Modern Trends in Data Mining

Trevor Hastie Stanford University November 2012.



How IBM built Watson, its "Jeopardy"-playing supercomputer by Dawn Kawamoto DailyFinance 02/08/2011



Learning from its mistakes According to David Ferrucci (PI of Watson DeepQA technology for IBM Research), Watson's software is wired for more than handling natural language processing.

"Machine learning allows the computer to become smarter as it tries to answer questions — and to learn as it gets them right or wrong."

Data Mining

For Today's Graduate, Just One Word: Statistics

By STEVE LOHR Published: August 5, 2009

MOUNTAIN VIEW, Calif. — At Harvard, Carrie Grimes majored in anthropology and archaeology and ventured to places like Honduras, where she studied Mayan settlement patterns by mapping where artifacts were found. But she was drawn to what she calls "all the computer and math stuff" that was part of the job.

Enlarge This Image



Thor Swift for The New York Times Carrie Grimes, senior staff engineer at Google, uses statistical analysis of data to help improve the company's search engine.

Multimedia



PND-in computer second with focus on antificial intelligence and here analysiss. Mr. 7 sectors here a research scientist as LB.M. who uses computing and modeling to exhapt including patients from text, wideo and audeo data.

17 University Bachelon's degree in computer science. Content Ph.B. in computer science M11. & MacArthur Tellow "People think of field archaeology as Indiana Jones, but much of what you really do is data analysis," she said.

Now Ms. Grimes does a different kind of digging. She works at Google,

where she uses statistical analysis of mounds of data to come up with ways to improve its search engine.

Ms. Grimes is an Internet-age statistician, one of many who are changing the image of the profession as a place for dronish number nerds. They are finding themselves increasingly in demand — and even cool.

"I keep saying that the sexy job in the next 10 years will be statisticians," said Hal Varian, chief economist at Google. "And I'm not kidding."



QUOTE OF THE DAY, NEW YORK TIMES, AUGUST 5, 2009 "I keep saying that the sexy job in the next 10 years will be statisticians (sic). And I'm not kidding." - HAL VARIAN, chief economist at Google.



- We have a collection of data pertaining to our business, industry, production process, monitoring device, etc.
- • ften the goals of data-mining are vague, such as *"look for patterns in the data"* not too helpful.
- In many cases a *"response"* or *"outcome"* can be identified as a good and useful target for prediction.
- Accurate prediction of this target can help the company make better decisions, and save a lot of money.
- Data-mining is particularly good at building such prediction models an area known as "supervised learning".

Example: Credit Risk Assessment

- Customers apply to a bank for a loan or credit card.
- They supply the bank with information such as age, income, employment history, education, bank accounts, existing debts, etc.
- The bank does further background checks to establish credit history of customer.
- Based on this information, the bank must decide whether to make the loan or issue the credit card.

- The bank has a large database of existing and past customers. Some of these defaulted on loans, others frequently made late payments etc. An outcome variable *"Status"* is defined, taking value *"good"* or *"default"*. Each of the past customers is scored with a value for status.
- Background information is available for all the past customers.
- Using supervised learning techniques, we can build a risk prediction model that takes as input the background information, and outputs a risk estimate (probability of default) for a prospective customer.

The California based company Fair-Isaac uses a generalized additive model + boosting methods in the construction of their credit risk scores.

- When a customer switches to another provider, we call this *"churn"*. Examples are cell-phone service and credit card providers.
- Based on customer information and usage patterns, we can predict
 - the probability of churn
 - the retention probability (as a function of time)
- This information can be used to evaluate
 - prospective customers to decide on acceptance
 - present customers to decide on intervention strategy

Risk assessment and survival models are used by US cell-phone companies such as AT&T to manage churn.

NETFL	NETFLIX								
Ne	Netflix Prize								
Home Rul	Home Rules Leaderboard Register Update Submit Download								
Leaderboard Display top 20 : leaders. Bank Team Name Bast Score % Improvement Last Submit Time									
1	The Ensemble	0.8553	10.10	2009-07-26 18:38:22					
2	BeliKor's Pragmatic Chaos	0.8554	10.09	2009-07-26 18:18:28					
Grand	Grand Prize - RMSE <= 0.8563								
3	Grand Prize Team	0.8571	9.91	2009-07-24 13:07:49					
4	Opera Solutions and Vandelay United	0.8573	9.89	2009-07-25 20:05:52					
5	Vandelay Industries !	0.8579	9.83	2009-07-26 02:49:53					
6	Pragmatic Theory	0.8582	9.80	2009-07-12 15:09:53					
7	BellKor in BlgChaos	0.8590	9.71	2009-07-26 12:57:25					
8	Dace_	0.8603	9.58	2009-07-24 17:18:43					
9	Opera Solutions	0.8611	9.49	2009-07-26 18:02:08					
10	BellKor	0.8612	9.48	2009-07-26 17:19:11					

Grand Prize: one million dollars, if beat Netflix's RMSE by 10%. Competition ends Sep 21, 2009 after ≈ 3 years, two leaders, 41305 teams! Ultimate winner is BellKor's Pragmatic Chaos.



Netflix users rate movies from 1-5. Based on a history of ratings, predict the rating a viewer will give to a new movie.

- Training data: sparse 400K (users) by 18K (movies) rating matrix, with 98.7% missing. About 100M movie/rater pairs.
- Quiz set of about 1.4M movie/viewer pairs, for which predictions of ratings are required (Netflix held them back)
- Probe set of about 1.4 million movie/rater pairs similar in composition to the quiz set, for which the ratings are known.
- Both winning teams used ensemble methods to achieve their results.

The Supervised Learning Problem

Starting point:

- Outcome measurement Y (also called dependent variable, response, target, output)
- Vector of *p* predictor measurements X (also called inputs, regressors, covariates, features, independent variables)
- In the *regression problem*, Y is quantitative (e.g price, **b**lood pressure, rating)
- In *classification*, Y takes values in a finite, unordered set (default yes/no, churn/retain, spam/email)
- We have training data $(x_1, y_1), \ldots, (x_N, y_N)$. These are observations (examples, instances) of these measurements.

Objectives

 \bullet n the basis of the training data we would like to:

- Accurately predict unseen test cases for which we know X but do not know Y.
- In the case of classification, predict the probability of an outcome.
- Understand which inputs affect the outcome, and how.
- Assess the quality of our predictions and inferences.

More Examples

- Predict whether someone will have a heart attack on the basis of demographic, diet and clinical measurements
- Determine whether an incoming email is "spam", based on frequencies of key words in the message
- Identify the numbers in a handwritten zip code, from a digitized image
- Estimate the probability that an insurance claim is fraudulent,
 based on client demographics, client history, and the amount
 and nature of the claim.
- Predict the type of cancer in a tissue sample using DNA expression values

Email or Spam?

- data from 4601 emails sent to an individual (named George, at HP labs, before 2000). Each is labeled as *"spam"* or *"email"*.
- goal: build a customized spam filter.
- input features: relative frequencies of 57 of the most commonly occurring words and punctuation marks in these email messages.

	george	you	hp	free		edu	remove
spam	0.00	2.26	●.●2	● .52	●.51	●.●1	0.28
email	1.27	1.27	0.90	●.●7	●.11	0.29	●.●1

Average percentage of words or characters in an email message equal to the indicated word or character. We have chosen the words and characters showing the largest difference between spam and email.



A sample of segmented and normalized handwritten digits, scanned from zip-codes on envelopes. Each image has 16×16 pixels of grayscale values ranging from 0 - 255.



Market Ma Microarray Cancer Data

Expression matrix of 6830 genes (rows) and 64 samples (columns), for the human tumor data (100 randomly chosen rows shown). The display is a heat map, ranging from bright green (under expressed) to bright red (over expressed).

Geal: predict cancer class based on expression values. Springer Series in Statistics

Trevor Hastie Robert Tibshirani Jerome Friedman

The Elements of Statistical Learning

Data Mining, Inference, and Prediction

Second Edition

🗋 Springer

Shameless self-promotion

All of the topics in this lecture are covered in the 2009 second edition of our 2001 book.

The book blends traditional linear methods with contemporary nonparametric methods, and many between the two.

Ideal "Bayes" Predictions

• For a *quantitative output* Y, the best prediction we can make when the input vector X = x is

$$f(x) = \operatorname{Ave}(Y|X = x)$$

- This is the conditional expectation deliver the Y-average of all those examples having X = x.
- This is best if we measure errors by average squared error $Ave(Y f(X))^2$.
- For a *qualitative output* Y taking values $1, 2, \ldots, M$, compute
 - Pr(Y = m | X = x) for each value of m. This is the conditional probability of class m at X = x.
 - Classify C(x) = j if Pr(Y = j | X = x) is the largest the majority vote classifier.

Implementation with Training Data

The ideal prediction formulas suggest a data implementation. To predict at X = x, gather all the training pairs (x_i, y_i) having $x_i = x$, then:

- For regression, use the mean of their y_i to estimate $f(x) = \operatorname{Ave}(Y|X = x)$
- For classification, compute the relative proportions of each class among these y_i , to estimate $\Pr(Y = m | X = x)$; Classify the new observation by majority vote.

Problem: in the training data, there may be N \bullet observations having $x_i = x$.

Nearest Neighbor Averaging

• Estimate $\operatorname{Ave}(Y|X = x)$ by

Averaging those y_i whose x_i are in a neighborhood of x.

- E.g. define the neighborhood to be the set of k observations having values x_i closest to x in euclidean distance $||x_i x||$.
- For classification, compute the class proportions among these k closest points.
- Nearest neighbor methods often outperform all other methods
 about one in three times especially for classification.





- Smooth version of nearestneighbor averaging
- At each point x, the function
 f(x) = Y(Y|X = x) is estimated by the weighted average of the y's.
- The weights die down smoothly with distance from the target point x (indicated by shaded orange region).

 * not to be confused with "kernel methods" as in SVMs

Structured Models

• When we have a lot of predictor variables, NN methods often fail because of the "curse of dimensionality"

It is hard to find nearby points in high dimensions!

- Near-neighbor models offer little interpretation.
- We can overcome these problems by assuming some structure for the regression function $\operatorname{Ave}(Y|X = x)$ or the probability function $\Pr(Y = k|X = x)$. Typical structural assumptions:
 - Linear Models
 - Additive Models
 - Low-order interaction models
 - Restrict attention to a subset of predictors
 - $-\ldots$ and many more



• Linear models assume

$$\operatorname{Ave}(Y|X=x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p$$

• For two class classification problems, linear logistic regression has the form

$$\log \frac{\Pr(Y = +1 | X = x)}{\Pr(Y = -1 | X = x)} = \mathbf{\beta}_{\bullet} + \mathbf{\beta}_{1} x_{1} + \mathbf{\beta}_{2} x_{2} + \ldots + \mathbf{\beta}_{p} x_{p}$$

• This translates t•

$$\Pr(Y = +1|X = x) = \frac{e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}}$$

Chapters 3 and 4 of deal with linear models.

Linear Model Complexity Control

With many inputs, linear regression can overfit the training data, leading to poor predictions on future data. Two general remedies are available:

- Variable selection: reduce the number of inputs in the model. For example, stepwise selection or *best subset selection*.
- Regularization: leave all the variables in the model, but when fitting the model, restrict their coefficients.
 - Ridge: $\sum_{j=1}^{p} \beta_{j}^{2} \leq s$. All the coefficients are non-zero, but are shrunk toward zero (and each other).
 - Lasso: $\sum_{j=1}^{p} |\mathbf{s}_j| \leq s$. Some coefficients drop out the model, others are shrink toward zero.

data. The red models are the candidates, and we need to choose s. variables, and shows the residual sum-of-squares on the training Each point corresponds to a linear model involving a subset of the









Lasso

Overfitting and Model Assessment

- In all cases above, the larger s, the better we will fit the training data.
- • •verfit models can perform poorly on test data (high variance).
- Underfit models can perform poorly on test data (high bias).

Model assessment aims to

- 1. Choose a value for a tuning parameter s for a technique.
- 2. Estimate the future prediction ability of the chosen model.
- For both of these purposes, the best approach is to evaluate the procedure on an independent test set, if one is available.
- If possible one should use different test data for (1) and (2) above: a *validation set* for (1) and a *test set* for (2)

K-Fold Cross-Validation

Primarily a method for estimating a tuning parameter s when data are scarce; we illustrate for the regularized linear regression models.

• Divide the data into K roughly equal parts (5 or 10) 1 2 3 4 5

-	-		-	e
Train	Train	Validation	Train	Train

- for each k = 1, 2, ..., K, fit the model with parameter s to the other K-1 parts, giving $\hat{\beta}^{-k}(s)$ and compute its error in predicting the kth part: $E_k(\lambda) = \sum_{i \in k \text{th part}} (y_i x_i^T \hat{\beta}^{-k}(s))^2$.
- This gives the overall cross-validation error $CV(s) = \frac{1}{K} \sum_{k=1}^{K} E_k(s)$
- do this for many values of s and choose the value of s that makes CV(s) smallest.

Cross-Validation Error Curve



- 10-fold CV error curve using lasso on some diabetes data (64 inputs, 442 samples).
- Thick curve is CV error curve
- Shaded region indicates standard error of CV estimate.
- Curve shows effect of overfitting — errors start to increase above s = 0.2.
- This shows a trade-off between bias and variance.

Modern Structured Models in Data Mining

The following is a list of some of the more important and currently popular prediction models in data mining.

- Linear Models (often heavily *regularized*)
- Generalized Additive Models
- Neural Networks
- Hierarchical Bayesian Prediction Models
- Trees, Random Forests and Boosted Tree Models hot!
- Support Vector and Kernel Machines hot!

Generalized Additive Models

Allow a compromise between linear models and more flexible local models (kernel estimates) when there are a many inputs $X = (X_1, X_2, \ldots, X_p).$

• Additive models for regression:

Ave
$$(Y|X = x) = \alpha_0 + f_1(x_1) + f_2(x_2) + \ldots + f_p(x_p).$$

• Additive models for classification:

$$\log \frac{\Pr(Y = +1 | X = x)}{\Pr(Y = -1 | X = x)} = \alpha_0 + f_1(x_1) + f_2(x_2) + \ldots + f_p(x_p).$$

Each of the functions $f_j(x_j)$ (one for each input variable), can be a *smooth function* (e.g. kernel estimate), *linear*, or *omitted*.



GAM fit to SPAM data

- Shown are the most important predictors.
- Many show nonlinear behavior.
- Overall error rate 5.3%.
- Functions can be reparametrized (e.g. log terms, quadratic, step-functions), and then fit by linear model.
- Produces a prediction per email $\Pr(\text{SPAM}|X = x)$

Neural Networks



Single (Hidden) Layer Perceptron

- Like a complex regression or logistic regression model more flexible, but less interpretable a *"black box"*.
- Hidden units Z_1, Z_2, \ldots, Z_m (4 here): $Z_j = \sigma(\alpha_{0j} + \alpha_j^T X)$ $\sigma(Z) = e^Z/(1 + e^Z)$ is the logistic sigmoid *activation* function.
- Output is a linear regression or logistic regression model in the Z_j .
- Complexity controlled by *m*, ridge regularization, and *early stopping* of the *backpropagation algorithm* for fitting the neural network.

Support Vector Machines



- Maximize the gap (margin) between the two classes on the training data.
- If not separable
 - enlarge the feature space via basis expansions (e.g. polynomials).
 - use a *"soft"* margin (allow limited overlap).
- Solution depends on a small number of points (*"support vectors"*) — 3 here.

Support Vector Machines



- Maximize the *soft margin* subject to a bound on the total overlap: $\sum_i \xi_i \leq B$.
- Even if data are separable, wider soft margin more stable.
- Primarily used for classification problems. Builds a linear classifier f(X) = β₀ + β^TX
 If f(X) > 0, classify as +1, else if f(X) < 0, classify as -1.
- Generalizations use kernels: $f(X) = \alpha_0 + \sum_{i=1}^N \alpha_i K(X, x_i)$

Classification and Regression Trees

- ✓ Can handle huge datasets
- \checkmark Can handle *mixed* predictors—quantitative and qualitative
- ✓ Easily ignore redundant variables
- ✓ Handle missing data elegantly
- \checkmark Small trees are easy to interpret
- **X** Large trees are hard to interpret
- \mathbf{X} •ften prediction performance is poor





Ensemble Methods and Boosting

Classification trees can be simple, but often produce noisy (bushy) or weak (stunted) classifiers.

- Bagging (Breiman, 1996): Fit many large trees to bootstrap-resampled versions of the training data, and classify by majority vote.
- Random Forests (Breiman 1999): Improvements over bagging.
- Boosting (Freund & Shapire, 1996): Fit many smallish trees to reweighted versions of the training data. Classify by weighted majority vote.

In general Boosting \succ Random Forests \succ Bagging \succ Single Tree.



Spam Data

Number of Trees

Modern Gradient Boosting (Friedman, 2001)

• Fits an additive model

 $F_m(X) = T_1(X) + T_2(X) + T_3(X) + \ldots + T_m(X)$

where each of the $T_j(X)$ is a *tree in* X.

- Can be used for regression, logistic regression and more. For example, gradient boosting for regression works by repeatedly fitting trees to the residuals:
 - 1. Fit a small tree $T_1(X)$ to Y.
 - 2. Fit a small tree $T_2(X)$ to the residual $Y T_1(X)$.
 - 3. Fit a small tree $T_3(X)$ to the residual $Y T_1(X) T_2(X)$. and so on.
- *m* is the tuning parameter, which must be chosen using a validation set (*m* too big will overfit).

Software

- *R* is free software for statistical modeling, graphics and a general programming environment. Works on PCs, Macs and Linux/Unix platforms. All the models here can be fit in R. R grew from its predecessor Splus, and both implement the S language developed at Bell Labs in the 80s.
- *SAS* and their *Enterprise Miner* can fit most of the models mentioned in this talk, with good data-handling capabilities, and high-end user interfaces.
- *Salford Systems* has commercial versions of trees, random forests and gradient boosting.
- SVM software is all over, but beware of patent infringements if put to commercial use.
- Many free versions of neural network software; *Google* will find.



- Many amazing tools are available, from the simplest linear models to complex boosting algorithms.
- Avoid unwarranted complexity; if linear models perform well, they are easier to manage than more complex models.
- Boosting provides a good benchmark for what performance might be achievable.
- A good software environment is essential; if R can manage your problem size, its a great environment.