

Genomics-microarrays Hastie et al:. DNA microarray data: expression matrix of 6830 genes (rows) and 64 samples (columns), for the human tumor data. Only a random sample of 100 rows are shown. The display is a heat map, ranging from bright green (negative, under expressed) to bright red (positive, over expressed). Missing values are gay. The rows and columns are displayed in a randomly chosen order.

732A99

Wide data

- Wide data $p \gg n$. Many variables, few data points.
 - GenomicsText
- Tall data: $p \ll n$. Few variables, many data points. Most of applications
 - · Economics, for ex. Currency exchange rates vs time
 - · Industry, Car performance characteristics vs probability of malfunctioning
 - Surveys, customer satisfaction vs survey answers
- Tall and Wide. Supermarket scanners. Many purchases, many products.

Text – document classification

Document	has('ball')	has('EU')	has('political_arena')	wordlen	Lex. Div.	Topic
Article1	Yes	No	Ne	4.1	5.4	Sports
Article2	No	No	No	6.5	13.4	Sports
:	:	:			:	- :
ArticleN	No	No	Yes	7.4	11.1	News

A problem with wide data

• Linear regression $\mu = w^T x, Y \sim N(\mu, \sigma^2)$

• ML solution $\widehat{w} = (X^T X)^{-1} X^T Y$

-X is $n \times p$, has rank n

 $-X^TX$ is $p \times p$, has rank n

 $- \rightarrow X^T X$ is not invertible!

· Solutions:

- Dimensionality reduction: PCA, PCR

- Shrinkage: Lasso, Ridge, Elastic network

- Forward variable selection

• Algorithms need sometimes be modified for wide data.

5

Classification: LDA

Standard LDA

$$\delta_k(x) = x^T \mathbf{\Sigma}^{-1} \mu_k - \frac{1}{2} \mu_k^T \mathbf{\Sigma}^{-1} \mu_k + \log \pi_k$$

$$\hat{\boldsymbol{\mu}}_c = \frac{1}{N_c} \sum_{i:y_i = c} \mathbf{x}_i, \quad \hat{\boldsymbol{\Sigma}}_c = \frac{1}{N_c} \sum_{i:y_i = c} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_c) (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_c)^T$$

$$\widehat{\Sigma} = \frac{1}{N} \sum_{c=1}^{k} N_c \, \widehat{\Sigma}_c$$

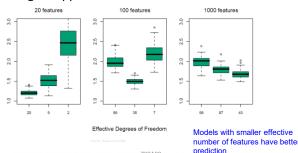
• $\rightarrow \Sigma^{-1}$ does not exist...

732A9

99

Effective amount of features for wide data

- Linear response generated with different p, n=100
- Ridge is applied with different λ



_

Classification: diagonal-covariance LDA

- Data is not enough to estimate dependences in covariance
- For wide data, we do diagonal-covariance LDA (naive Bayes): $\Sigma = diag(\sigma_1^2,...\sigma_p^2)$
- · Discriminant function

$$\delta(x^{new}) = -\sum_{j=1}^{p} \frac{\left(x_{j}^{new} - \bar{x}_{kj}\right)^{2}}{s_{j}^{2}} + 2\log \pi_{k}$$

$$-s_i^2 = \frac{1}{n} \sum_i n_i var(x_i | Y = C_i)$$

$$- \bar{x}_{kj} = mean(x_j | Y = C_k), \bar{x}_j = mean(x_j)$$

- Classify to the highest discriminant function value
- Drawback: all features are in the model → difficult to use in interpretations.

732A99

Classification: NSC

Nearest Shrunken Centroids

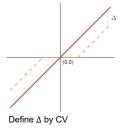
· Idea: Shrink classwise means towards overall mean

1. Compute
$$d_{kj} = \frac{\bar{x}_{kj} - \bar{x}_j}{m_k(s_j + s_0)}$$

2. Shrink
$$d'_{kj} = sign(d_{kj})(|d_{kj}| - \Delta)_+$$

3. Set
$$x'_{kj} = \bar{x}_j + m_k(s_j + s_0)d'_{kj}$$

Only features with nonzero d'_{kj} contribute to classification! \rightarrow insignificant features are shrunk!



32A99

132M33

9

NSC: example

- Package pamr
 - pamr.train()
 - pamr.cv

data0-read.csv2("voice.csv")
data0-data0
data0-as, data.frame(scale(data))
data3(pulity=as.factor(data0\$Quality)
library(pam')
rownames(data)=1:nrow(data)
x+t(data[,311])
y=data[[311]]
mydata1[ist(x=x,y=as.factor(y),geneid=as.character(1:nrow(x)), genenames=rownames(x))
model-pamr.train(mydata,threshold=seq(0,4,0.1))
pamr.plotcen(model, mydata, threshold=1)
pamr.plotcen(model, mydata, threshold=2.5)
cat(paste(colnames(data)[as.numeric(a[,1]]), collapse='\n'))
cvmodel-pamr.cv(model,mydata)
print(cwmodel)
pamr.plotcv(cwmodel)
732A99

NSC: example

 LSVT Voice Rehabilitation Data Set

- Target: Quality of voice rehabilitation
 - 1=acceptable, 2=not acceptable
- Features: Properties of the signal (voice)
- n=126, p=309



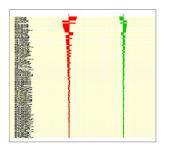
732A99

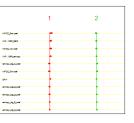
10

10

NSC: example

• Centroid plot, $\Delta = 1$ and $\Delta = 2.5$





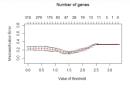
732A99

NSC: example

> pamr.listgenes(model.mydata,threshold=2.5)

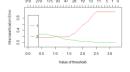
[1] id 1-score 2-score
[1] id 1-score 2-score
[2] immediate mec.2nd.cc
[3] immediate mec.2nd.cc
[4] immediate mec.2nd.cc
[5] immediate mec.2nd.cc
[6] immediate mec.2nd.cc
[6] immediate mec.2nd.cc
[7] immediate mec.2nd.cc
[8] immediate mec.2nd.cc
[9] immediate mec.2nd.cc
[9]

MFCC_2nd. coef IMF..NSR_SEO MFCC_1st.coef IMF..NSR_entropy entropy_log_2_coef entropy_log_5_coef entropy_log_6_coef entropy_log_3_coef



• Confusion matrix optimal Δ

	Pred 1	Pred 2
True 1	33	9
True 2	5	79



13

Regularized logistic regression

· Usual logistic regression

$$p(Y = C_i | x) = \frac{e^{w_{i0} + w_i^T x}}{\sum_{i=1}^K e^{w_{j0} + w_j^T x}} = softmax(w_{i0} + w_i^T x)$$

• Lp -Regularization:

$$\max_{w} \sum_{i=1}^{n} \log p(Y_{i}|x_{i}) - \frac{\lambda}{2} \sum_{k=1}^{K} ||w_{i}||^{p}$$

- Parameter redunancy is solved
- L1 regularization: some w are shrunk to 0
- · Numerical optimization is used to solve
- · R: LiblineaR() in package LiblineaR

732A99

RDA

Regularized discriminant analysis

• Another way of solving singularity of Σ

$$-\gamma$$
 is some constant

$$\hat{\Sigma}(\gamma) = \gamma \hat{\Sigma} + (1 - \gamma) diag(\hat{\Sigma})$$

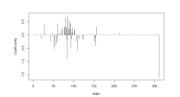
- $\gamma = 0 \rightarrow \text{diagonal-covariance LDA}$
- γ is chosen by CV
- R: rda() in klaR

14

L1 logistic regression

· Voice rehabilitation

W=mode12\$W plot(t(W), type="h", ylab="Coefficients")



	Pred 1	Pred 2
rue 1	41	1
rue 2	0	84

Overfitted?

SVM

- Support Vector Machine do not suffer from $p\gg n$ problem
 - Largest margin can be found even if the data is perfectly separable

A99 17

17

19

Elastic net

L1 regularization

$$\min_{w} -\log p(D|\mathbf{w}) + \lambda ||\mathbf{w}||_{1}$$
$$||\mathbf{w}||_{1} = \sum_{i} |w_{i}|$$

- For p>n, LASSO can extract at most n nonzero components
 - Severe regularization if $p \gg n$
- L1 regularization

 selects some feature among the correlated ones
- L2 regularization

 w's of the correlated variables are shrunk towards each other are nonzero

732A99 19

Computational shortcuts p>>n

- SVD decomposition $X = UDV^T = RV^T$
- If model is linear in parameters and has quadratic penalties:
 - Transform data observations from X into R
 - Minimize loss (minus log likelihood) with R instead of X and get $m{ heta}$
 - Original parameters $\boldsymbol{w} = V\boldsymbol{\theta}$
- Can be applied to many methods
- Example: ridge regression

/32A99 18

18

Elastic net

Combine L1 and L2 to diminish effect of L1 regularization.

• Elastic net regularization:

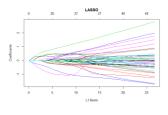
$$\min -\log p(D|\mathbf{w}) + \lambda(\alpha ||\mathbf{w}||_1 + (1-\alpha) ||\mathbf{w}||_2)$$

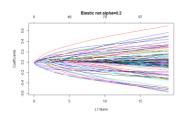
- α is set ad hoc or chosen by CV
- Elastic net may select more than n features
- R: glmnet() in glmnet package
 - Specify "family" for classification or regression

732A99 20

Elastic net

Voice rehabilitation



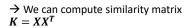


32A99 21

21

When features are not available

- Sometimes it is difficult to define or use the feature set
 - Molecule
 - Text document
 - possible, but can be very high dimensional
- ..but a proximity measure K(x, x') is easier to define
 - Ex: How much one document is different from another one





2A99 23

Comparative analysis

• Gene expression data

Methods	CV errors (SE) Out of 144	Test errors Out of 54	Number of Genes Used
1. Nearest shrunken centroids	35 (5.0)	17	6,520
2. L ₂ -penalized discriminant analysis	25 (4.1)	12	16,063
3. Support vector classifier	26 (4.2)	14	16,063
4. Lasso regression (one vs all)	30.7 (1.8)	12.5	1,429
 k-nearest neighbors 	41 (4.6)	26	16,063
6. L ₂ -penalized multinomial	26 (4.2)	15	16,063
7. L ₁ -penalized multinomial	17 (2.8)	13	269
8. Elastic-net penalized multinomial	22 (3.7)	11.8	384

732A99

22

When features are not available

- Many methods can use K instead of X
 - Note: p is not involved in calculations!!
- SVM: kernel trick → K can be used directly
- K-Nearest neighbors
 - Transform similarity into distance $d_{ij}^2 = K(x_i, x_i) + K(x_i, x_j) 2K(x_i, x_j)$
 - Use distances to find neighbors
- Can also be done for
 - Logistic and multinomial regression with L2 penalty
 - LDA
 - PCA: kernel PCA

732A99 24

Kernel PCA

- Usual PCA
 - Center X
 - Find $Su_i = \lambda_i u_i$, $S = \frac{1}{n} X^T X$, $S = [p \times p]$
 - ${m u}_i$ has dimension p
 - Project data on PCs: Z = X U
- Problems: X is unknown, and it can be p can be very large

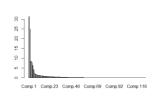
25

25

Kernel PCA in R

Use kpca() in kernlab

library(kernlab)
k <- as.kernelMatrix(crossprod(t(x)))
res=kpca(k)
barplot(res@eig)
plot(res@eid,1], res@rotated[,2], xlab="PC1",
ylab="PC2",</pre>



732A99 2

Kernel PCA

· Kernel PCA: Equivalent formulation

1. Solve $K'a_i = \lambda'_i a_i$, i = 1, ...M

-
$$K = ||K(x_i, x_i), i, j = 1, ...n||$$

- Centering
$$K' = K - \mathbf{1}_n K - K \mathbf{1}_n + \mathbf{1}_n K \mathbf{1}_n$$

$$-\lambda_i = \lambda_i'/n$$

2. Scores for PC_i : $z_i(x) = \sum_{i=1}^n a_{in}K(x, x_n)$

• There are at most *n* eigenvectors even if *p>>n*

/32A

26

Feature assessment

- Which features are important?
 - Ex: Which protein values differ between normal and cancer samples
- P-values in our predictive models can not be computed (too few observations)
- → Traditional hypothesis testing is used

732A99

28

27

Feature assessment

· Individual gene: t-test Hoi: treatment has no effect on gene j H_{1i} : treatment has an effect on gene j

$$t = \frac{\bar{x}_{2j} - \bar{x}_{1j}}{se_i}$$

- Alternatively, nonparametric tests (permutation tests) can be used to compare two populations
- Testing hypothesis for all genes?→multiple hypothesis testing
- · Control family-wise error rate
 - Bonferroni correction: $\alpha' = \alpha/M$
 - Ex: α=0.05, M=12000→α' ≈ 10⁻⁶

In practice, no genes with such small p-values

mean: \bar{x}_{2j}

mean: \bar{x}_{1j}

0

29

Feature assessment

- Alternative: false discovery rate (FDR)
 - Can not be exactly computed in practice

	Called nonsignif	Called signif	Total
H0 true	U	V	M0
H0 false	Т	S	M1
Total	M-R	R	M

$$FDR = E\left(\frac{V}{R}\right)$$

Feature assessment

Hypothesis testing Voice Rehabilitation

- Feature "MFCC 2nd.coef"

res=t.test(MFCC_2nd.coef~Quality,data=data, alternative="two.sided") res\$p.value

> res\$p.value [1] 1.21246e-11

res=oneway_test(MFCC_2nd.coef~as.factor(Qu ality), data=data,paired=FALSE) pvalue(res)

> pvalue(res) [1] 3.166942e-09

30

Feature assessment

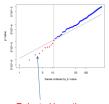
- Benjamini-Hochberg method (BH method)
 - Shown that $FDR(BH) < \alpha$ for independent hypotheses
 - − → we can control FDR!

Algorithm 18.2 Benjamini-Hochberg (BH) Method. 1. Fix the false discovery rate α and let $p_{(1)} \leq p_{(2)} \leq \cdots \leq p_{(M)}$ denote the ordered p-values

2. Define

 $L = \max \{ j : p_{(j)} < \alpha \cdot \frac{j}{M} \}.$

3. Reject all hypotheses H_{0j} for which $p_j \leq p_{(L)}$, the BH rejection threshold.

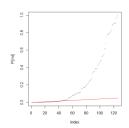


732499

Rejected hypotheses 33

Feature assessment

Voice rehabilitation



```
"Det judice feets, collapse-(in ) )
"Det judice feets, collapse-(in ) )
"Det judice, judice feets, collapse-(in ) )
"Det judice, judice feets, collapse-(in ) ]
"Det judice, judic
```

732A99

33