Multivariate Analysis, I 1st part

Aldo Solari

INFN School of Statistics

Paestum, June 6, 2019



About me

Aldo Solari Associate Professor at the University of Milano-Bicocca PhD in Statistics at the University of Padua Homepage: aldosolari.github.io



Outline

What this lecture is not

A detailed math and statistics class, a programming course

What this lecture tries to be

An overview of statistical methods for prediction

Part 1 (1.5h)

- Statistical learning: supervised and unsupervised
- Model bias and variance and the risk of overfitting
- Cross-validation

Part 2 (1.5h)

- The model versus the modeling process
- Ensemble learning: bagging, random forest and boosting
- Regularized regression: ridge and lasso

If anything is unclear, just ask



Bibliography

- Azzalini, Scarpa (2012) *Data analysis and data mining, an introduction*. New York: Oxford University Press
- **Bishop** (2006) *Pattern Recognition and Machine Learning.* Springer
- Gareth, Witten, Hastie, Tibshirani (2014) An Introduction to Statistical Learning, with Applications in R. Springer
- Hastie, Tibshirani, Friedman (2009) The Elements of Statistical Learning. Data Mining, Inference and Prediction. Springer
- Hastie, Tibshirani, Wainwright (2015) *Statistical Learning with Sparsity. The Lasso and Generalizations.* Chapman and Hall/CRC
- Kuhn, Johnson (2013) *Applied Predictive Modelling*. Springer
- Kuhn, Johnson (2019) Feature Engineering and Selection: A Practical Approach for Predictive Models. Chapman and Hall/CRC



Evolution of multivariate statistics

Classic

• Multivariate Analysis

Books by Anderson (1958) and Mardia, Kent & Bibby (1979)

• Statistical Modeling

Nelder & Wedderburn (1972) paper on GLM

Computer-age

- Data Mining \approx process of discovering patterns in data
- Machine Learning \approx algorithms that can learn from data

Modern

• Statistical Learning

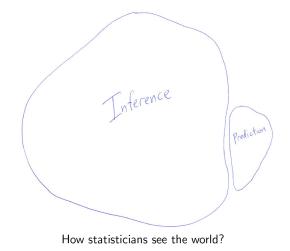
Hastie, Tibshirani & Friedman (2001) book

• Data Science



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

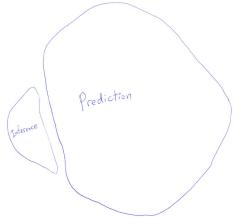
Statistics versus machine learning





From Robert Tibshirani

Statistics versus machine learning



How machine learners see the world?



From Robert Tibshirani

<ロ> <目> <目> <目> <目> <目> <日> <日> <日</p>

The two cultures

Breiman (2001)

Breiman distinguished between

Statistical modeling

- emphasis on probability models
- the goal is explanation

Machine learning

- emphasis on algorithms
- the goal is prediction accuracy

Today: Data Science vs. Statistics: Two Cultures? Carmichael and Marron (2018)



・ロト ・ 戸 ト ・ ヨ ト ・ ヨ ト - ヨ - -

To explain or to predict? Shmueli (2010)

True model

$$Y = \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

Underspecified model

$$Y = \gamma_1 X_1 + \epsilon$$

Explanation requires estimating the coefficients of the true model, but *a wrong model can sometimes predict better* :

- when the predictors X_1 and X_2 are highly correlated
- when the data are very noisy
- when β_2 is small



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Outline

1 Statistical learning

- **2** Model bias and variance
- **3** Cross-validation
- **4** Nonparametric methods



Statistical learning

Unsupervised learning

- the data consists of variables X_1, \ldots, X_p ; no variable has a special status
- the goal is
 - dimensionality reduction
 - clustering
 - etc.

Supervised learning

- the data consists of *response* Y and *predictors* X_1, \ldots, X_p
- the goal is **prediction** of the response:
 - Y continuous : regression problem
 - Y binary/multiclass: classification problem



< □ > < 同 > < E > < E > E

Data matrix

	variable X_1		variable X_j		variable X_p
observation 1	x ₁₁	•••	<i>x</i> _{1<i>j</i>}	•••	x _{1p}
observation 2	x ₂₁	•••	<i>x</i> _{2j}	•••	<i>x</i> ₂ <i>p</i>
		•••	•••	• • •	
observation <i>i</i>	x _{i1}	• • •	X _{ij}		X _{ip}
•••		•••	•••	•••	
observation <i>n</i>	X _{n1}		x _{nj}		X _{np}

- *n* = number of observations
- *p* = number of variables

Dimensionality reduction

$$\begin{array}{cc} X \mapsto Y & q$$

such that the transformed data Y retains most of the information

Principal Components Analysis (PCA)

Linear transformation

$$Y_q = \tilde{X} V_q$$

where

- \tilde{X} is the centered matrix
- the columns of V_q are the q eigenvectors with largest eigenvalues of the variance/covariance matrix $n^{-1}\tilde{X}^{\top}\tilde{X}$



Image compression with PCA



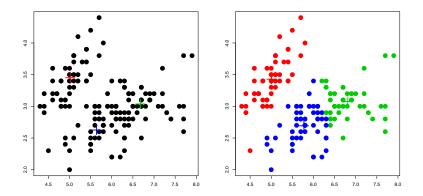




 $Y_q V_q^{\mathsf{T}} + 1\bar{x}^{\mathsf{T}} \quad q = 10$



Cluster analysis



K-means algorithm



E

1

Supervised learning

Machine Learning		Statistics
target variable	Y	response variable
attribute, feature	X	predictor
hypothesis	$Y = f(X) + \varepsilon$	model
instances, examples	$(y_1, x_1), \ldots, (y_n, x_n)$	observations
learning	$\hat{f} = \arg\min_{f \in \mathcal{F}} \sum_{i=1}^{''} \operatorname{loss}(y_i, f(x_i))$	estimation
classification generalization error	$\hat{y} = \hat{f}(x)$ $\mathbb{E}[\log(Y, \hat{f}(X))]$	prediction prediction error

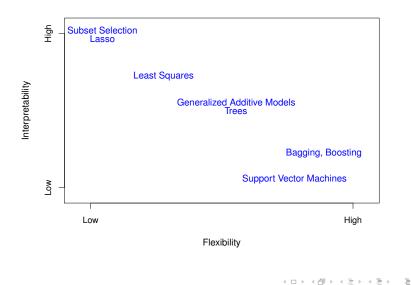


Models overview

~
シンシン
1
X
X
X
X



Flexibility versus interpretability





SUMMARY: statistical learning

- Learning: unsupervised / supervised
- Goal: to explain / to predict
- Prediction problem: regression / classification
- Models: interpretable / black-box



Important concepts

- Model bias and variance and the risk of overfitting
- Cross-validation
- Nonparametric methods



Outline

1 Statistical learning

2 Model bias and variance

3 Cross-validation

4 Nonparametric methods



Training and test

Training set

$$(x_1, y_1), \ldots, (x_i, y_i), \ldots, (x_n, y_n)$$

Test set

$$(x_1, y_1^*), \ldots, (x_i, y_i^*), \ldots, (x_n, y_n^*)$$

where x_1, \ldots, x_n are treated as fixed. In matrix notation:

$$\mathbf{y}_{n\times 1} = \begin{bmatrix} y_1 \\ \cdots \\ y_i \\ \cdots \\ y_n \end{bmatrix}, \quad \mathbf{X}_{n\times p} = \begin{bmatrix} x_1 \\ \cdots \\ x_i \\ \cdots \\ x_n \end{bmatrix} = \begin{bmatrix} x_{11} & \cdots & x_{1j} & \cdots & x_{1p} \\ \cdots & \cdots & \cdots & \cdots \\ x_{i1} & \cdots & x_{ij} & \cdots & x_{ip} \\ \cdots & \cdots & \cdots & \cdots \\ x_{n1} & \cdots & x_{nj} & \cdots & x_{np} \end{bmatrix}$$

Here y_1, \ldots, y_n are realizations of the random variables

$$Y_i = f(x_i) + \varepsilon_i, \quad i = 1, \dots, n$$

- *f* is the **regression function** (signal)
- ε is the **error** (noise) $\varepsilon_1, \ldots, \varepsilon_n$ i.i.d. with $\mathbb{E}(\varepsilon) = 0$ and $\mathbb{V}ar(\varepsilon) = \sigma^2$



Training set

0.550-0 ° ° ° ° ° ° ° ° ° 0 0 0.525 00 000 Ο 0 Ο 0 0.500 0 0 0 > 0 0.475 0.450 0 0.425 2 ż х

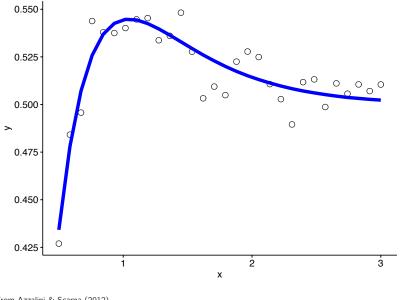
(日) (四) (三) (三) (三)

E

590

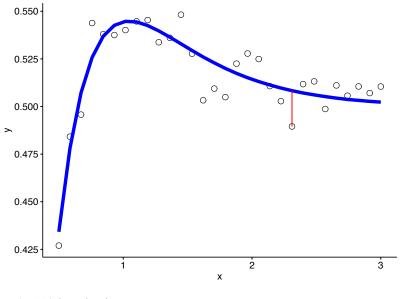
From Azzalini & Scarpa (2012)

Regression function



From Azzalini & Scarpa (2012)

(日) (四) (三) (三) (三) E 590 Error



From Azzalini & Scarpa (2012)

・ロト・「「「・」」・「」・ 「」・ (ロト

Mean squared error

- Loss function = mean squared error (MSE)
- Suppose we have estimated f by \hat{f} by using the training set
- We can compute the MSE on the training set

$$MSE_{Tr} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{f}(x_i))^2$$

• But the goal is to have a low MSE on the test set

$$MSE_{Te} = \frac{1}{n} \sum_{i=1}^{n} (y_i^* - \hat{f}(x_i))^2$$



Polynomial regression

• Polynomial regression model of degree *d*:

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \ldots + \beta_d x^d$$

• Use the training set to estimate *f* :

$$\hat{f}(x) = \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2 + \ldots + \hat{\beta}_d x^d$$

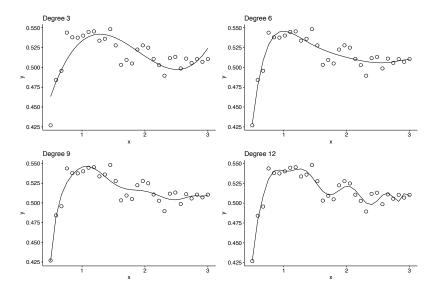
- The degree *d* is a **tuning parameter**: small and large *d* corresponds to low and high flexibility, respectively
- Goal: find the degree *d* that minimizes the prediction error

$$\operatorname{Err} = \mathbb{E}(\operatorname{MSE}_{\operatorname{Te}}) = \mathbb{E}\left\{\frac{1}{n}\sum_{i=1}^{n}[Y_{i}^{*} - \hat{f}(x_{i})]^{2}\right\}$$

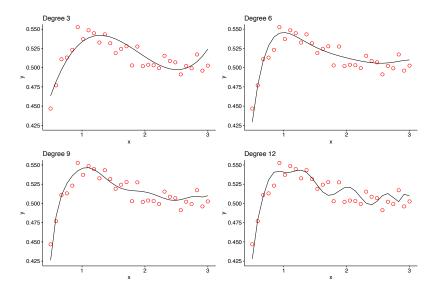


3

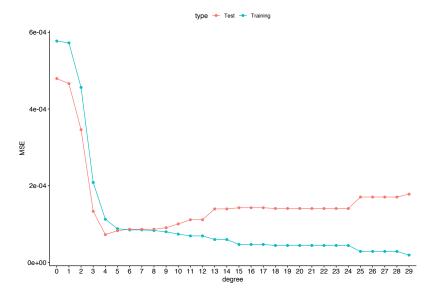
《日》《卽》《臣》《臣》







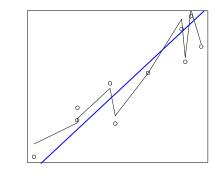
▲ ● ● ● ● ● ● ●





Overfitting

>



х

mistaking noise for signal

≣

Prediction error

Sources of error:

• Irreducible error

Can we ever predict Y from X with zero error? No. Even the true regression function f cannot do this

• Estimation bias

What happens if our fitted function \hat{f} belongs to a model class that is far from the true f? E.g. we choose to fit a regression line in a setting where the true relationship is far from linear?

• Estimation variance

What happens if our fitted function \hat{f} is itself quite variable? In other words, over different copies of the training data, we end up constructing substantially different \hat{f} ?



Reducible and irreducible error

The prediction error

$$\operatorname{Err} = \mathbb{E}(\operatorname{MSE}_{\operatorname{Te}}) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}\left\{ [Y_{i}^{*} - \hat{f}(x_{i})]^{2} \right\}$$

depends on two quantities, the **reducible error** and the **irreducible error**:

$$\begin{split} \mathbb{E}\{[Y_i^* - \hat{f}(x_i)]^2\} &= \mathbb{E}\{[f(x_i) - \hat{f}(x_i)]^2\} + \mathbb{V}\mathrm{ar}(\varepsilon_i^*) \\ &= \underbrace{\mathbb{E}\{[f(x_i) - \hat{f}(x_i)]^2\}}_{\text{Reducible}} + \underbrace{\sigma^2}_{\text{Irreducible}} \end{split}$$



Bias-variance decomposition

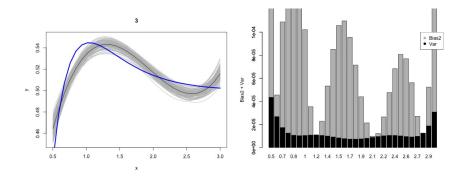
The reducible error can be decomposed into (squared) **bias** and **variance** of \hat{f}

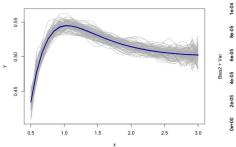
$$\mathbb{E}\{[f(x_i) - \hat{f}(x_i)]^2\} = \underbrace{[f(x_i) - \mathbb{E}\hat{f}(x_i)]^2}_{\text{Bias}^2} + \underbrace{\mathbb{Var}[\hat{f}(x_i)]}_{\text{Variance}}$$

Summing up, we have

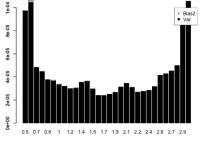
$$\operatorname{Err} = \frac{1}{n} \sum_{i=1}^{n} [f(x_i) - \mathbb{E}\hat{f}(x_i)]^2 + \frac{p\sigma^2}{n} + \sigma^2$$



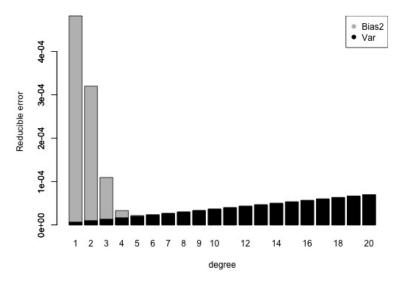




12



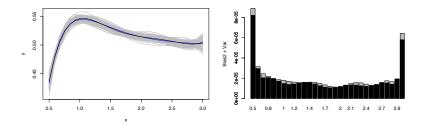
▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●





▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Best polynomial



Degree d = 5



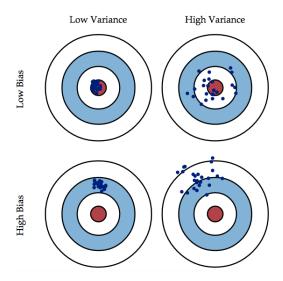
▲ロト ▲ □ ト ▲ 三 ト ▲ 三 三 - のへで

Bias-variance trade-off

Bias and variance are conflicting entities, and we cannot minimise both simultaneously:

- Models \hat{f} with low bias tend to have high variance
- Models \hat{f} with low variance tend to have high bias
- Even if our prediction is unbiased, i.e. $\mathbb{E}\hat{f}(x_i) = f(x_i)$, we can still incur a large error if it is highly variable
- On the other hand, $\hat{f}(x_i) = 0$ has 0 variance but will be terribly biased

To predict well, we must therefore choose a **trade-off** between bias and variance





SUMMARY: bias-variance tradeoff

- Data: training set / test set
- Signal/noise: regression function / error
- Performance: MSE on the same data / MSE on new data
- **Overfitting**: mistaking noise for signal
- **Prediction error**: reducible + irreducible
- **Reducible error**: bias² + variance
- Tradeoff: allow some bias if it decreases more variance



Outline

1 Statistical learning









Optimism

Optimism is the expected difference of the test error and training error

$$\mathrm{Opt} = \mathbb{E}(\mathrm{MSE}_{\mathrm{Te}}) - \mathbb{E}(\mathrm{MSE}_{\mathrm{Tr}})$$

Optimism is important because an estimate of optimism leads to an estimate of prediction error

$$\mathbb{E}(\widehat{\mathrm{MSE}}_{\mathrm{Te}}) = \mathrm{MSE}_{\mathrm{Tr}} + \widehat{\mathrm{OptF}}$$

For the linear model

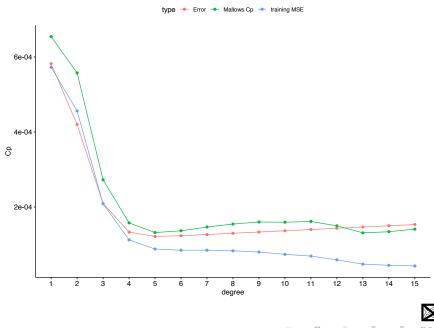
$$Opt = \frac{2\sigma^2 p}{n}$$

it gives Mallows' Cp

$$Cp = MSE_{Tr} + \frac{2\hat{\sigma}^2 p}{n}$$



<ロト < 四ト < 三ト < 三ト < 三ト = 三</p>



AIC and BIC

• AIC is given by

$$AIC = -2\ell(\hat{\beta}, \hat{\sigma}^2) + 2p$$

where $\boldsymbol{\ell}$ is the loglikelihood. For the linear model

$$-2\ell(\hat{\beta},\hat{\sigma}^2) = n\log(\mathrm{MSE}_{\mathrm{Tr}})$$

- For linear models, Cp and AIC are proportional to each other, and the lowest Cp corresponds to the lowest AIC
- BIC is given by

$$BIC = -2\ell(\hat{\beta}, \hat{\sigma}^2) + \log(n)p$$

• Since log(n) > 2 for any n > 7, the BIC statistic generally results in the selection of smaller models than AIC



◆□▶ ◆檀▶ ◆臣▶ ◆臣▶ ─ 臣・

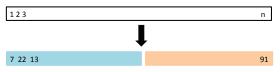
Cross-validation

- Training an algorithm and evaluating its statistical performance on the same data yields an overoptimistic result
- **Cross-validation** (CV) is based on the idea of **splitting the data**: part of data (the training set) is used for training the algorithm, and the remaining data (the validation set) are used for evaluating the performance of the algorithm
- The major interest of CV lies in the universality of the data splitting heuristics. Therefore, CV is a **non-parametric method** which can be applied to any algorithm in any framework. This universality is not shared by e.g. Cp, which is specific to linear regression



Validation set approach

• A simple approach is to randomly divide the *n* observations into two parts: a training set and a *validation* or *hold-out* set



- The model is fitted on the training set T ⊂ {1,..., n}, and the fitted model f^{-V} is used to predict the responses for the observations in the validation set V = {1,..., n} \ T
- This results in the estimate of the expected test error

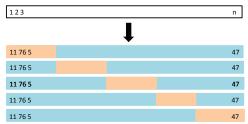
$$\widehat{\operatorname{Err}} = \frac{1}{\#V} \sum_{i \in V} [y_i - \hat{f}^{-V}(x_i)]^2$$

• This scheme reduces the sample size used for fitting the model, but this is not a problem when *n* is very large. If *n* is not very large, however, the validation estimate of the test error can be highly variable



K-fold CV

• Split the data into equal parts V_1, \ldots, V_K :



 Use observations i ∉ V_k for training the model and i ∈ V_k for evaluating it:

$$\frac{1}{\#V_k} \sum_{i \in V_k} [y_i - \hat{f}^{-V_k}(x_i)]^2$$

• Take the average to estimate the expected test error:

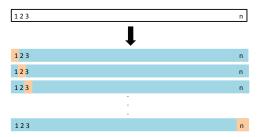
$$\widehat{\operatorname{Err}} = \frac{1}{K} \sum_{k=1}^{K} \left[\frac{1}{\#V_k} \sum_{i \in V_k} (y_i - \hat{f}^{-V_k}(x_i))^2 \right]$$



≣

Leave-one-out CV

Each data point is left out and used for validation:



For i = 1, ..., n:

- Hold out the *i*th training observation (x_i, y_i)
- Use n 1 observations for training the model f⁻ⁱ and the hold-out observation (x_i, y_i) to evaluate it by (y_i f⁻ⁱ(x_i))²

Take the average to estimate the expected test error:

$$\widehat{\operatorname{Err}} = \frac{1}{n} \sum_{i=1}^{n} [y_i - \hat{f}^{-i}(x_i)]^2$$



Generalized CV

For the linear model, there is a shortcut for computing LOOCV:

$$\frac{1}{n}\sum_{i=1}^{n}\left(y_{i}-\hat{f}^{-i}(x_{i})\right)^{2}=\frac{1}{n}\sum_{i=1}^{n}\left(\frac{y_{i}-\hat{f}(x_{i})}{1-h_{i}}\right)^{2}$$

where h_i is the *i*th diagonal element of the projection matrix

$$\mathbf{H}_{n \times n} = \mathbf{X} (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top}$$

In generalized cross-validation we compute

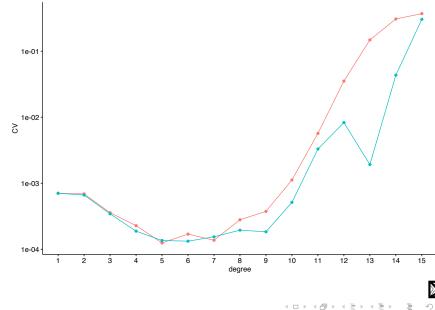
$$\widehat{\mathrm{Err}} = \frac{\mathrm{MSE}_{\mathrm{Tr}}}{\left(1 - \frac{p}{n}\right)^2}$$

where we approximate each h_i by their average $\frac{1}{n}\sum_{i=1}^n h_i = \frac{p}{n}$



 \exists





≣

CV bias-variance tradeoff

- A common choice for *K* other than *K* = *n* is to choose *K* = 5 or *K* = 10
- Bias: K-fold CV with K = 5 or 10 gives a biased (upward) estimate of 𝔅(MSE_{Te}) because it uses less information (4/5 or 9/10 of the observations). LOOCV has very low bias (it uses n − 1 observations)
- Variance: Usually, LOOCV has high variance because it is an average of *n* extremely correlated quantities (because the fits \hat{f}^{-i} and \hat{f}^{-j} are based on n-2 common observations), and *K*-fold CV with K = 5 or 10 has less variance because it is an average of quantities that are less correlated. Remember that the variance of the sum of highly correlated quantities is larger than that with midly correlated quantities:

 $\operatorname{Var}(A+B) = \operatorname{Var}(A) + \operatorname{Var}(B) + 2\operatorname{Cov}(A, B)$

 However, drawing a general conclusion on CV is nearly an impossible task because of the variety of frameworks



SUMMARY: optimism and cross-validation

- Optimism: test MSE training MSE
- Cp, AIC and BIC: training MSE + penality for complexity
- Cross-validation: model-free approach based on data split
- Types: validation set, K-fold, leave-one-out, generalized

Outline

1 Statistical learning

- **2** Model bias and variance
- **G** Cross-validation



4 Nonparametric methods



Nonparametric methods

- Nonparametric methods do not make explicit assumptions about the functional form of *f* (e.g. polynomials)
- Leave data to speak for themselves in a free way
- **Advantage**: by avoiding the assumption of a particular functional form for *f*, they have a wider range of possible shapes for *f*
- **Disadvantage**: since they do not reduce the problem of estimating *f* to a small number of parameters, a very large number of observations is required



k-nearest neighbors

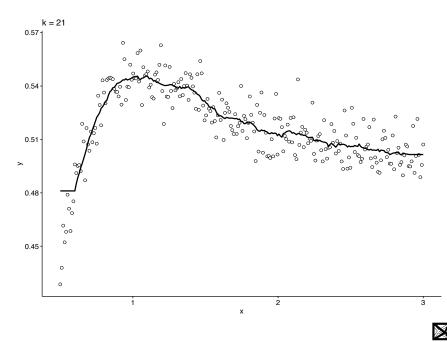
- A very simple and quite commonly used method k-nearest neighbors (kNN)
- Suppose we want to make a prediction at some x₁^{*}. Define the neighbourhood N_k(x₁^{*}) to be the set of k training observations having values x_i closest to x₁^{*} in Euclidean norm ||x_i - x₁^{*}||

$$\hat{f}(x_1^*) = \frac{1}{k} \sum_{i \in N_k(x_1^*)} y_i$$

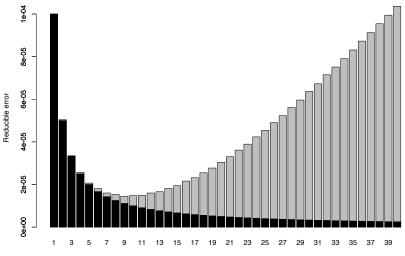
- The number k is a **tuning parameter**: small and large k corresponds to a more and less flexible fit, respectively
- Since we are computing a distance, usually we center and scale the predictors



(日)



コト 4 酉 ト 4 目 ト 4 目 - クタマ



k



kNN for classification

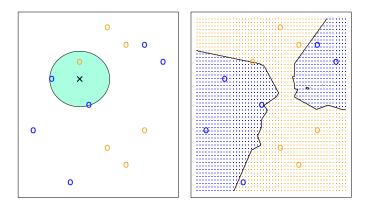


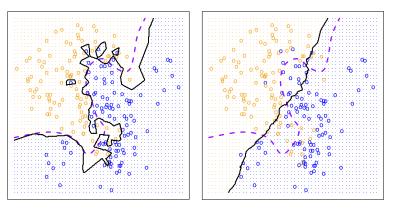
Figure from Gareth, Witten, Hastie & Tibshirani (2013)



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三 のへ⊙

KNN: K=1

KNN: K=100



Too flexible

Not flexible enough

・ロト ・ 同 ト ・ 三 ト ・

글▶ 글





SUMMARY: nonparametric methods

• Nonparametric: very flexible but more observations needed

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

- **kNN**: look at the neighbors
- **Tuning parameter**: *k* defines flexibility
- Neighbors distance: center and scale the predictors
- **Curse of dimensionality**: closeness difficult in high dimensions