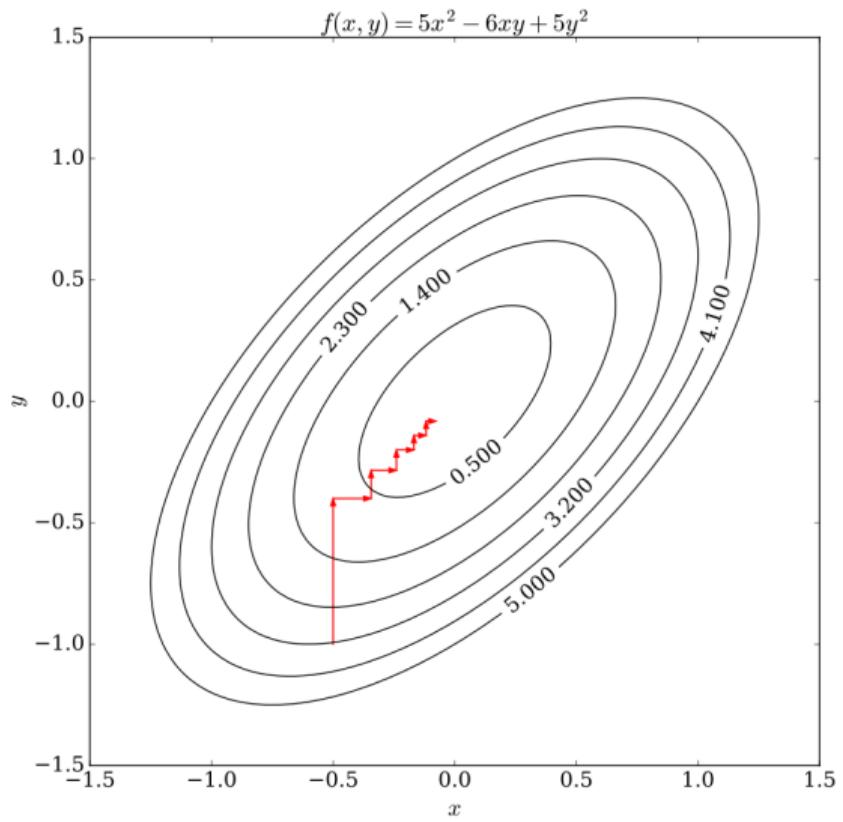




$C = 10000$



$C = 0.01$



Example III: Regression trees

- Classification and Regression Trees, Leo Breiman, Jerome Friedman, Charles J. Stone, and R.A. Olshen.
- Main idea: recursively subdivide regressors into subspaces and compute the mean of the regressor in that subspace.
- Result is a flexible step function.
- Works surprisingly well in practice (very popular in data mining), but nearly no theoretical result (asymptotic normality?).
- Cross-validation to determine penalty λ on the number of leaves in the tree.
- “Bet on sparsity” principle.

Recursive partitions I

- Compute

$$f(\mathbf{x}) = \bar{y}$$

with

$$MSE(f(\mathbf{x})) = \sum_{i=1}^N (y_i - f(\mathbf{x}))^2$$

- For a regressor x_k and a threshold t , find:

$$\bar{y}|x_k \leq t \text{ and } \bar{y}|x_k > t$$

and define

$$f_{k,t}(\mathbf{x}) = \begin{cases} \bar{y}|x_k \leq t \text{ if } x_k \leq t \\ \bar{y}|x_k > t \text{ if } x_k > t \end{cases}$$

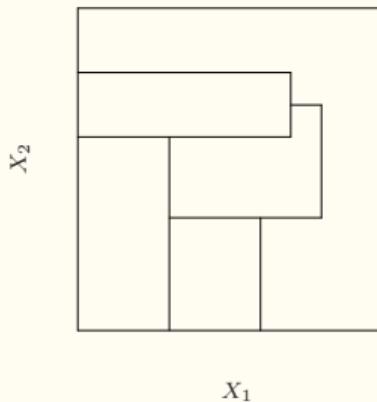
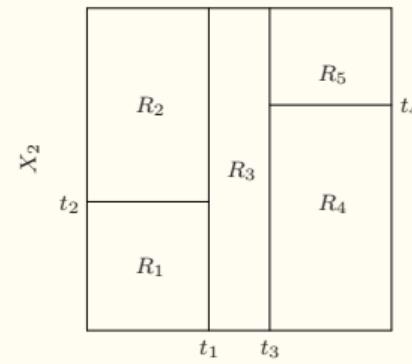
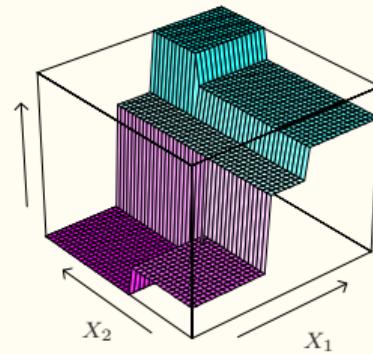
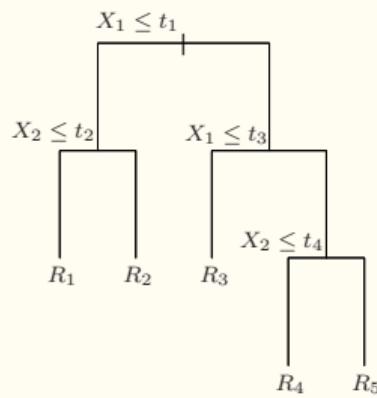
Recursive partitions II

- Select regressor x_k and a threshold t that minimize:

$$(x_k^*, t^*) = \arg \min_{x_k, t} MSE(f_{k,t}(\mathbf{x}))$$

- Keep iterating until

$$\min \{MSE(f_{k,t}(\mathbf{x})) + \lambda \#leaves\}$$


 X_1

 X_1


Improvements

- Pruning an existing tree.
- Boosting: regression tree on new data $(\varepsilon - f(\mathbf{x}), \mathbf{x})$ and iterate.
- Bagging (Bootstrap AGGREGatING): regression trees on bootstraped data, estimate is average over bootstraped samples.
- Random forest.

Other algorithms

- Multi-task LASSO.
- Least Angle Regression (LARS).
- LARS LASSO.
- Gaussian Process Regression (GPR).
- Dantzig Selector.
- Ensemble methods.
- Super learners/stacking (asymptotic properties).