

# Chapters 1 & 2. Introduction & Overview

Wei Pan

Division of Biostatistics, School of Public Health, University of Minnesota,  
Minneapolis, MN 55455

Email: [weip@biostat.umn.edu](mailto:weip@biostat.umn.edu)

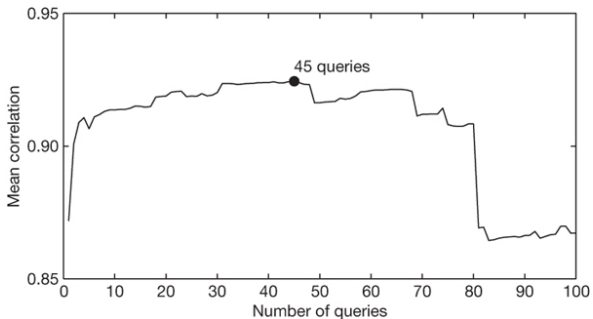
PubH 7475/8475

©Wei Pan

# Big Data and ML/DL/AI

- ▶ *Big Data is on the rise, bringing big questions* (WSJ, 11-29-2012)
- ▶ *Big data: the next frontier for innovation, competition, and productivity* (McKinsey report 05-2011)
- ▶ *Big Data's big problem: little talent* (WSJ, 04-29-2012)
- ▶ Example: Google Flu Trends (GFT); “Nowcast”  
Ginsberg et al (2009). Detecting influenza epidemics using search engine query data. *Nature*, 457:1012–1014. <https://www.nature.com/articles/nature07634#MOESM269>

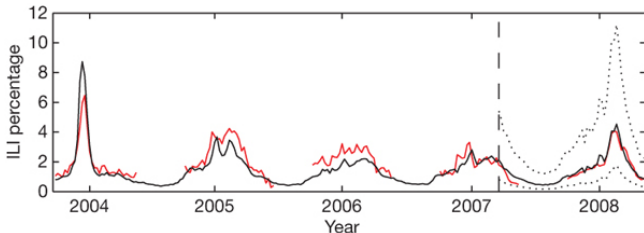
An evaluation of how many top-scoring queries to include in the ILLI-related query fraction.



J Ginsberg *et al. Nature* **000**, 1-3 (2008) doi:10.1038/nature07634

nature

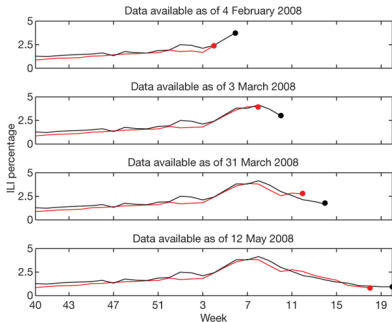
A comparison of model estimates for the mid-Atlantic region (black) against CDC-reported ILI percentages (red), including points over which the model was fit and validated.



J Ginsberg *et al. Nature* **000**, 1-3 (2008) doi:10.1038/nature07634

nature

ILI percentages estimated by our model (black) and provided by the CDC (red) in the mid-Atlantic region, showing data available at four points in the 2007-2008 influenza season.



J Ginsberg *et al. Nature* **000**, 1-3 (2008) doi:10.1038/nature07634

- ▶ Impressive!
- ▶ But now? GFT was long gone...  
Lazar et al. (2014). The Parable of Google Flu: Traps in Big Data Analysis. *Science*, 343(6176): 1203-5.  
All-data!
- ▶ Yang et al. (2015). Accurate estimation of influenza epidemics using Google search data via ARGO. *PNAS*, 112: 14473-8.

$$y_t = y_{t-2}\beta + \sum_{g \in G} g\beta_g + \epsilon, \quad \sum_g |\beta_g| < T.$$

- ▶ Aiken et al. (2020). Real-time estimation of disease activity in emerging outbreaks using internet search information. *PLOS Comp Biol*, 16(8): e1008117.

# Data Science

- ▶ “Data Science” (Cleveland 2001/2014)
- ▶ How is this related to statistics?  
Breiman (2001). Statistical modeling: the two cultures. *Stat Sci*, 16:199-231.  
Inference vs prediction!  
Change and expand the subjects
- ▶ Computing:  
Hadoop (or RHadoop), MapReduce, Spark, ...
- ▶ You do not need to do everything ...  
DeltaRho (formerly, Tessera): interface b/w R and Hadoop...  
<http://deltarho.org/>  
R packages `datadr`, `trelliscope`  
Based on “Divide and Recombine” (D&R) (Guha et al 2012).
- ▶ So ...still need to go back to the basics of ...!

- ▶ *Harvard Business Review* Oct 30, 2019: “AI Can Outperform Doctors...”  
*“Medical artificial intelligence (AI) can perform with expert-level accuracy and deliver cost-effective care at scale. IBM’s Watson diagnoses heart disease better than cardiologists do. Chatbots dispense medical advice for the United Kingdom’s National Health Service in lieu of nurses. Smartphone apps now detect skin cancer with expert accuracy. Algorithms identify eye diseases just as well as specialized physicians. Some forecast that medical AI will pervade 90% of hospitals and replace as much as 80% of what doctors currently do.”*
- ▶ “...So Why Don’t Patients Trust It?”



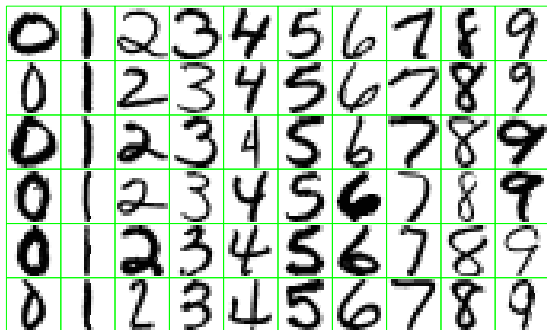
# Introduction

- ▶ Focus: prediction or discovery.  
Approach: build a model  $\hat{f}(x)$ .
- ▶ Types: supervised vs unsupervised vs semi-supervised learning.  
Training data: with vs without known response values vs a mixture of both.
- ▶ Supervised learning: classification vs regression.  
Training data:  $(Y_i, X_i)$ 's;  $Y_i$  is categorical (e.g. binary) vs quantitative.  
 $X_i$ : typically multivariate and mixed types.  
Tuning and test data:  $(Y_i, X_i)$ 's;  
Future use: only  $X_i$ 's.

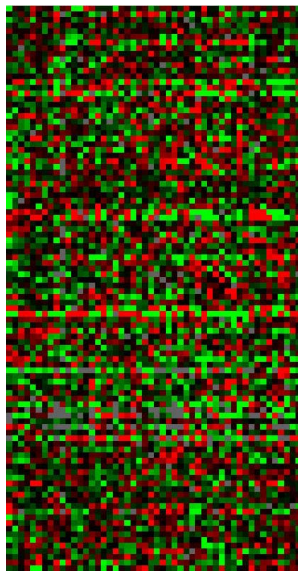
# Examples

- ▶ Example 1.  $X_i^0$ : an email;  $Y_i = 0$  or  $1$ , indicating whether it is a junk email;  $i = 1, \dots, 4601$ .
- ▶ Feature extraction: e.g. use some key words in emails as  $X_i$ ; manually.  
Automated: word/sentence embedding by DL?
- ▶ A classification problem: use a 0-1 loss, build a model  $\hat{f}(x) \in \{0, 1\}$ , calculate misclassification rate,...
- ▶ Loss function: here a false positive is much more costly than a false negative.

- ▶ Example 2. Predict prostate specific antigen (PSA) using some lab measurements.
- ▶ A regression problem.
- ▶ Example 3. Handwritten digit recognition.
- ▶  $X_i^0$ : a 16 by 16 black/white image (= a 16 by 16 binary matrix);  $Y_i \in \{1, 2, \dots, 9\}$ .
- ▶  $X_i$ : maybe (vectorized)  $X_i^0$ , or better, its summary stat's, e.g. marginal histograms or numbers of "crossing changes" ..., manually  
automated: CNNs/DL.



**FIGURE 1.2.** *Examples of handwritten digits from U.S. postal envelopes.*



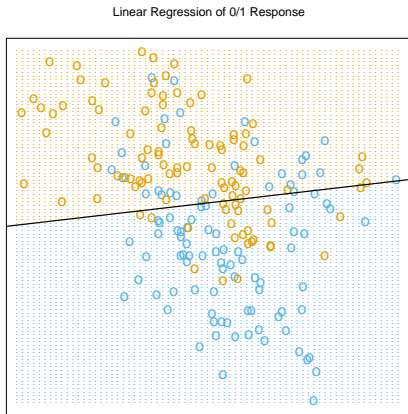
SEW029104  
 SEW030102  
 SE07161  
 GNL  
 H4888806  
 AD32334  
 RAG219AGE  
 SD02772  
 EST  
 SEW037402  
 H4888804  
 SEW449284  
 EST  
 SD471015  
 MFBPNC70  
 EST\_Chr 1  
 SD37461  
 DNAPOC1M8  
 SD33913  
 SEW01489  
 SD18117  
 SEW470409  
 SEW47261  
 H4888808  
 SEW037696  
 CH  
 MFC0CH040  
 SD47116  
 EST\_Chr 6  
 SEW036310  
 SD48817  
 SD305167  
 EST\_Chr 3  
 SD12764  
 SD3814  
 PFPBC  
 SEW036603  
 SEW010141  
 SEW07608  
 EST\_Chr 31  
 SD114241  
 SD377419  
 SD39117  
 SEW001620  
 SEW07684  
 SEW010534  
 H4888805  
 SEW020464  
 SEW03591  
 SEW006716  
 SEW07676  
 HYPOTHET  
 WAD0804  
 SEW031854  
 EST\_Chr 15  
 SEW0376384  
 SD38506  
 EST\_Chr 5  
 SD48821  
 SD46536  
 SEW02715  
 EST\_Chr 2  
 SEW02266  
 SD30034  
 EST\_Chr 15  
 SD38483  
 SD48148  
 SD297905  
 EST  
 SEW486740  
 GNAL150C  
 EST  
 SEW060311  
 SEW037197  
 SD30379  
 EST  
 SE03609  
 SEW416621  
 EST\_LINEN  
 TLPL1L1UP1  
 SEW43842  
 SD381079  
 SEW03862  
 SEW417270  
 SEW042471  
 EST\_Chr 15  
 SEW031026  
 SD380345  
 SEW038182  
 SD381406  
 SD37113  
 SEW060699  
 EST\_Chr 19  
 SEW035130  
 SD380307  
 SD379960  
 SEW130368  
 SD301802  
 SD31384  
 SD42364

SEW029104  
 SEW030102  
 SE07161  
 GNL  
 H4888806  
 AD32334  
 RAG219AGE  
 SD02772  
 EST  
 SEW037402  
 H4888804  
 SEW449284  
 EST  
 SD471015  
 MFBPNC70  
 EST\_Chr 1  
 SD37461  
 DNAPOC1M8  
 SD33913  
 SEW01489  
 SD18117  
 SEW470409  
 SEW47261  
 H4888808  
 SEW037696  
 CH  
 MFC0CH040  
 SD47116  
 EST\_Chr 6  
 SEW036310  
 SD48817  
 SD305167  
 EST\_Chr 3  
 SD12764  
 SD3814  
 PFPBC  
 SEW036603  
 SEW010141  
 SEW07608  
 EST\_Chr 31  
 SD114241  
 SD377419  
 SD39117  
 SEW001620  
 SEW07684  
 SEW010534  
 H4888805  
 SEW020464  
 SEW03591  
 SEW006716  
 SEW07676  
 HYPOTHET  
 WAD0804  
 SEW031854  
 EST\_Chr 15  
 SEW0376384  
 SD38506  
 EST\_Chr 5  
 SD48821  
 SD46536  
 SEW02715  
 EST\_Chr 2  
 SEW02266  
 SD30034  
 EST\_Chr 15  
 SD38483  
 SD48148  
 SD297905  
 EST  
 SEW486740  
 GNAL150C  
 EST  
 SEW060311  
 SEW037197  
 SD30379  
 EST  
 SE03609  
 SEW416621  
 EST\_LINEN  
 TLPL1L1UP1  
 SEW43842  
 SD381079  
 SEW03862  
 SEW417270  
 SEW042471  
 EST\_Chr 15  
 SEW031026  
 SD380345  
 SEW038182  
 SD381406  
 SD37113  
 SEW060699  
 EST\_Chr 19  
 SEW035130  
 SD380307  
 SD379960  
 SEW130368  
 SD301802  
 SD31384  
 SD42364

- ▶ Example 4. Microarray gene expression data.
- ▶  $X_i$ : 6830 genes' expression levels; quantitative;  
 $Y_i$ : tumor types.
- ▶ A typical “small  $n$ , large  $p$ ” problem:  $n = 64$  vs  $p = 6830$ .
- ▶ A classification problem.
- ▶ Can be an unsupervised learning problem: finding subtypes of cancer.  
only use  $X_i$ 's to find new class labels  $Y_i^*$ ; clustering analysis.
- ▶ Can be a semi-supervised learning problem: some known and possibly novel subtypes of cancer.

# Overview

- ▶ Consider two popular, yet simple and extreme methods: LR vs NN;  
parametric vs non-parametric.
- ▶ Q: Is a non-parametric method better than a parametric one?  
or reverse?
- ▶ Consider simulated data:  $(Y_i, X_i)$ ,  $Y_i = 0$  or  $1$  and  $X_i$   
bivariate; 100 obs's in each class (as training data).
- ▶ LR:  $E(Y_i|X_i) = Pr(Y_i = 1|X_i) = \beta_0 + X_i'\beta$ ;  
Use LS to estimate  $\beta$ 's  $\implies \hat{Y}_i = \widehat{Pr}(Y_i = 1|X_i)$ ;  
 $\tilde{Y}_i = I(\hat{Y}_i \geq 0.5)$ .
- ▶ Decision boundary:  $\hat{Y}(x) = \hat{\beta}_0 + x'\hat{\beta} = 0.5$ , linear.



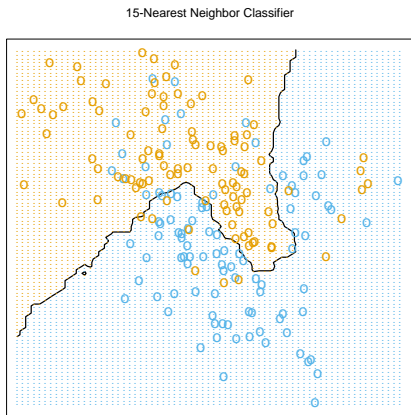
**FIGURE 2.1.** A classification example in two dimensions. The classes are coded as a binary variable (BLUE = 0, ORANGE = 1), and then fit by linear regression. The line is the decision boundary defined by  $x^T \hat{\beta} = 0.5$ . The orange shaded region denotes that part



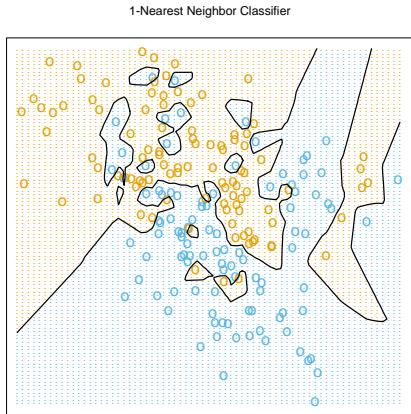
- ▶ kNN:  $N_{k(x)}$  is the  $k$  nearest training data points that are closest to  $x$ ,

$$\hat{Y}(x) = \frac{1}{k} \sum_{X_i \in N_k(x)} Y_i = \widehat{Pr}(Y_i = 1 | X_i).$$

- ▶ Idea: using local “smoothness” to estimate the population mean by ...
- ▶ Key: choice of  $k$ , or how much “smoothness” is to be assumed; do not know!  
Modeling assumption: larger  $k$ , higher or lower model complexity?
- ▶ Try a few values of  $k$ , then ...

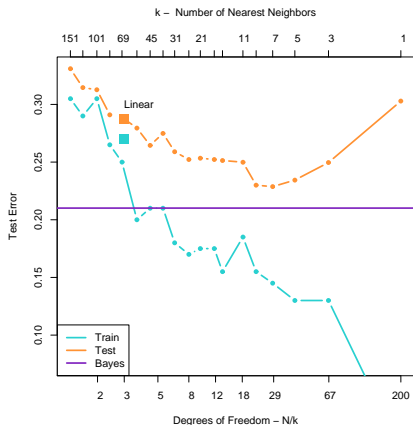


**FIGURE 2.2.** The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (BLUE = 0, ORANGE = 1) and then fit by 15-nearest-neighbor averaging as in (2.8). The pre-



**FIGURE 2.3.** The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (BLUE = 0, ORANGE = 1), and then

- ▶ Key Q: which kNN (and LR) to use?
- ▶ Key: cannot use the training data to compare models!  
Why not? too optimistic, favoring ...  
Recall: how to estimate the noise variance in linear regression?
- ▶ How? use a separate test dataset, or CV, or some model selection criterion (if any).  
Key: test data should **not** be used in model building!  
Q: how about AIC, BIC ,...
- ▶ Previous example: *generate a new test dataset* with  $n = 10,000$ .



**FIGURE 2.4.** Misclassification curves for the simulation example used in Figures 2.1, 2.2 and 2.3. A single training sample of size 200 was used, and a test sample of size 10,000. The orange curves are test and the blue are training error for  $h$  nearest neighbor classification.

- ▶ Q: is there a best classifier?
- ▶ Ideal situation: if we know the data distribution, then use the Bayes rule:

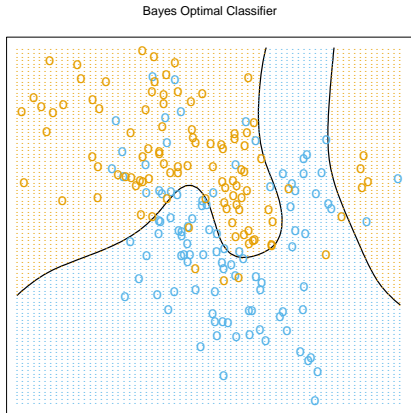
$$k_0 = \arg \max_k Pr(k|x).$$

- ▶ An example: 1) prior  $\pi_k = Pr(k)$ ; 2) PDF of class  $k$ ,  $f_k(x) = f(x|k)$ , then

$$Pr(k|x) = \frac{\pi_k f_k(x)}{\sum_i \pi_i f_i(x)}.$$

If  $f_k$  is assumed to be Normal, then LDA or QDA.  
LR and kNN are also estimating  $Pr(k|x)$ .

- ▶ Bayes rule: offering a theoretical lower bound of the test error rate; often unknown.
- ▶ Previous example: R code example 2.1.



**FIGURE 2.5.** *The optimal Bayes decision boundary for the simulation example of Figures 2.1, 2.2 and 2.3. Since the generating density is known for each class, this boundary can be calculated exactly (Exercise 2.2).*

- ▶ Q: for real data, often cannot generate new data; how to evaluate models?
- ▶ Use sample splitting: divide the original whole dataset into two parts, (e.g. 1/2 or 2/3) for training and (the remaining) for test.  
efficient?
- ▶ Use cross-validation (CV); read §7.10
- ▶ K-fold CV: Divide the data  $D$  into almost equally sized and none-overlapping  $D_1, \dots, D_K$ , then

$$\text{CVerr} = \sum_{j=1}^K \sum_{(Y_i, X_i) \in D_j} L[Y_i, \hat{f}(X_i | D - D_j)] / n.$$

- ▶ Leave-One-Out-CV (LOOCV):  $K = n$ .
- ▶ Remarks: 1) not necessarily larger  $K$ , the better; CV related to AIC/BIC; 2) maybe better to use bootstrap (§7.11).
- ▶ Previous example: R code example 2.1.

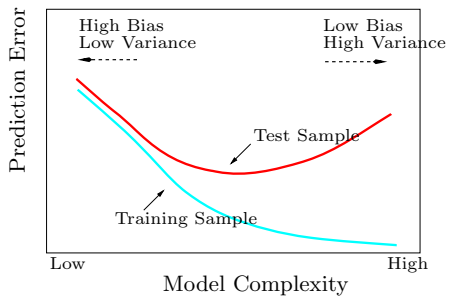


- ▶ Example 1: goal: need to train a model and estimate its test error.  
Way 1: splitting the whole sample into a training subset and a test subset;  
Way 2: applying CV to the whole sample.
- ▶ Example 2: goal: need to tune/select a model.  
Way 1: splitting the whole sample into training and tuning subsets;  
Way 2: applying CV to the whole sample.
- ▶ Example 3: goal: need to tune/select a model and estimate the test error.  
Way 1: splitting the whole sample into training, tuning and test subsets;  
Way 2: splitting the whole sample into training and test subsets, then using CV on the training subset.
- ▶ After applying CV to a training dataset (to select the tuning parameter value), one would often refit the model to the whole training dataset with the selected tuning parameter value. Why?

- ▶ Key: celebrated bias-variance trade-off!
- ▶ Suppose  $\hat{f}$  is any estimate of  $f$ ,

$$\begin{aligned}MSE &= E[(\hat{f} - f)^2] = E[(\hat{f} - E(\hat{f}) + E(\hat{f}) - f)^2] \\ &= E[(\hat{f} - E\hat{f})^2] + E[(E\hat{f} - f)^2] \\ &= \text{Var} + \text{Bias}^2.\end{aligned}$$

- ▶ Very very useful: helps explain
  - i) Complex models vs simple models;
  - ii) Nonparametrics vs parametrics; ...
- ▶ Perhaps the most important plot in the course:



**FIGURE 2.11.** *Test and training error as a function of model complexity.*

- ▶ Remark: there is an emerging *theory* of a double-dip (W-shaped) generalized/test error curve, instead of the classic single-dip (U- or V-shaped) one.

DL;

Ref: Belkin et al. (2018).

<https://arxiv.org/abs/1812.11118>

- ▶ Q: If the test error rates are 0.1 and 0.2 for two methods/models, is the first one better?
- ▶ Generalization error: for a new/future data point  $(X^*, Y^*)$ ,  
 $GE(\hat{f}) = E[L(Y^*, \hat{f}(X^*))] = E[I(Y^* \neq \hat{f}(X^*))] = P(Y^* \neq \hat{f}(X^*)) = p.$
- ▶ Given a test dataset  $\{(X_i, Y_i) : i = 1, 2, \dots, n\}$ ,  
 $TE(\hat{f}) = \sum_{i=1}^n L(Y_i, \hat{f}(X_i))/n = \hat{p}.$   
 $\hat{p} \sim N(p, p(1-p)/n).$   
 $Var(\hat{p}) = \hat{p}(1-\hat{p})/n.$
- ▶ To compare  $\hat{p}_1$  and  $\hat{p}_2$ , need to consider their var.  
e.g. construct their 95% CIs
- ▶ BUT, ...  
use a paired-t-test, or McNemar's test.