

Predictive Modeling in Insurance, in the context of (possibly) Big Data

A. Charpentier (UQAM & Université de Rennes 1)

Statistical & Actuarial Sciences Joint Seminar

& Center of studies in Asset Management (CESAM)

<http://freakonometrics.hypotheses.org>

Predictive Modeling in Insurance, in the context of (possibly) Big Data

A. Charpentier (UQAM & Université de Rennes 1)

Professor of Actuarial Sciences, Mathematics Department, UQAM
(previously Economics Department, Univ. Rennes 1 & ENSAE Paristech
actuary in Hong Kong, IT & Stats FFSA)

ACTINFO-Covéa Chair, Actuarial Value of Information

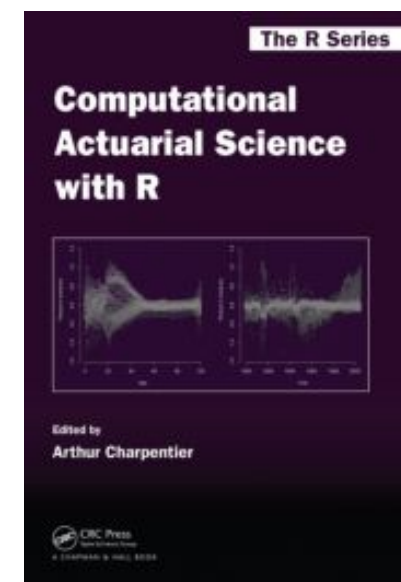
Data Science for Actuaries program, Institute of Actuaries

PhD in Statistics (KU Leuven), Fellow Institute of Actuaries

MSc in Financial Mathematics (Paris Dauphine) & ENSAE

Editor of the freakonometrics.hypotheses.org's blog

Editor of Computational Actuarial Science, CRC



Actuarial Science, an ‘American Perspective’

FUNDAMENTAL CONCEPTS OF ACTUARIAL SCIENCE

CHARLES L. TROWBRIDGE,
F.S.A., M.A.A.A., E.A.

Revised Edition



ACTUARIAL EDUCATION
AND RESEARCH FUND

CONTENTS

Preface for the Actuarial Education & Research Fund	vii	iv Fundamental Concepts of Actuarial Science	
Author's Preface	ix	“Time until Termination” Random Variables	17
I Introduction	1	“Number of Claims” Random Variables	18
Purpose	1	“Claim Amount” Random Variables	19
Audience	2	“Total Claims” Random Variables	19
Geographical Range	2	The Rate of Interest as a Random Variable	20
Brief History of the Actuarial Profession	3	The Importance of Expected Values	21
Evolution	4	Actuarial Interest in Human Mortality	21
Following Chapters	5	The Concept of Credibility	22
II Economics of Risk	7	Summary	23
Introduction	7	References	24
Avoidance or Mitigation of Economic Risk	8	IV The Time Value of Money	27
Financial Security Systems	9	Introduction	27
Classification of Financial Security Systems	9	Time Preference	28
Financial Security Systems as	10	Productivity of Capital	29
Transfer Mechanisms	10	The Uncertain Future	30
The Philosophic Base—Utilitarianism	11	The Level of Interest Rates	31
Utility Theory and Risk Aversion	11	The Actuary's Relationship to the	
The Actuarial Role	12	Time Value of Money	31
Summary	12	Summary	33
References	13	References	34
III Random Variables	15	V Individual Model	35
Introduction	15	Introduction	35
		A Generalized Individual Model	36
		The Concept of Reserves	40
		More Sophisticated Applications of the	
		Generalized Individual Model	41
		Summary	41
		References	42
		VI Collective Models	43
		Introduction	43
		Employee Benefit Plans	44
		Group Model	44
		Defined Benefit Pension Plan Model	46
		The Social Insurance Model	48

Source: Trowbridge (1989) *Fundamental Concepts of Actuarial Science*.

Actuarial Science, a ‘European Perspective’

Contents

*There are 10^{11} stars in the galaxy. That used to be a huge number. But it's only a hundred billion. It's less than the national deficit! We used to call them astronomical numbers. Now we should call them economical numbers—
Richard Feynman (1918–1988)*

1 Utility theory and insurance	1
1.1 Introduction	1
1.2 The expected utility model	2
1.3 Classes of utility functions	5
1.4 Stop-loss reinsurance	8
1.5 Exercises	14
2 The individual risk model	17
2.1 Introduction	17
2.2 Mixed distributions and risks	18
2.3 Convolution	25
2.4 Transforms	28
2.5 Approximations	30
2.5.1 Normal approximation	30
2.5.2 Translated gamma approximation	32
2.5.3 NP approximation	33
2.6 Application: optimal reinsurance	35
2.7 Exercises	36
3 Collective risk models	41
3.1 Introduction	41
3.2 Compound distributions	42
3.2.1 Convolution formula for a compound cdf	44
3.3 Distributions for the number of claims	45
3.4 Properties of compound Poisson distributions	47
3.5 Panjer's recursion	49
3.6 Compound distributions and the Fast Fourier Transform	54
3.7 Approximations for compound distributions	57
3.8 Individual and collective risk model	59
3.9 Loss distributions: properties, estimation, sampling	61
3.9.1 Techniques to generate pseudo-random samples	62
3.9.2 Techniques to compute ML-estimates	63

xv

xvi

Contents

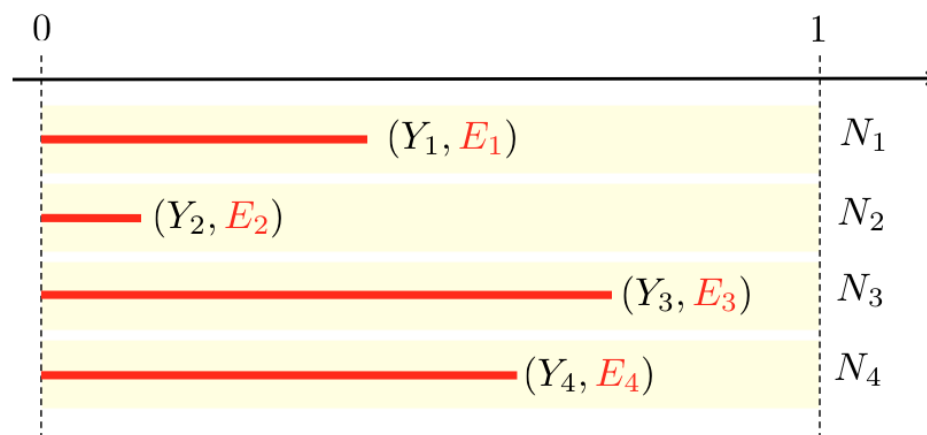
3.9.3 Poisson claim number distribution	63
3.9.4 Negative binomial claim number distribution	64
3.9.5 Gamma claim severity distributions	66
3.9.6 Inverse Gaussian claim severity distributions	67
3.9.7 Mixtures/combinations of exponential distributions	69
3.9.8 Lognormal claim severities	71
3.9.9 Pareto claim severities	72
3.10 Stop-loss insurance and approximations	73
3.10.1 Comparing stop-loss premiums in case of unequal variances	76
3.11 Exercises	78
4 Ruin theory	87
4.1 Introduction	87
4.2 The classical ruin process	89
4.3 Some simple results on ruin probabilities	91
4.4 Ruin probability and capital at ruin	95
4.5 Discrete time model	98
4.6 Reinsurance and ruin probabilities	99
4.7 Beekman's convolution formula	101
4.8 Explicit expressions for ruin probabilities	106
4.9 Approximation of ruin probabilities	108
4.10 Exercises	111
5 Premium principles and Risk measures	115
5.1 Introduction	115
5.2 Premium calculation from top-down	116
5.3 Various premium principles and their properties	119
5.3.1 Properties of premium principles	120
5.4 Characterizations of premium principles	122
5.5 Premium reduction by coinsurance	125
5.6 Value-at-Risk and related risk measures	126
5.7 Exercises	133
6 Bonus-malus systems	135
6.1 Introduction	135
6.2 A generic bonus-malus system	136
6.3 Markov analysis	138
6.3.1 Loimaranta efficiency	141
6.4 Finding steady state premiums and Loimaranta efficiency	142
6.5 Exercises	146
7 Ordering of risks	149
7.1 Introduction	149
7.2 Larger risks	152
7.3 More dangerous risks	154
7.3.1 Thicker-tailed risks	154
7.3.2 Stop-loss order	159
7.3.3 Exponential order	160

Source: Dhaene et al. (2004) *Modern Actuarial Risk Theory*.

Exemples of Actuarial Problems: Ratemaking and Pricing

$$\mathbb{E}[S|\mathbf{X}] = \mathbb{E}\left[\sum_{i=1}^N Z_i \mid \mathbf{X}\right] = \underbrace{\mathbb{E}[N|\mathbf{X}]}_{\text{annual frequency}} \cdot \underbrace{\mathbb{E}[Z_i|\mathbf{X}]}_{\text{individual cost}}$$

- **censoring / incomplete datasets** (exposure + delay to report claims)



We observe Y and E , but the variable of interest is N .

$$Y_i \sim \mathcal{P}(E_i \cdot \lambda_i) \text{ with } \lambda_i = \exp[\beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta} + Z_i].$$

Exemples of Actuarial Problems: Pricing and Classification

Econometric models on classes

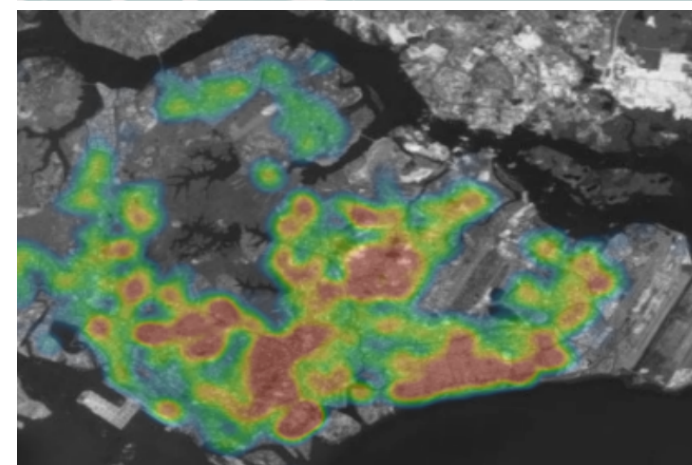
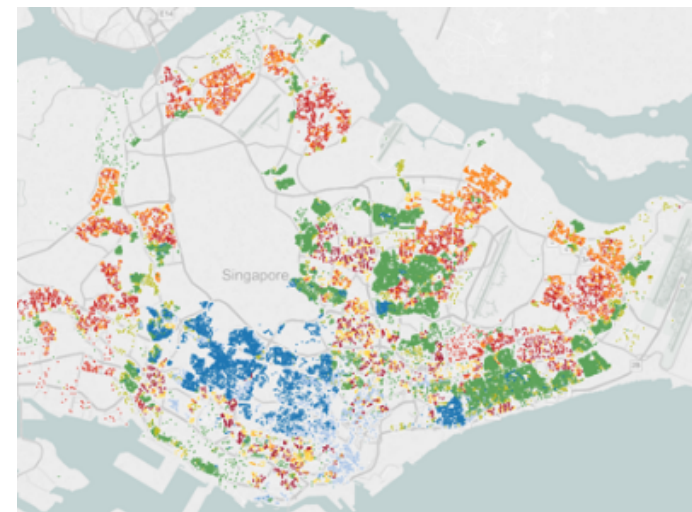
$$Y_i \sim \mathcal{B}(p_i) \text{ with } p_i = \frac{\exp[\beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta}]}{1 + \exp[\beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta}]}$$

or on counts

$$Y_i \sim \mathcal{P}(\lambda_i) \text{ with } \lambda_i = \exp[\beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta}]$$

- **(too) large datasets** \mathbf{X} can be large (or complex)

factors with a lot of modalities, spatial data, text information, etc.



Exemples of Actuarial Problems: Pricing and Classification

How to avoid overfit? How to group modalities?

How to choose between (very) correlated features?

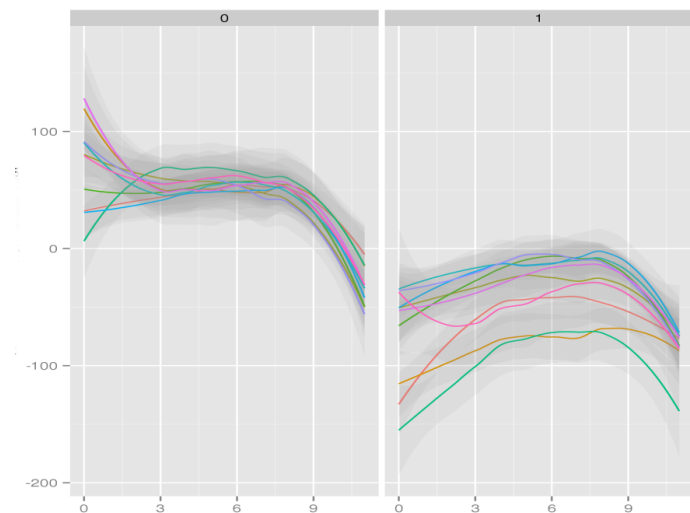
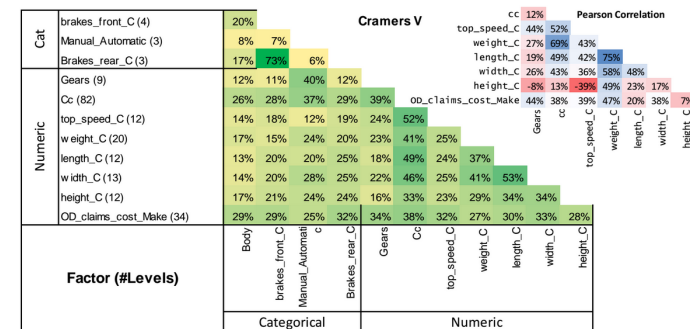
- **model selection** issues

Historically **Bailey (1963)** ‘margin method’ $\hat{n} = \mathbf{r}\mathbf{c}^T$, with row (\mathbf{r}) and column (\mathbf{c}) effects, and constraints

$$\sum_i n_{i,j} = \sum_i r_i \cdot c_j \text{ and } \sum_j n_{i,j} = \sum_j r_i \cdot c_j$$

Related to Poisson regression,

$$N \sim \mathcal{P}(\exp[\beta_0 + \mathbf{r}^T \boldsymbol{\beta}_R + \mathbf{c}^T \boldsymbol{\beta}_C])$$



Exemples of Actuarial Problems: Claims Reserving and Predictive Models

- **predictive modeling** issues

In all those cases, the goal is to get a predictive model, $\hat{y} = \hat{m}(\mathbf{x})$ given some features \mathbf{x} .

Recall that the main interest in insurance is either

- a probability $m(\mathbf{x}) = \mathbb{P}[Y = 1 | \mathbf{X} = \mathbf{x}]$
- an expected value $m(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$

but sometimes, we need the (conditional) distribution of \hat{Y} .

Observed Increments, in Currency Units									
AY/DY	0	1	2	3	4	5	6	7	8
1999	1217	1041	653	3	186	6	150	174	237
2000	1233	1374	134	235	24	167	166	280	249
2001	1264	1229	322	147	167	141	152	236	253
2002	1301	1538	270	214	158	185	180	265	284
2003	1527	1732	296	198	264	158	205	302	323
2004	1655	1683	261	160	189	155	201	297	318
2005	1753	1414	246	174	180	180	148	192	283
2006	1765	1237	328	170	176	144	187	276	296
2007	1900	1825	407	211	218	179	233	343	367

↓

Observed Cumulative Data, in Currency Units									
AY/DY	0	1	2	3	4	5	6	7	8
1999	1217	2258	2911	2915	3101	3107	3257	3431	3668
2000	1233	2607	2741	2977	3000	3167	3333	3613	3863
2001	1264	2493	2815	2962	3129	3270	3422	3658	3910
2002	1301	2839	3109	3322	3481	3666	3846	4111	4395
2003	1527	3259	3555	3753	4017	4175	4380	4681	5005
2004	1655	3338	3599	3759	3948	4103	4304	4601	4918
2005	1753	3166	3413	3587	3767	3915	4107	4390	4693
2006	1765	3001	3329	3499	3675	3819	4007	4283	4578
2007	1900	3725	4133	4344	4562	4741	4973	5316	5683

↑

Development Factors									
	0	1	2	3	4	5	6	7	8
	1.9603	1.1093	1.0511	1.0502	1.0393	1.0491	1.0689	1.0690	1

Source: Swedish data—Data from: Swedish Financial Supervisory Authority.
 Note: The figures in gray are estimated using an ordinary development approach. AY: Accident year; DY: Development year. Currency in 1000 units.

History of Actuarial Models (in one slide)

Bailey (1963) or Taylor (1977) considered **deterministic models**, $n_{i,j} = r_i \cdot c_j$ or $n_{i,j} = r_i \cdot d_{i+j}$. Some additional constraints are given to get an identifiable model.

Then some **stochastic** version of those models were introduced, see [Hachemeister \(1975\)](#) or [de Vylder \(1985\)](#), e.g.

$$N_{i,j} \sim \mathcal{P}(\exp[\mathcal{P}(\exp[\beta_0 + \mathbf{R}^\top \boldsymbol{\beta}_R + \mathbf{C}^\top \boldsymbol{\beta}_C])])$$

or

$$\log N_{i,j} \sim \mathcal{N}(\beta_0 + \mathbf{R}^\top \boldsymbol{\beta}_R + \mathbf{C}^\top \boldsymbol{\beta}_C, \sigma^2)$$

All those techniques are econometric-based techniques. Why not consider some **statistical learning** techniques?

Statistical Learning and Philosophical Issues

From *Machine Learning and Econometrics*, by Hal Varian :

“**Machine learning** use data to predict some variable as a function of other covariables,

- may, or may not, care about insight, importance, patterns
- may, or may not, care about inference (how y changes as some x change)

Econometrics use statistical methodes for prediction, inference and causal modeling of economic relationships

- hope for some sort of insight (inference is a goal)
- in particular, causal inference is goal for decision making.”

→ machine learning, ‘new tricks for econometrics’

Statistical Learning and Philosophical Issues

Remark machine learning can also learn from econometrics, especially with non i.i.d. data (time series and panel data)

Remark machine learning can help to get better predictive models, given good datasets. No use on several data science issues (e.g. selection bias).

Machine Learning and ‘Statistics’

Machine learning and statistics seem to be very similar, they share the same goals—they both focus on data modeling—but their methods are affected by their cultural differences.

“The goal for a statistician is to predict an interaction between variables with some degree of certainty (we are never 100% certain about anything). Machine learners, on the other hand, want to build algorithms that predict, classify, and cluster with the most accuracy, see [Why a Mathematician, Statistician & Machine Learner Solve the Same Problem Differently](#)

Machine learning methods are about algorithms, more than about asymptotic statistical properties.

Machine Learning and 'Statistics'

See also nonparametric inference: “Note that the non-parametric model is not none-parametric: parameters are determined by the training data, not the model. [...] non-parametric covers techniques that do not assume that the structure of a model is fixed. Typically, the model grows in size to accommodate the complexity of the data.” see [wikipedia](#)

Validation is not based on mathematical properties, but on properties out of sample: we must use a [training sample](#) to train (estimate) model, and a [testing sample](#) to compare algorithms.

Goldilock Principle: the Mean-Variance Tradeoff

In statistics and in machine learning, there will be **parameters** and **meta-parameters** (or **tunning parameters**). The first ones are estimated, the second ones should be chosen.

See **Hill estimator** in extreme value theory. X has a Pareto distribution above some threshold u if

$$\mathbb{P}[X > x | X > u] = \left(\frac{u}{x}\right)^{\frac{1}{\xi}} \text{ for } x > u.$$

Given a sample \mathbf{x} , consider the Pareto-QQ plot, i.e. the scatterplot

$$\left\{ -\log \left(1 - \frac{i}{n+1} \right), \log x_{i:n} \right\}_{i=n-k, \dots, n}$$

for points exceeding $X_{n-k:n}$.

Goldilock Principle: the Mean-Variance Tradeoff

The slope is ξ , i.e.

$$\log X_{n-i+1:n} \approx \log X_{n-k:n} + \xi \left(-\log \frac{i}{n+1} - \log \frac{n+1}{k+1} \right)$$

Hence, consider estimator $\hat{\xi}_k = \frac{1}{k} \sum_{i=0}^{k-1} \log x_{n-i:n} - \log x_{n-k:n}$.

Standard **mean-variance tradeoff**,

- k large: bias too large, variance too small
- k small: variance too large, bias too small

Goldilock Principle: the Mean-Variance Tradeoff

Same holds in **kernel regression**, with bandwidth h (length of neighborhood)

$$\hat{m}_h(x) = \frac{\sum_{i=1}^n K_h(x - x_i) \cdot y_i}{\sum_{i=1}^n K_h(x - x_i)}$$

for some kernel $K(\cdot)$.

Standard **mean-variance tradeoff**,

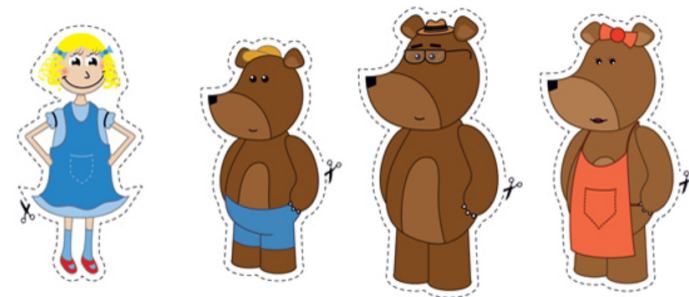
- h large: bias too large, variance too small
- h small: variance too large, bias too small

Goldilock Principle: the Mean-Variance Tradeoff

More generally, we estimate $\hat{\theta}_h$ or $\hat{m}_h(\cdot)$

Use the **mean squared error** for $\hat{\theta}_h$

$$\mathbb{E} \left[\left(\theta - \hat{\theta}_h \right)^2 \right]$$



or **mean integrated squared error** $\hat{m}_h(\cdot)$,

$$\mathbb{E} \left[\int (m(\mathbf{x}) - \hat{m}_h(\mathbf{x}))^2 d\mathbf{x} \right]$$

In statistics, derive an asymptotic expression for these quantities, and find h^* that minimizes those.

Goldilock Principle: the Mean-Variance Tradeoff

In **classical statistics**, the MISE can be approximated by

$$\frac{h^4}{4} \left(\int x^2 K(\mathbf{x}) d\mathbf{x} \right)^2 \int \left(m''(\mathbf{x}) + 2m'(\mathbf{x}) \frac{f'(\mathbf{x})}{f(\mathbf{x})} \right) d\mathbf{x} + \frac{1}{nh} \sigma^2 \int K^2(x) dx \int \frac{d\mathbf{x}}{f(\mathbf{x})}$$

where f is the density of \mathbf{x} 's. Thus the optimal h is

$$h^* = n^{-\frac{1}{5}} \left(\frac{\sigma^2 \int K^2(x) dx \int \frac{d\mathbf{x}}{f(\mathbf{x})}}{\left(\int x^2 K(\mathbf{x}) d\mathbf{x} \right)^2 \int \left(\int m''(\mathbf{x}) + 2m'(\mathbf{x}) \frac{f'(\mathbf{x})}{f(\mathbf{x})} \right)^2 d\mathbf{x}} \right)^{\frac{1}{5}}$$

(hard to get a simple rule of thumb... up to a constant, $h^* \sim n^{-\frac{1}{5}}$)

In **statistics learning**, use bootstrap, or cross-validation to get an optimal h ...

Randomization is too important to be left to chance!

Consider some sample $\mathbf{x} = (x_1, \dots, x_n)$ and some statistics $\hat{\theta}$. Set $\hat{\theta}_n = \hat{\theta}(\mathbf{x})$

Jackknife used to reduce bias: set $\hat{\theta}_{(-i)} = \hat{\theta}(\mathbf{x}_{(-i)})$, and $\tilde{\theta} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{(-i)}$

If $\mathbb{E}(\hat{\theta}_n) = \theta + O(n^{-1})$ then $\mathbb{E}(\tilde{\theta}_n) = \theta + O(n^{-2})$.

See also **leave-one-out** cross validation, for $\hat{m}(\cdot)$

$$\text{mse} = \frac{1}{n} \sum_{i=1}^n [y_i - \hat{m}_{(-i)}(x_i)]^2$$

Bootstrap estimate is based on bootstrap samples: set $\hat{\theta}_{(b)} = \hat{\theta}(\mathbf{x}_{(b)})$, and

$\tilde{\theta} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{(b)}$, where $\mathbf{x}_{(b)}$ is a vector of size n , where values are drawn from $\{x_1, \dots, x_n\}$, with replacement. And then use the law of large numbers...

See [Efron \(1979\)](#).

Statistical Learning and Philosophical Issues

From (y_i, \mathbf{x}_i) , there are different stories behind, see [Freedman \(2005\)](#)

- the **causal story** : $x_{j,i}$ is usually considered as independent of the other covariates $x_{k,i}$. For all possible \mathbf{x} , that value is mapped to $m(\mathbf{x})$ and a noise is attached, ε . The goal is to recover $m(\cdot)$, and the residuals are just the difference between the response value and $m(\mathbf{x})$.
- the **conditional distribution story** : for a linear model, we usually say that Y given $\mathbf{X} = \mathbf{x}$ is a $\mathcal{N}(m(\mathbf{x}), \sigma^2)$ distribution. $m(\mathbf{x})$ is then the conditional mean. Here $m(\cdot)$ is assumed to really exist, but no causal assumption is made, only a conditional one.
- the **explanatory data story** : there is no model, just data. We simply want to summarize information contained in \mathbf{x} 's to get an accurate summary, close to the response (i.e. $\min\{\ell(\mathbf{y}_i, m(\mathbf{x}_i))\}$) for some loss function ℓ .

Machine Learning vs. Statistical Modeling

In **machine learning**, given some dataset (\mathbf{x}_i, y_i) , solve

$$\hat{m}(\cdot) = \operatorname{argmin}_{m(\cdot) \in \mathcal{F}} \left\{ \sum_{i=1}^n \ell(y_i, m(\mathbf{x}_i)) \right\}$$

for some **loss functions** $\ell(\cdot, \cdot)$.

In **statistical modeling**, given some probability space $(\Omega, \mathcal{A}, \mathbb{P})$, assume that y_i are realization of i.i.d. variables Y_i (given $\mathbf{X}_i = \mathbf{x}_i$) with distribution F_i . Then solve

$$\hat{m}(\cdot) = \operatorname{argmax}_{m(\cdot) \in \mathcal{F}} \{ \log \mathcal{L}(m(\mathbf{x}); \mathbf{y}) \} = \operatorname{argmax}_{m(\cdot) \in \mathcal{F}} \left\{ \sum_{i=1}^n \log f(y_i; m(\mathbf{x}_i)) \right\}$$

where $\log \mathcal{L}$ denotes the **log-likelihood**.

Loss Functions

Fitting criteria are based on **loss functions** (also called **cost functions**). For a quantitative response, a popular one is the quadratic loss,

$$\ell(y, m(\mathbf{x})) = [y - m(\mathbf{x})]^2.$$

Recall that

$$\left\{ \begin{array}{l} \mathbb{E}(Y) = \operatorname{argmin}_{m \in \mathbb{R}} \{ \|Y - m\|_{\ell_2} \} = \operatorname{argmin}_{m \in \mathbb{R}} \{ \mathbb{E}([Y - m]^2) \} \\ \operatorname{Var}(Y) = \min_{m \in \mathbb{R}} \{ \mathbb{E}([Y - m]^2) \} = \mathbb{E}([Y - \mathbb{E}(Y)]^2) \end{array} \right.$$

The empirical version is

$$\left\{ \begin{array}{l} \bar{y} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n \frac{1}{n} [y_i - m]^2 \right\} \\ s^2 = \min_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n \frac{1}{n} [y_i - m]^2 \right\} = \sum_{i=1}^n \frac{1}{n} [y_i - \bar{y}]^2 \end{array} \right.$$

Model Evaluation

In linear models, the R^2 is defined as the proportion of the variance of the the response y that can be obtained using the predictors.

But maximizing the R^2 usually yields **overfit** (or **unjustified optimism** in Berk (2008)).

In linear models, consider the adjusted R^2 ,

$$\bar{R}^2 = 1 - [1 - R^2] \frac{n - 1}{n - p - 1}$$

where p is the number of parameters, or more generally $\text{trace}(\mathbf{S})$ when some smoothing matrix is considered

$$\hat{y} = \hat{m}(x) = \sum_{i=1}^n \mathbf{S}_{x,i} y_i = \mathbf{S}_x^\top \mathbf{y}$$

where \mathbf{S}_x is some vector of weights (called **smoother vector**), related to a $n \times n$ smoother matrix, $\hat{y} = \mathbf{S} \mathbf{y}$ where prediction is done at points x_i 's.

Model Evaluation

Alternatives are based on the Akaike Information Criterion (**AIC**) and the Bayesian Information Criterion (**BIC**), based on a penalty imposed on some criteria (the logarithm of the variance of the residuals),

$$AIC = \log \left(\frac{1}{n} \sum_{i=1}^n [y_i - \hat{y}_i]^2 \right) + \frac{2p}{n}$$

$$BIC = \log \left(\frac{1}{n} \sum_{i=1}^n [y_i - \hat{y}_i]^2 \right) + \frac{\log(n)p}{n}$$

In a more general context, replace p by $\text{trace}(\mathbf{S})$.

Model Evaluation

One can also consider the expected prediction error (with a probabilistic model)

$$\mathbb{E}[\ell(Y, \hat{m}(\mathbf{X}))]$$

We cannot claim (using the law of large number) that

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{m}(\mathbf{x}_i)) \stackrel{a.s.}{\not\rightarrow} \mathbb{E}[\ell(Y, m(\mathbf{X}))]$$

since \hat{m} depends on (y_i, \mathbf{x}_i) 's.

Natural option : use two (random) samples, a **training** one and a **validation** one.

Alternative options, use cross-validation, leave-one-out or k -fold.

Underfit / Overfit and Variance - Mean Tradeoff

Goal in predictive modeling: reduce uncertainty in our predictions.

Need more data to get a better knowledge.

Unfortunately, reducing the error of the prediction on a dataset does not generally give a good **generalization** performance

→ need a training and a validation dataset

Overfit, Training vs. Validation and Complexity (Vapnik Dimension)

complexity \longleftrightarrow polynomial degree

Overfit, Training vs. Validation and Complexity (Vapnik Dimension)

complexity \longleftrightarrow number of neighbors (k)

Logistic Regression

Assume that $\mathbb{P}(Y_i = 1) = \pi_i$,

$$\text{logit}(\pi_i) = \mathbf{x}_i^\top \boldsymbol{\beta}, \text{ where } \text{logit}(\pi_i) = \log\left(\frac{\pi_i}{1 - \pi_i}\right),$$

or

$$\pi_i = \text{logit}^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta}) = \frac{\exp[\mathbf{x}_i^\top \boldsymbol{\beta}]}{1 + \exp[\mathbf{x}_i^\top \boldsymbol{\beta}]}.$$

The log-likelihood is

$$\log \mathcal{L}(\boldsymbol{\beta}) = \sum_{i=1}^n y_i \log(\pi_i) + (1 - y_i) \log(1 - \pi_i) = \sum_{i=1}^n y_i \log(\pi_i(\boldsymbol{\beta})) + (1 - y_i) \log(1 - \pi_i(\boldsymbol{\beta}))$$

and the first order conditions are **solved numerically**

$$\frac{\partial \log \mathcal{L}(\boldsymbol{\beta})}{\partial \beta_k} = \sum_{i=1}^n X_{k,i} [y_i - \pi_i(\boldsymbol{\beta})] = 0.$$

Predictive Classifier

To go from a score

$$\hat{s}(\mathbf{x}) = \frac{\exp[\mathbf{x}^T \hat{\boldsymbol{\beta}}]}{1 + \exp[\mathbf{x}^T \hat{\boldsymbol{\beta}}]}$$

to a class:

if $\hat{s}(\mathbf{x}) > s$, then $\hat{Y}(\mathbf{x}) = 1$ (or ●) and $\hat{s}(\mathbf{x}) \leq s$, then $\hat{Y}(\mathbf{x}) = 0$ (or ●).

Plot $TP(s) = \mathbb{P}[\hat{Y} = 1|Y = 1]$ against $FP(s) = \mathbb{P}[\hat{Y} = 1|Y = 0]$

Why a Logistic and not a Probit Regression?

Bliss (1934) suggested a model such that

$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = H(\mathbf{x}^\top \boldsymbol{\beta}) \text{ where } H(\cdot) = \Phi(\cdot)$$

the c.d.f. of the $\mathcal{N}(0, 1)$ distribution. This is the **probit** model.

This yields a latent model, $y_i = \mathbf{1}(y_i^* > 0)$ where

$$y_i^* = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i \text{ is a nonobservable score.}$$

In the logistic regression, we model the **odds ratio**,

$$\frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y \neq 1 | \mathbf{X} = \mathbf{x})} = \exp[\mathbf{x}^\top \boldsymbol{\beta}]$$

$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = H(\mathbf{x}^\top \boldsymbol{\beta}) \text{ where } H(\cdot) = \frac{\exp[\cdot]}{1 + \exp[\cdot]}$$

which is the c.d.f. of the **logistic** variable, see **Verhulst (1845)**

Table 3.2 Transformation of percentages to probits

%	0	1	2	3	4	5	6	7	8	9
0	—	2.67	2.05	3.12	3.25	3.36	3.45	3.52	3.59	3.66
10	3.72	3.77	3.82	3.87	3.92	3.96	4.01	4.05	4.08	4.12
20	4.16	4.10	4.23	4.20	4.20	4.33	4.30	4.39	4.42	4.45
30	4.48	4.50	4.53	4.50	4.59	4.61	4.64	4.67	4.69	4.72
40	4.75	4.77	4.80	4.82	4.85	4.87	4.90	4.92	4.95	4.97
50	5.00	5.03	5.05	5.08	5.10	5.13	5.15	5.18	5.20	5.23
60	5.25	5.28	5.31	5.33	5.30	5.39	5.41	5.44	5.47	5.50
70	5.62	5.65	5.68	5.61	5.64	5.67	5.71	5.74	5.77	5.81
80	5.84	5.88	5.92	5.95	5.99	6.04	6.08	6.13	6.18	6.23
90	6.28	6.34	6.41	6.48	6.55	6.64	6.75	6.88	7.05	7.33
—	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
99	7.33	7.37	7.41	7.46	7.51	7.58	7.65	7.75	7.88	8.00

Soit p la population : représentons par dp l'accroissement infiniment petit qu'elle reçoit pendant un temps infiniment court dt . Si la population croissait en progression géométrique, nous aurions l'équation $\frac{dp}{dt} = mp$. Mais comme la vitesse d'accroissement de la population est retardée par l'augmentation même du nombre des habitants, nous devrons retrancher de mp une fonction inconnue de p ; de manière que la formule à intégrer devienne

$$\frac{dp}{dt} = mp - \varphi(p).$$

L'hypothèse la plus simple que l'on puisse faire sur la forme de la fonction φ , est de supposer $\varphi(p) = np^2$. On trouve alors pour intégrale de l'équation ci-dessus

$$t = \frac{1}{m} [\log. p - \log. (m - np)] + \text{constante},$$

et il suffira de trois observations pour déterminer les deux coefficients constants m et n et la constante arbitraire.

En résolvant la dernière équation par rapport à p , il vient

$$p = \frac{mp' e^{mt}}{np' e^{mt} + m - np'} \dots \dots \dots (1)$$

en désignant par p' la population qui répond à $t = 0$, et par e la base des logarithmes népériens. Si l'on fait $t = \infty$, on voit que la valeur de p correspondante est $P = \frac{m}{n}$. Telle est donc la limite supérieure de la population.

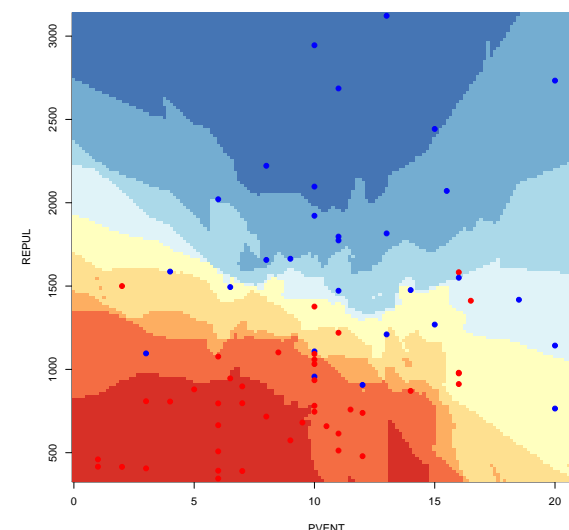
k -Nearest Neighbors (a.k.a. k -NN)

In pattern recognition, the k -Nearest Neighbors algorithm (or k -NN for short) is a non-parametric method used for classification and regression. (Source: [wikipedia](#)).

$$\mathbb{E}[Y|\mathbf{X} = \mathbf{x}] \sim \frac{1}{k} \sum_{d(\mathbf{x}_i, \mathbf{x}) \text{ small}} y_i$$

For k -Nearest Neighbors, the class is usually the **majority** vote of the k closest neighbors of \mathbf{x} .

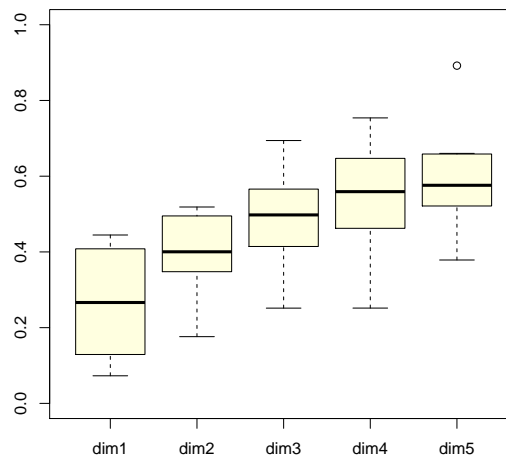
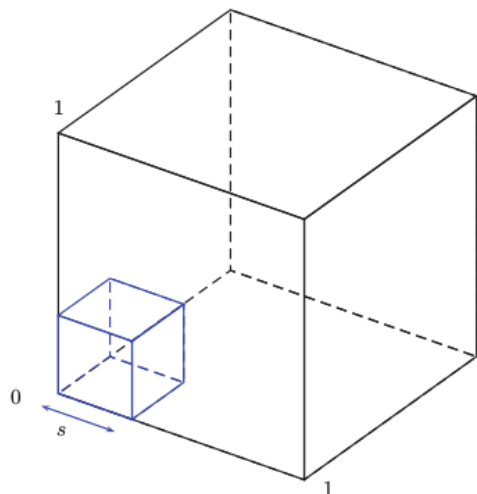
Distance $d(\cdot, \cdot)$ should not be sensitive to units:
normalize by standard deviation



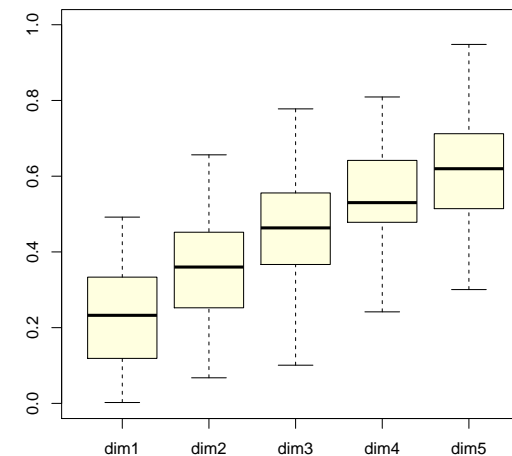
k -Nearest Neighbors and Curse of Dimensionality

The higher the dimension, the larger the distance to the closest neighbor

$$\min_{i \in \{1, \dots, n\}} \{d(\mathbf{a}, \mathbf{x}_i)\}, \mathbf{x}_i \in \mathbb{R}^d.$$



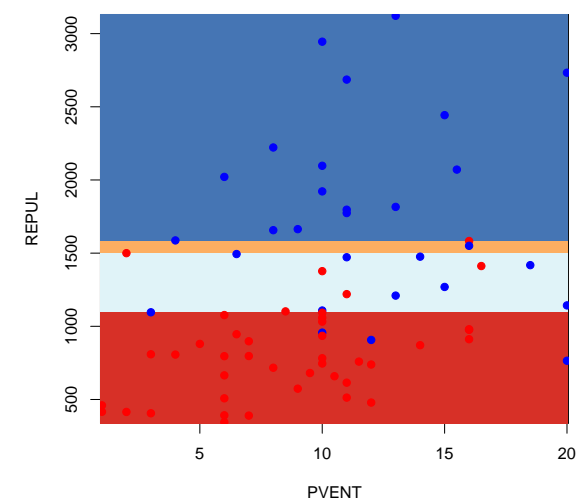
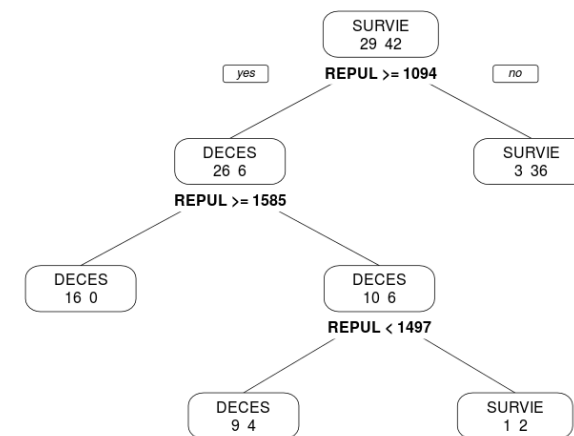
$n = 10$



$n = 100$

Classification (and Regression) Trees, CART

one of the predictive modelling approaches used in statistics, data mining and machine learning [...]
 In tree structures, leaves represent class labels and branches represent conjunctions of features that lead to those class labels. (Source: [wikipedia](https://en.wikipedia.org/wiki/Decision_tree_learning)).



Bagging

Bootstrapped Aggregation (Bagging) , is a machine learning ensemble meta-algorithm designed to improve the stability and accuracy of machine learning algorithms used in statistical classification (Source: [wikipedia](#)).

It is an ensemble method that creates multiple models of the same type from different sub-samples of the same dataset [**bootstrap**]. The predictions from each separate model are combined together to provide a superior result [**aggregation**].

→ can be used on any kind of model, but interesting for trees, see [Breiman \(1996\)](#)

Bootstrap can be used to define the concept of **margin**,

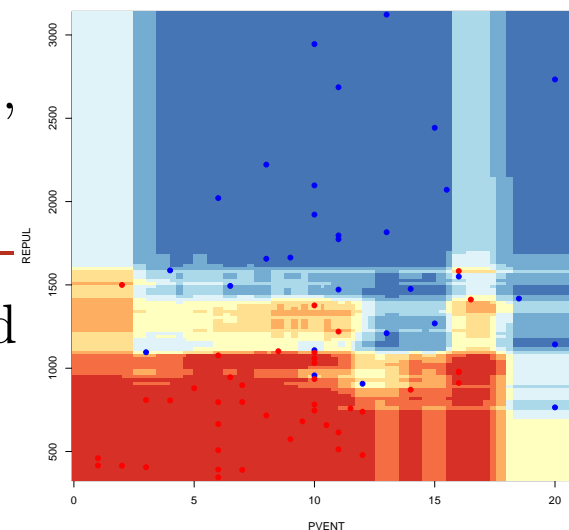
$$\text{margin}_i = \frac{1}{B} \sum_{b=1}^B \mathbf{1}(\hat{y}_i = y_i) - \frac{1}{B} \sum_{b=1}^B \mathbf{1}(\hat{y}_i \neq y_i)$$

Remark Probability that i th row is not selection $(1 - n^{-1})^n \rightarrow e^{-1} \sim 36.8\%$, cf training / validation samples (2/3-1/3)

Random Forests

Strictly speaking, when bootstrapping among observations, and aggregating, we use a bagging algorithm.

In the **random forest** algorithm, we combine Breiman's **bagging** idea and the random selection of features, introduced independently by **Ho (1995)** and **Amit & Geman (1997)**.



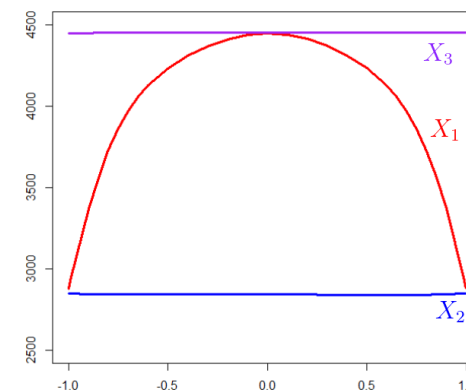
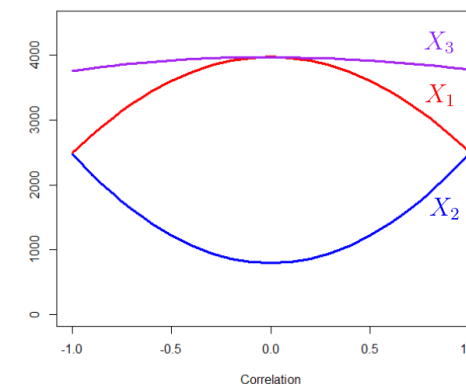
Bagging, Forests, and Variable Importance

Given some random forest with M trees, set

$$VI(X_k) = \frac{1}{M} \sum_m \sum_t \frac{N_t}{N} \Delta i(t)$$

where the first sum is over all trees, and the second one is over all nodes where the split is done based on variable X_k . But difficult to interpret with correlated features. Consider model $y = \beta_0 + \beta_1 x_1 + \beta_3 x_3 + \varepsilon$, and we consider a model based on $\mathbf{x} = (x_1, x_2, x_3)$ where x_1 and x_2 are correlated.

Compare AIC vs. Variable Importance as a function of r



Support Vector Machine

SVMs were developed in the 90's based on previous work, from [Vapnik & Lerner \(1963\)](#), see [Vailant \(1984\)](#)

Assume that points are **linearly separable**, i.e. there is ω and b such that

$$Y = \begin{cases} +1 & \text{if } \omega^T \mathbf{x} + b > 0 \\ -1 & \text{if } \omega^T \mathbf{x} + b < 0 \end{cases}$$

Problem: infinite number of solutions, need a **good** one, that separate the data, (somehow) far from the data.

The distance from \mathbf{x}_0 to the hyperplane $\omega^T \mathbf{x} + b$ is

$$d(\mathbf{x}_0, H_{\omega,b}) = \frac{\omega^T \mathbf{x}_0 + b}{\|\omega\|}$$

Support Vector Machine

Define **support vectors** as observations such that

$$|\boldsymbol{\omega}^\top \mathbf{x}_i + b| = 1$$

The margin is the distance between hyperplanes defined by support vectors.

The distance from support vectors to $H_{\boldsymbol{\omega},b}$ is $\|\boldsymbol{\omega}\|^{-1}$, and the margin is then $2\|\boldsymbol{\omega}\|^{-1}$.

→ the algorithm is to minimize the inverse of the margins s.t. $H_{\boldsymbol{\omega},b}$ separates ± 1 points, i.e.

$$\min \left\{ \frac{1}{2} \boldsymbol{\omega}^\top \boldsymbol{\omega} \right\} \text{ s.t. } y_i(\boldsymbol{\omega}^\top \mathbf{x}_i + b) \geq 1, \forall i.$$

Support Vector Machine

Now, what about the **non-separable case**?

Here, we **cannot** have $y_i(\omega^\top \mathbf{x}_i + b) \geq 1 \forall i$.

→ introduce **slack variables**,

$$\begin{cases} \omega^\top \mathbf{x}_i + b \geq +1 - \xi_i & \text{when } y_i = +1 \\ \omega^\top \mathbf{x}_i + b \leq -1 + \xi_i & \text{when } y_i = -1 \end{cases}$$

where $\xi_i \geq 0 \forall i$. There is a classification error when $\xi_i > 1$.

The idea is then to solve

$$\min \left\{ \frac{1}{2} \omega^\top \omega + C \mathbf{1}^\top \mathbf{1}_{\xi > 1} \right\}, \text{ instead of } \min \left\{ \frac{1}{2} \omega^\top \omega \right\}$$

Support Vector Machines, with a Linear Kernel

So far, $d(\mathbf{x}_0, H_{\omega, b}) = \min_{\mathbf{x} \in H_{\omega, b}} \{\|\mathbf{x}_0 - \mathbf{x}\|_{\ell_2}\}$

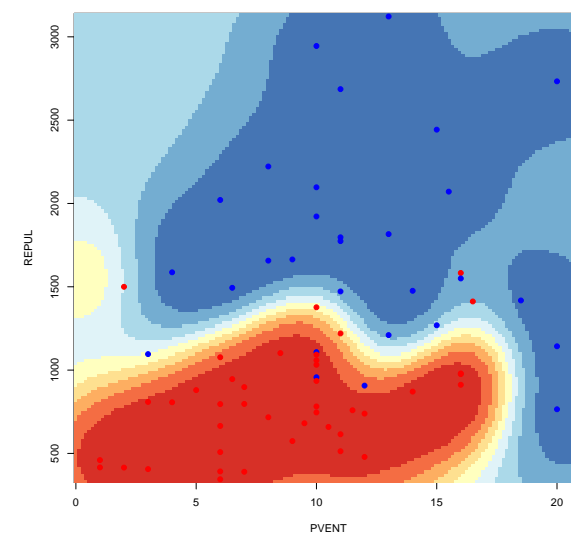
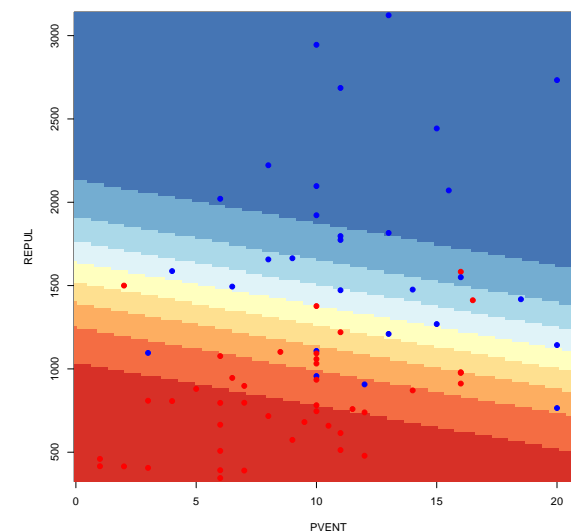
where $\|\cdot\|_{\ell_2}$ is the Euclidean (ℓ_2) norm,

$$\begin{aligned} \|\mathbf{x}_0 - \mathbf{x}\|_{\ell_2} &= \sqrt{(\mathbf{x}_0 - \mathbf{x}) \cdot (\mathbf{x}_0 - \mathbf{x})} \\ &= \sqrt{\mathbf{x}_0 \cdot \mathbf{x}_0 - 2\mathbf{x}_0 \cdot \mathbf{x} + \mathbf{x} \cdot \mathbf{x}} \end{aligned}$$

More generally, $d(\mathbf{x}_0, H_{\omega, b}) = \min_{\mathbf{x} \in H_{\omega, b}} \{\|\mathbf{x}_0 - \mathbf{x}\|_k\}$

where $\|\cdot\|_k$ is some kernel-based norm,

$$\|\mathbf{x}_0 - \mathbf{x}\|_k = \sqrt{k(\mathbf{x}_0, \mathbf{x}_0) - 2k(\mathbf{x}_0, \mathbf{x}) + k(\mathbf{x}, \mathbf{x})}$$



Regression?

In statistics, regression analysis is a statistical process for estimating the relationships among variables [...]. In a narrower sense, regression may refer specifically to the estimation of continuous response variables, as opposed to the discrete response variables used in classification. (Source: [wikipedia](#)).

Here **regression** is opposed to classification (as in the CART algorithm). y is either a continuous variable $y \in \mathbb{R}$ or a counting variable $y \in \mathbb{N}$.

In many cases in econometric and actuarial literature we *simply* want a good fit for the **conditional expectation**, $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$.

Linear, Non-Linear and Generalized Linear

$$(Y|\mathbf{X} = \mathbf{x}) \sim \mathcal{N}(\theta_{\mathbf{x}}, \sigma^2) \quad (Y|\mathbf{X} = \mathbf{x}) \sim \mathcal{N}(\theta_{\mathbf{x}}, \sigma^2) \quad (Y|\mathbf{X} = \mathbf{x}) \sim \mathcal{L}(\theta_{\mathbf{x}}, \varphi)$$

$$\mathbb{E}[Y|\mathbf{X} = \mathbf{x}] = \mathbf{x}^\top \boldsymbol{\beta} \quad \mathbb{E}[Y|\mathbf{X} = \mathbf{x}] = h(\mathbf{x}) \quad \mathbb{E}[Y|\mathbf{X} = \mathbf{x}] = h(\mathbf{x})$$

Regression Smoothers, *natura non facit saltus*

In statistical learning procedures, a key role is played by **basis functions**. We will see that it is common to assume that

$$m(\mathbf{x}) = \sum_{j=0}^k \beta_j h_j(\mathbf{x}),$$

where h_0 is usually a constant function and h_j defined basis functions.

For instance, $h_m(x) = x^j$ for a polynomial expansion with a single predictor, or $h_j(x) = (x - s_j)_+$ for some knots s_j 's (for linear splines, but one can consider quadratic or cubic ones).

Regression Smoothers: Polynomial or Spline

Stone-Weierstrass theorem every continuous function defined on a closed interval $[a, b]$ can be uniformly approximated as closely as desired by a polynomial function

Use also spline functions, e.g. piecewise linear

$$h(x) = \beta_0 + \sum_{j=1}^k \beta_j (x - s_j)_+$$

Linear Model

Consider some linear model $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$ for all $i = 1, \dots, n$.

Assume that ε_i are i.i.d. with $\mathbb{E}(\varepsilon) = 0$ (and finite variance). Write

$$\underbrace{\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}}_{\mathbf{y}, n \times 1} = \underbrace{\begin{pmatrix} 1 & x_{1,1} & \cdots & x_{1,k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & \cdots & x_{n,k} \end{pmatrix}}_{\mathbf{X}, n \times (k+1)} \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}}_{\boldsymbol{\beta}, (k+1) \times 1} + \underbrace{\begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}}_{\boldsymbol{\varepsilon}, n \times 1}.$$

Assuming $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbb{I})$, the maximum likelihood estimator of $\boldsymbol{\beta}$ is

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin}\{\|\mathbf{y} - \mathbf{X}^\top \boldsymbol{\beta}\|_{\ell_2}\} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

... under the assumption that $\mathbf{X}^\top \mathbf{X}$ is a full-rank matrix.

What if $\mathbf{X}_i^\top \mathbf{X}$ cannot be inverted? Then $\hat{\boldsymbol{\beta}} = [\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{X}^\top \mathbf{y}$ does not exist, but $\hat{\boldsymbol{\beta}}_\lambda = [\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}]^{-1} \mathbf{X}^\top \mathbf{y}$ always exist if $\lambda > 0$.

Ridge Regression

The estimator $\hat{\beta} = [\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}]^{-1} \mathbf{X}^\top \mathbf{y}$ is the **Ridge** estimate obtained as solution of

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n [y_i - \beta_0 - \mathbf{x}_i^\top \beta]^2 + \lambda \underbrace{\|\beta\|_{\ell_2}}_{\mathbf{1}^\top \beta^2} \right\}$$

for some tuning parameter λ . One can also write

$$\hat{\beta} = \underset{\beta; \|\beta\|_{\ell_2} \leq s}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^\top \beta\|_{\ell_2} \}$$

Remark Note that we solve $\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \operatorname{objective}(\beta) \}$ where

$$\operatorname{objective}(\beta) = \underbrace{\mathcal{L}(\beta)}_{\text{training loss}} + \underbrace{\mathcal{R}(\beta)}_{\text{regularization}}$$

Going further on sparsity issues

In several applications, k can be (very) large, but a lot of features are just noise: $\beta_j = 0$ for many j 's. Let s denote the number of relevant features, with $s \ll k$, cf [Hastie, Tibshirani & Wainwright \(2015\)](#),

$$s = \text{card}\{\mathcal{S}\} \text{ where } \mathcal{S} = \{j; \beta_j \neq 0\}$$

The model is now $y = \mathbf{X}_{\mathcal{S}}^{\top} \boldsymbol{\beta}_{\mathcal{S}} + \varepsilon$, where $\mathbf{X}_{\mathcal{S}}^{\top} \mathbf{X}_{\mathcal{S}}$ is a full rank matrix.

Going further on sparsity issues

Define $\|\mathbf{a}\|_{\ell_0} = \sum \mathbf{1}(|a_i| > 0)$. Ici $\dim(\boldsymbol{\beta}) = s$.

We wish we could solve

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}; \|\boldsymbol{\beta}\|_{\ell_0} \leq s}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^T \boldsymbol{\beta}\|_{\ell_2} \}$$

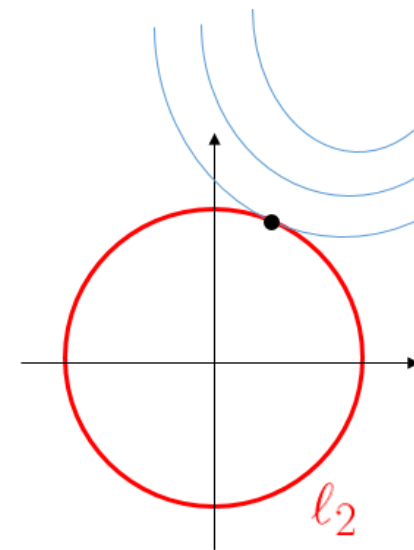
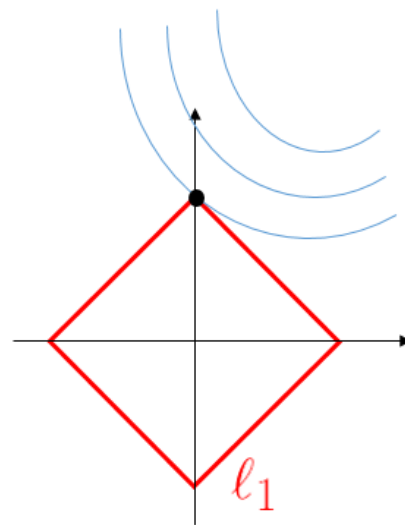
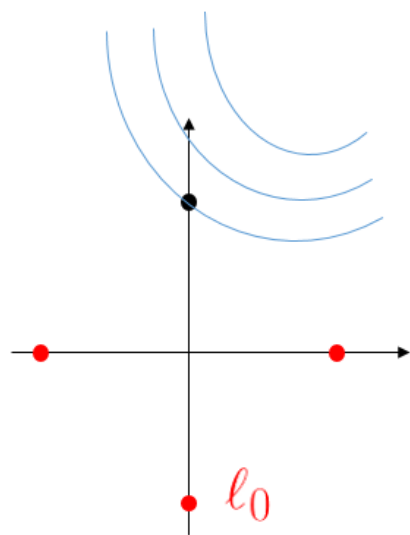
Problem: it is usually not possible to describe all possible constraints, since $\binom{s}{k}$ coefficients should be chosen here (with k (very) large).

Idea: solve the dual problem

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}; \|\mathbf{Y} - \mathbf{X}^T \boldsymbol{\beta}\|_{\ell_2} \leq h}{\operatorname{argmin}} \{ \|\boldsymbol{\beta}\|_{\ell_0} \}$$

where we might convexify the ℓ_0 norm, $\|\cdot\|_{\ell_0}$.

Regularization l_0 , l_1 et l_2



Going further on sparsity issues

On $[-1, +1]^k$, the convex hull of $\|\beta\|_{\ell_0}$ is $\|\beta\|_{\ell_1}$

On $[-a, +a]^k$, the convex hull of $\|\beta\|_{\ell_0}$ is $a^{-1}\|\beta\|_{\ell_1}$

Hence,

$$\hat{\beta} = \underset{\beta; \|\beta\|_{\ell_1} \leq \tilde{s}}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^T \beta\|_{\ell_2} \}$$

is equivalent (Kuhn-Tucker theorem) to the Lagrangian optimization problem

$$\hat{\beta} = \operatorname{argmin} \{ \|\mathbf{Y} - \mathbf{X}^T \beta\|_{\ell_2} + \lambda \|\beta\|_{\ell_1} \}$$

LASSO *Least Absolute Shrinkage and Selection Operator*

$$\hat{\beta} \in \operatorname{argmin}\{\|\mathbf{Y} - \mathbf{X}^T \beta\|_{\ell_2} + \lambda \|\beta\|_{\ell_1}\}$$

is a convex problem (several algorithms^{*}), but not strictly convex (no unicity of the minimum). Nevertheless, predictions $\hat{\mathbf{y}} = \mathbf{x}^T \hat{\beta}$ are unique

^{*} MM, minimize majorization, coordinate descent [Hunter \(2003\)](#).

Optimal LASSO Penalty

Use cross validation, e.g. K -fold,

$$\hat{\beta}_{(-k)}(\lambda) = \operatorname{argmin} \left\{ \sum_{i \notin \mathcal{I}_k} [y_i - \mathbf{x}_i^\top \beta]^2 + \lambda \|\beta\| \right\}$$

then compute the sum of the squared errors,

$$Q_k(\lambda) = \sum_{i \in \mathcal{I}_k} [y_i - \mathbf{x}_i^\top \hat{\beta}_{(-k)}(\lambda)]^2$$

and finally solve

$$\lambda^* = \operatorname{argmin} \left\{ \bar{Q}(\lambda) = \frac{1}{K} \sum_k Q_k(\lambda) \right\}$$

Note that this might overfit, so [Hastie, Tibshiriani & Friedman \(2009\)](#) suggest the largest λ such that

$$\bar{Q}(\lambda) \leq \bar{Q}(\lambda^*) + \operatorname{se}[\lambda^*] \quad \text{with} \quad \operatorname{se}[\lambda]^2 = \frac{1}{K^2} \sum_{k=1}^K [Q_k(\lambda) - \bar{Q}(\lambda)]^2$$

Boosting

Boosting is a machine learning ensemble meta-algorithm for reducing bias primarily and also variance in supervised learning, and a family of machine learning algorithms which convert weak learners to strong ones. (Source: [wikipedia](#))

The heuristic is simple: we consider an iterative process where we keep modeling the errors.

Fit model for \mathbf{y} , $m_1(\cdot)$ from \mathbf{y} and \mathbf{X} , and compute the error, $\boldsymbol{\varepsilon}_1 = \mathbf{y} - m_1(\mathbf{X})$.

Fit model for $\boldsymbol{\varepsilon}_1$, $m_2(\cdot)$ from $\boldsymbol{\varepsilon}_1$ and \mathbf{X} , and compute the error, $\boldsymbol{\varepsilon}_2 = \boldsymbol{\varepsilon}_1 - m_2(\mathbf{X})$, etc. Then set

$$m(\cdot) = \underbrace{m_1(\cdot)}_{\sim \mathbf{y}} + \underbrace{m_2(\cdot)}_{\sim \boldsymbol{\varepsilon}_1} + \underbrace{m_3(\cdot)}_{\sim \boldsymbol{\varepsilon}_2} + \cdots + \underbrace{m_k(\cdot)}_{\sim \boldsymbol{\varepsilon}_{k-1}}$$

Boosting

With (very) general notations, we want to solve

$$m^* = \operatorname{argmin}\{\mathbb{E}[\ell(Y, m(\mathbf{X}))]\}$$

for some loss function ℓ .

It is an iterative procedure: assume that at some step k we have an estimator $m_k(\mathbf{X})$. Why not constructing a new model that might improve our model,

$$m_{k+1}(\mathbf{X}) = m_k(\mathbf{X}) + h(\mathbf{X}).$$

What $h(\cdot)$ could be?

Boosting

In a perfect world, $h(\mathbf{X}) = \mathbf{y} - m_k(\mathbf{X})$, which can be interpreted as a residual.

Note that this residual is the gradient of $\frac{1}{2}[y - m(\mathbf{x})]^2$

A gradient descent is based on Taylor expansion

$$\underbrace{f(\mathbf{x}_k)}_{\langle f, \mathbf{x}_k \rangle} \sim \underbrace{f(\mathbf{x}_{k-1})}_{\langle f, \mathbf{x}_{k-1} \rangle} + \underbrace{(\mathbf{x}_k - \mathbf{x}_{k-1})}_{\alpha} \underbrace{\nabla f(\mathbf{x}_{k-1})}_{\langle \nabla f, \mathbf{x}_{k-1} \rangle}$$

But here, it is different. We claim we can write

$$\underbrace{f_k(\mathbf{x})}_{\langle f_k, \mathbf{x} \rangle} \sim \underbrace{f_{k-1}(\mathbf{x})}_{\langle f_{k-1}, \mathbf{x} \rangle} + \underbrace{(f_k - f_{k-1})}_{\beta} \underbrace{?}_{\langle f_{k-1}, \nabla \mathbf{x} \rangle}$$

where ? is interpreted as a ‘gradient’.

Boosting

Here, f_k is a $\mathbb{R}^d \rightarrow \mathbb{R}$ function, so the gradient should be in such a (big) functional space \rightarrow want to approximate that function.

$$m_k(\mathbf{x}) = m_{k-1}(\mathbf{x}) + \operatorname{argmin}_{f \in \mathcal{F}} \left\{ \sum_{i=1}^n \ell(Y_i, m_{k-1}(\mathbf{x}) + f(\mathbf{x})) \right\}$$

where $f \in \mathcal{F}$ means that we seek in a class of **weak learner functions**.

If learner are too strong, the first loop leads to some fixed point, and there is no learning procedure, see linear regression $y = \mathbf{x}^\top \boldsymbol{\beta} + \varepsilon$. Since $\varepsilon \perp \mathbf{x}$ we cannot learn from the residuals.

In order to make sure that we learn **weakly**, we can use some **shrinkage parameter** ν (or collection of parameters ν_j).

Boosting with Piecewise Linear Spline & Stump Functions

Take Home Message

- Similar goal: getting a predictive model, $\hat{m}(\mathbf{x})$
- Different/Similar tools: minimize loss/maximize likelihood

$$\hat{m}(\cdot) = \operatorname{argmin}_{m(\cdot) \in \mathcal{F}} \left\{ \sum_{i=1}^n \ell(y_i, m(\mathbf{x}_i)) \right\} \text{ vs. } \hat{m}(\cdot) = \operatorname{argmax}_{m(\cdot) \in \mathcal{F}} \left\{ \sum_{i=1}^n \log f(y_i; m(\mathbf{x}_i)) \right\}$$

- Try to remove the noise and avoid overfit using cross-validation,

$$\ell(y_i, \hat{m}_{(-i)}(\mathbf{x}_i))$$

- Use computational tricks (bootstrap) to increase robustness
- Nice tools to select interesting features (LASSO, variable importance)