

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

سرطان روده بزرگ (Colon Cancer)

پژوهش درس داده کاوی

استاد درس: دکتر الهام عباسی

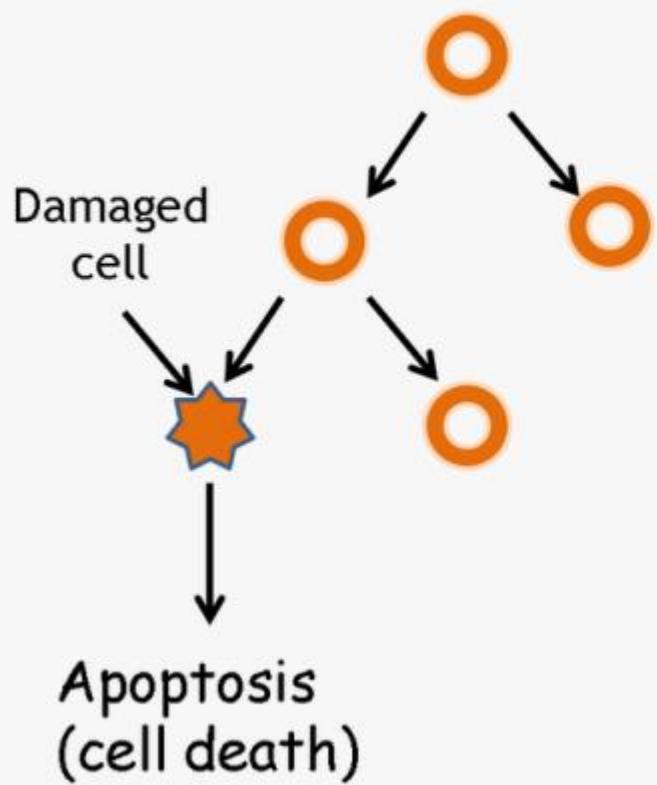
اسفند ۹۹

فهرست

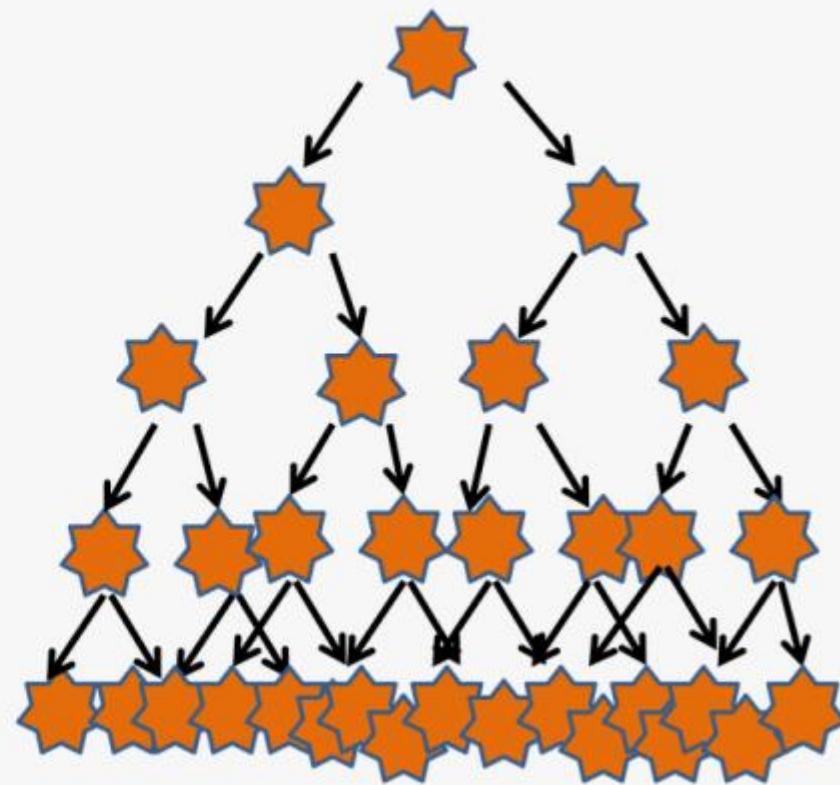


سرطان چیست؟

Normal Cell Division



Cancer Cell Division



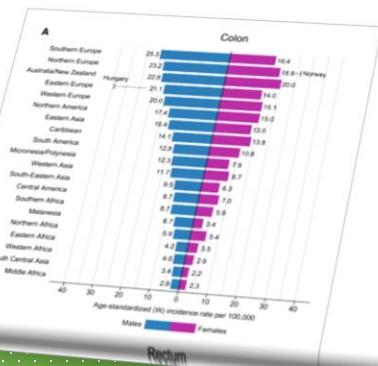
سرطان رو ده بزرگ

TABLE 1. New Cases and Deaths for 36 Cancers and All Cancers Combined in 2020

CANCER SITE	NO. OF NEW CASES (% OF ALL SITES)	NO. OF NEW DEATHS (% OF ALL SITES)
Female breast	2,261,419 (11.7)	684,996 (6.9)
Lung	2,206,771 (11.4)	1,796,144 (18.0)
Prostate	1,414,259 (7.3)	375,304 (3.8)
Nonmelanoma of skin ^a	1,198,073 (6.2)	63,731 (0.6)
Colon	1,148,515 (6.0)	576,858 (5.8)
Stomach	1,089,103 (5.6)	768,793 (7.7)
Liver	905,677 (4.7)	830,180 (8.3)
Rectum	732,210 (3.8)	339,022 (3.4)
Cervix uteri	604,127 (3.1)	341,831 (3.4)
Esophagus	604,100 (3.1)	544,076 (5.5)
Thyroid	586,202 (3.0)	43,646 (0.4)
Bladder	573,278 (3.0)	212,536 (2.1)
Non-Hodgkin lymphoma	544,352 (2.8)	259,793 (2.6)
Pancreas	495,773 (2.6)	466,003 (4.7)
Leukemia	474,519 (2.5)	311,594 (3.1)
All sites excluding nonmelanoma skin	18,094,716	9,894,402
All sites	19,292,789	9,958,133

Colorectal cancer can be considered a marker of socioeconomic development, and, in countries undergoing major transition, incidence rates tend to rise uniformly with increasing HDI.^{92,93} Incidence rates have been steadily rising in many countries in Eastern Europe, South Eastern and South Central Asia, and South America.^{22,94} The increase in formerly low-risk and lower HDI countries likely reflects changes in lifestyle factors and diet, ie, shifts toward an increased intake of animal-source foods and a more sedentary lifestyle, leading to decreased physical activity and increased prevalence of excess body weight, which are independently associated with colorectal cancer risk.⁹⁵ Additional risk factors include heavy alcohol consumption, cigarette smoking, and consumption of red or processed meat, whereas calcium supplements and adequate consumption of whole grains appear to decrease risk.⁹⁶ Primary prevention remains the cornerstone of colorectal cancer control.

More than 1.9 million new colorectal cancer (including anus) cases and 935,000 deaths were estimated to occur in 2020, representing about one in 10 cancer cases and deaths (Table 1). Overall, colorectal cancer ranks third in terms of incidence, but second in terms of mortality (Fig. 4). Incidence rates are approximately 4-fold higher in transitioning countries compared with transitioning countries, but there is less variation in the mortality rates because of higher fatality in transitioning countries (Fig. 7). There is an approximately 9-fold variation in colon cancer incidence rates by world regions, with the highest rates in European regions, Australia/New Zealand, and Northern America, with Hungary and Norway ranking first in men and women, respectively (Fig. 10A). Rectal cancer incidence rates have a similar regional distribution, although rates in Eastern Asia rank among the highest (Fig. 10B). Rates of both colon and rectal cancer incidence tend to be low in most regions of Africa and in South Central Asia.





ریاضیات سرطان

۸۲% (۲۹ رای)
لایک

ثبت‌نام رایگان

نمایش تیزر



دانش ما و سرطان

زبان برنامه نویسی R

پیش پردازش

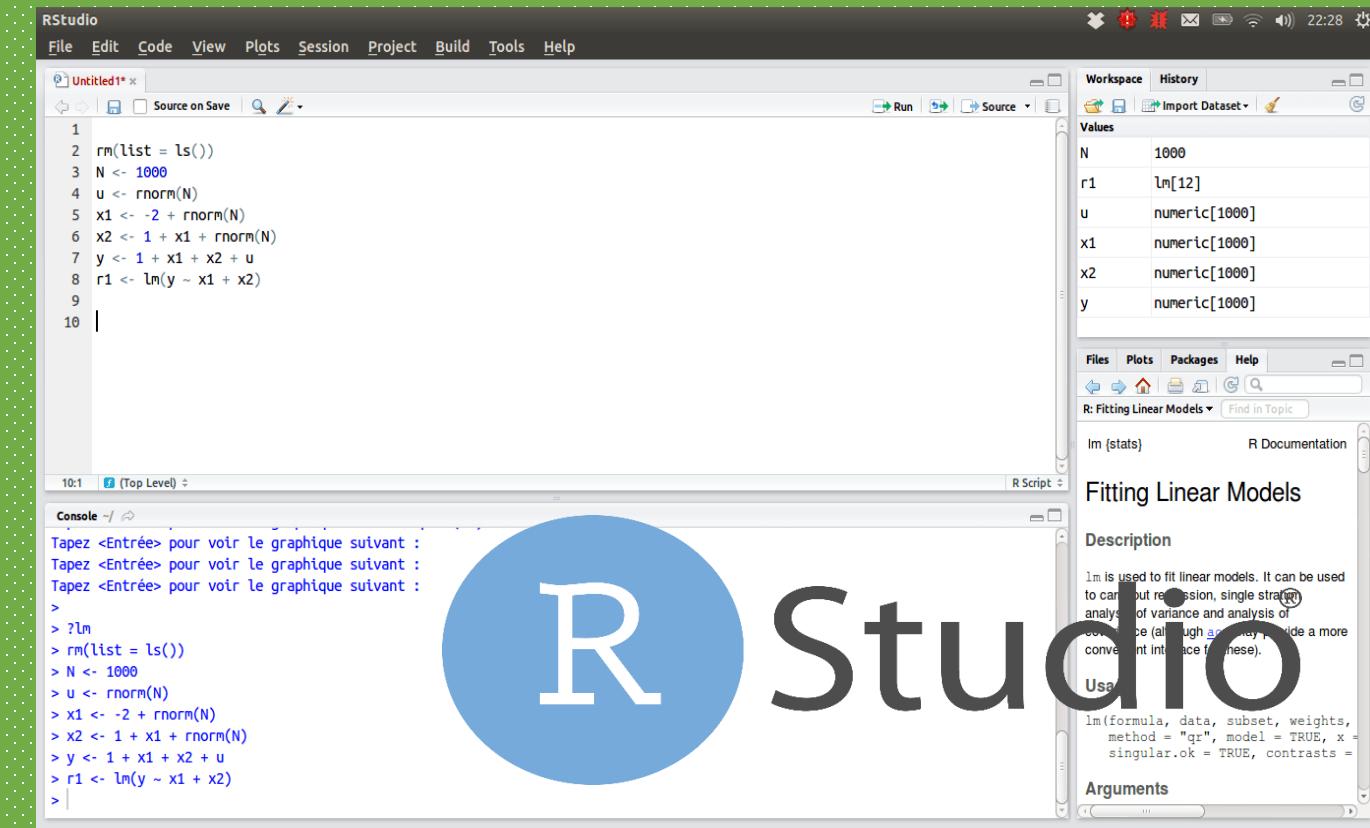
مدلسازی

ارزیابی

ارائه

خواندن داده

ابزار کار



Studio



خواندن داده

National Center for Biotechnology Information (NCBI)

NCBI Resources How To Sign in to NCBI

GEO DataSets GEO DataSets (Colon cancer healthy) AND "Homo sapiens"[porgn:_txid9606] Search Create alert Advanced Help

COVID-19 is an emerging, rapidly evolving situation.

[Public health information \(CDC\)](#) | [Research information \(NIH\)](#) | [SARS-CoV-2 data \(NCBI\)](#) | [Prevention and treatment information \(HHS\)](#)

Entry type Summary 20 per page Sort by Default order Send to: Filters: [Manage Filters](#)

Organism Top Organisms [Tree]
Homo sapiens (26)

Search results Items: 1 to 20 of 26 << First < Prev Page 1 of 2 Next > Last >>

Study type Expression profiling by array Filters activated: Expression profiling by array. [Clear all](#) to show 171 items.

Expression data from human non-affected and cancerous colon

1. (Submitter supplied) Innate lymphoid cells (ILCs) have the ability to sense and amplify inflammatory signals, and have been shown to exert bi-directional regulation with T helper cells in mice. However, how crosstalk between ILCs and CD4+ T cells influences immune function in humans is unknown. We observed that human intestinal ILCs co-localize with T cells in **healthy** and colorectal cancer tissue and display elevated HLA-DR expression in tumor and tumor-adjacent areas. [more...](#)

Organism: **Homo sapiens**
Type: **Expression profiling by array**
Platform: GPL23126 9 Samples
Download data: CEL, CHP
Series Accession: GSE145626 ID: 200145626
PubMed Full text in PMC Similar studies Analyze with GEO2R

Plasma Long Noncoding RNA and mRNA Expression Profile of Crohn's Disease identified by Microarray

2. (Submitter supplied) We performed a genome-wide analysis of lncRNA expression to identify novel targets in Crohn's disease (CD). Samples obtained from CD patients and control were analyzed using Arraystar human 8×60K lncRNA/mRNA v3.0 microarrays chips to find differentially expressed lncRNAs and mRNAs; The results were confirmed by quantitative reverse transcription-polymerase chain reaction (qRT-PCR). The differentially expressed lncRNAs and mRNAs were identified through fold-change filtering. [more...](#)

Organism: **Homo sapiens**
Type: **Expression profiling by array; Non-coding RNA profiling by array**

Find related data Database: Select Find items

Search details (("colonic neoplasms" [MeSH Terms] OR Colon cancer[All Fields]) AND healthy[All Fields] AND "Homo sapiens"[porgn] AND "Expression profiling by array"[Filter])

Recent activity Turn Off Clear

(Colon cancer healthy) AND "Homo sapiens" [porgn] (171) GEO DataSets

Colon cancer healthy (183) GEO DataSets

• اینزار کار

• خواندن داده

• پیش پردازش

• مدلسازی

• ارزیابی

• ارائه

- اینزار کار
- خواندن داده
- پیش پردازش
- مدلسازی
- آرژیمادی
- ارائه

National Taiwan University



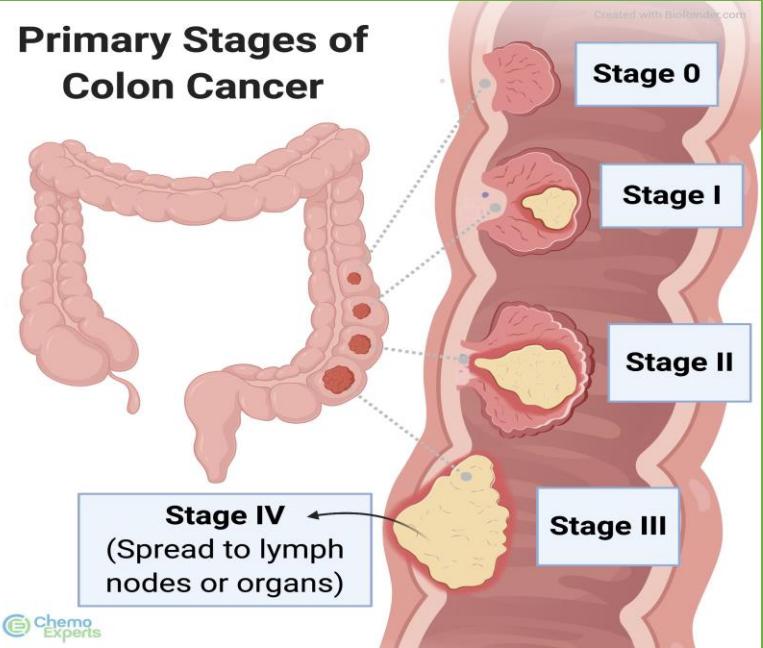
colon-cancer

- Source: [\[AU99a\]](#)
- Preprocessing: Instance-wise normalization to mean zero and variance one. Then feature-wise normalization to mean zero and variance one. [\[SKS03a\]](#)
- # of classes: 2
- # of data: 62
- # of features: 2,000
- Files:
 - [colon-cancer.bz2](#)

Cancer ~ 2000 Gene Expression

خواندن داده

Primary Stages of Colon Cancer



پیش پردازش

داده های نرمال و هم توزیع شده

خواندن داده

پیش پردازش

مدلسازی

ارزیابی

ارائه

ابزار کار

The screenshot shows a Notepad++ window displaying a large dataset. The file is named 'colon-cancer' and has a size of 1727461 bytes. The data consists of 63 columns and 1727461 rows. The columns represent various features of the colon cancer samples, such as age, sex, and tumor characteristics. The data is presented in a tabular format with each row representing a sample and each column representing a feature. The Notepad++ interface includes a menu bar, toolbars, and a status bar at the bottom.

```
length:1727461 lines:63 Ln:1 Col:1 Pos:1 Unix (LF) UTF-8 INS
```

پیش پردازش

اصلاح عددها و تبدیل cancer به فاکتور

- اینوار کار
- خواندن داده
- پیش پردازش
- مدلسازی
- آرژیمادی
- ارائه

The screenshot shows the RStudio interface with the following details:

- Project - RStudio**: The main window title.
- File Bar**: Contains options: File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help.
- Toolbar**: Includes icons for New Project, Open Project, Save, Print, Go to file/function, and Addins.
- File Tab Bar**: Shows multiple files: main.R, note 1.Rmd, Paper I.R, server.R, Paper II.R, ui.R, and Paper V.R.
- Code Editor**: Displays the R code for data preprocessing. The code reads a dataset named 'colon' from a file 'colon-cancer', converts its columns into doubles, and then converts the 'cancer' column into a factor.

```
1 ##Load Data
2 library(readr)
3 name <- c(paste0("Gene", 1:2000))
4 colon <- read_table2("colon-cancer", col_names = c("cancer", name))
5
6 ##Preprocessing
7 for (i in 2:2001) {
8   if(i<=10){
9     colon[i] <- as.double(substr(colon[[i]], 3, 20))
10  }else if(i<=100){
11    colon[i] <- as.double(substr(colon[[i]], 4, 20))
12  }else if(i<=1000){
13    colon[i] <- as.double(substr(colon[[i]], 5, 20))
14  }else{
15    colon[i] <- as.double(substr(colon[[i]], 6, 20))
16  }
17}
18 colon$cancer <- as.factor(colon$cancer)
```

- Status Bar**: Shows the current line (21:1), the top level, and the script type (R Script).

پیش پردازش

اصلاح عددها و تبدیل cancer به فاکتور

پیش پردازش

مدلسازی

ارزیابی

خواندن داده

ارائه

اینار کار

	cancer	Gene1	Gene2	Gene3	Gene4	Gene5	Gene6	Gene7	Gene8	Gene9	Gene10	Gene11	Gen
1	-1	2.080750	1.099070	0.927763	1.029080	-0.130763	1.265460	-0.436286	0.728881	2.107980	1.359870	0.265471	
2	1	1.109460	0.786453	0.445560	-0.146323	-0.996316	0.555759	0.290734	-0.145259	1.132660	0.559093	-1.469130	
3	-1	-0.676530	1.693100	1.559250	1.559980	-0.982179	-1.358510	-1.313990	-0.455067	0.295214	0.290694	0.415632	
4	1	0.534396	1.677540	1.489030	0.778605	-0.183776	-1.116850	-1.487560	-0.579511	0.292683	1.345480	-0.687898	
5	-1	-1.018900	0.511080	0.755641	1.013820	0.529899	0.160440	-0.087055	1.295290	0.458736	0.714082	0.727290	
6	1	-1.185370	-0.514473	-0.566634	1.224720	0.619244	-0.684713	-0.798129	1.368770	-0.697007	-1.006190	0.808748	
7	-1	1.779050	0.423947	0.820696	2.525690	0.666921	0.661346	0.425365	0.165247	1.967910	1.140080	-0.779230	
8	1	-0.889638	-0.315453	-0.073131	1.157500	-0.311039	-0.364472	-1.621640	1.193000	0.689805	0.203786	1.907460	
9	-1	-0.659694	-0.184388	-0.540022	1.122420	0.562609	-2.988310	-2.349810	-1.325010	-0.017002	0.109081	1.907460	1.005260
10	1	-1.225800	-0.212615	-0.588923	1.335410	-0.356505	0.354394	0.699607	0.190782	-0.139117	-0.551846	0.292320	
11	-1	-0.377282	-2.620490	-2.763910	0.612038	-0.155718	-1.456070	0.683292	-1.031570	0.718861	-2.313070	-0.092949	
12	1	1.639210	0.663551	0.324924	1.456460	-0.170587	-0.443486	-0.832590	-0.928334	1.320070	0.591865	0.158091	

Showing 1 to 13 of 62 entries, 2001 total columns

مدلسازی

• اینزارت کار

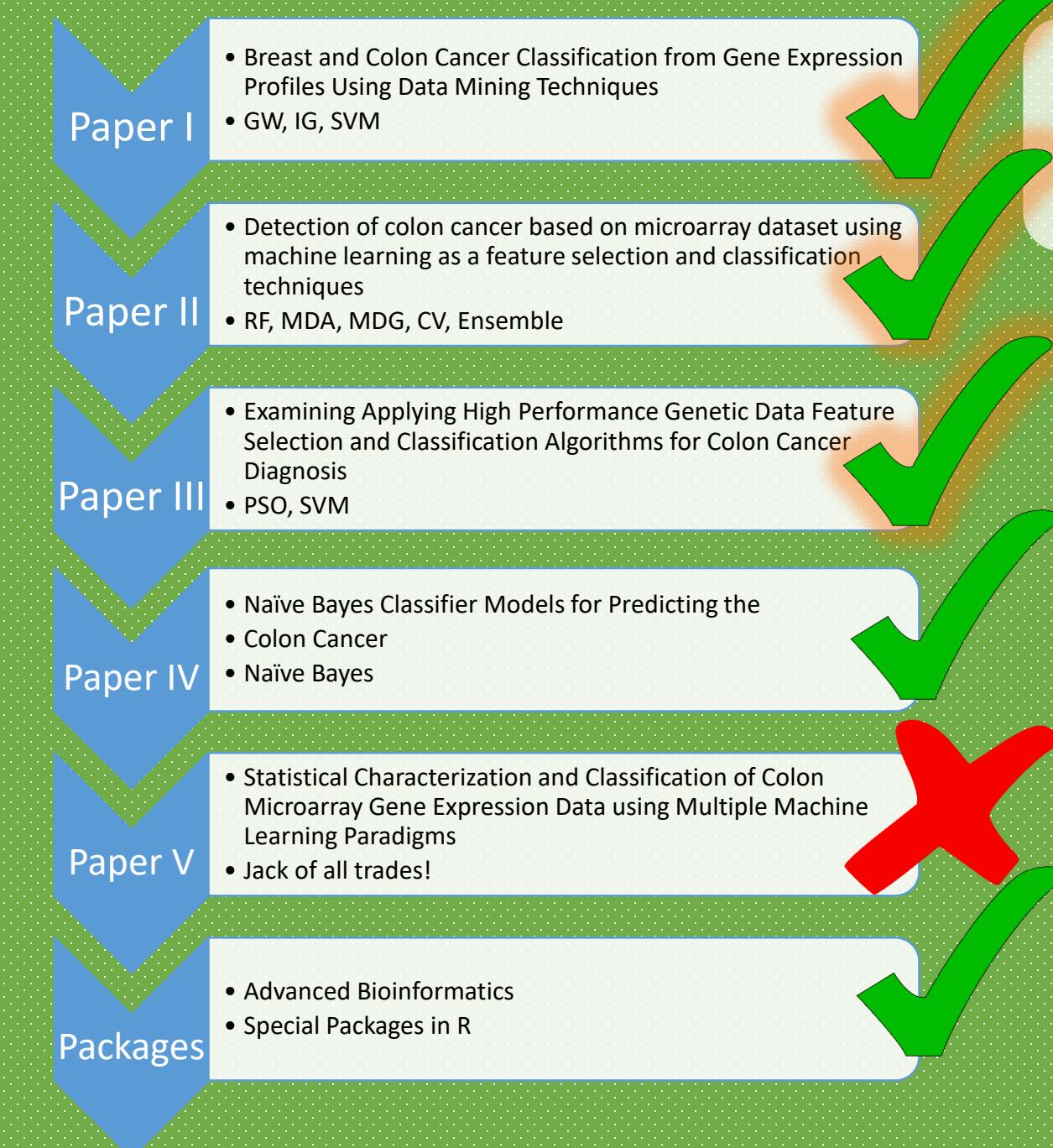
• خواندن داده

• پیش پردازش

• مدلسازی

• ارزیابی

• ارائه



مدلسازی

- اینوار کار
- Detection of colon cancer based on microarray dataset using machine learning as a feature selection and classification techniques
- RF, MDA, MDG, CV, Ensemble

• خواندن داده

• پیش پردازش

• مدلسازی

• ارزیابی

• ارائه

The screenshot shows the RStudio interface with several files open in the left pane: main.R, Paper.I.R, Paper.II.R, note1.Rmd, decisionTree.Rmd, and Paper.V.R. The main.R file contains R code for training a random forest model on a colon cancer dataset. The code includes loading the dataset, splitting it into training and testing sets, fitting a random forest model, and calculating various performance metrics like accuracy, sensitivity, and specificity. A large green checkmark is overlaid on the top right of the slide.

```
23  
24 #val <- train[i*(1:5)]  
25 #ktrain <- train[-i*(1:5)]  
26 ##train_test_split  
27 library(randomForest)  
28 set.seed(51)  
29 ind <- sample(1:62, size = 50)  
30 train <- colon[ind,]  
31 test <- colon[-ind,]  
32 #tree <- data.frame()  
33 krf <- randomForest(cancer~, data = train,  
34                         xtest = test[-1], ytest = test[[1]],  
35                         mtry=2000, importance =TRUE)  
36 #conf <- krf$confusion[1:2,1:2]  
37 #error <- 1-sum(diag(tabl))/50  
38 #tree <- rbind(tree,c(i,error))  
39  
40 tp <- krf$confusion[2,2]  
41 tn <- krf$confusion[1,1]  
42 fn <- krf$confusion[1,2]  
43 fp <- krf$confusion[2,1]  
44 acc <- (tp+tn)/(tp+tn+fp+fn)  
45 sens <- tp/(tp+fn)  
46 spec <- tn/(tn+fp)  
47 f.score <- 2*(acc*sens)/(acc+sens)  
48 cat(paste0(" | Real+ | Real- ", "\n",  
49             "predict+", "tp, " | ", fp, " ", "\n",  
50             "predict-", "fn, " | ", tn, " ", "\n",  
51             "*****", "\n",  
52             "Accuracy = ", acc, "\n",  
53             "Sensitivity = ", sens, "\n",  
54             "Specificity = ", spec, "\n",  
55             "f-score = ", f.score, "\n",  
56             "*****"))
```

The right pane shows the Global Environment, History, Connections, and Tutorial tabs. The Global Environment tab lists variables like prob_pc, sedgo1, sens, spec, spl, split, ssedgo1, st, and target. The Project tab shows a folder structure with files like colon-cancer, main.R, note 1.nb.html, note-1_files, Paper.I.R, Paper.V.R, texput.log, note1.Rmd, note1_files, note1.tex, note1.log, and note1.aux. The bottom status bar shows the current time as 43:25 and the file as (Top Level).

مدلسازی

فاعده بیز

- Naïve Bayes Classifier Models for Predicting the Colon Cancer
- Naïve Bayes



• خواندن داده

• پیش پنداش

• مدلسازی

• ارزیابی

• ارائه

$$\Pr(C|X) = \frac{\Pr(X|C)\Pr(C)}{\Pr(X)}$$

که با جایگذاری ویژگی ها به جای X به فرمول زیر می رسمیم:

$$\Pr(C|x_1, x_2, \dots, x_{2000}) = \frac{\Pr(x_1, x_2, \dots, x_{2000}|C)\Pr(C)}{\Pr(x_1, x_2, \dots, x_{2000})}$$

با فرض استقلال هر دو ویژگی از هم داریم:

$$\Pr(C|x_1, x_2, \dots, x_{2000}) = \Pr(C)\prod_{i=1}^{2000} \Pr(x_i|C)$$

تا اینجا همه چیز مطابق چیزی بود که در درس داشتیم اما قسمت متفاوت رفتار با داده های پیوسته است که پیشنهاد مقایله این است که برای آنها احتمال شرطی را به صورت زیر حساب می کنیم:

$$\Pr(X_i = x_i | C = c_j) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

حال می توانیم دست به کد شویم:

مدلسازی

- ناïve Bayes Classifier Models for Predicting the Colon Cancer
- ناïve Bayes



• خواندن داده

• پیش پردازش

• مدلسازی

• ارزیابی

• ارائه

```
cancer probability
cprob <- sum(train$cancer==1)/dim(train) [1]
hprob <- 1-cprob
ctrain <- subset(train,train$cancer==1)
htrain <- subset(train,train$cancer== -1)
hmean <- sapply(htrain[,2:2001], mean)
hsd <- sapply(htrain[,2:2001], sd)
cmean <- sapply(ctrain[,2:2001], mean)
csd <- sapply(ctrain[,2:2001], sd)
target = c()
for (i in 1:nrow(test)) {
p1 = cprob
p2 = hprob
for (j in 1:2000) {
  p1 = p1*dnorm(test[[i,j+1]],cmean[j],csd[j],log = F)
  p2 = p2*dnorm(test[[i,j+1]],hmean[j],hsd[j],log = F)
}
if(p1>p2){
  target = c(target,1)
}else{
  target = c(target,-1)
}
```

مدلسازی

- اینار کار
- ناïve Bayes Classifier Models for Predicting the Colon Cancer
- ناïve Bayes

• خواندن داده

• پیش پردازش

• مدلسازی

• ارزیابی

• ارائه



The screenshot shows an RStudio interface with several tabs open: main.R, Paper I.R, Paper II.R, note1.Rmd (which is the active tab), and Paper V.R. The code in note1.Rmd is as follows:

```
159 - if(pi>p2){  
160   target = c(target,1)  
161 }else{  
162   target = c(target,-1)  
163 }  
164 }  
165 ...  
166  
167 <h4>نایو بیزیس مدل</h4>  
168 ...  
169 ...{r echo=FALSE}  
170 tp <- sum(target==test$cancer & target==1)  
171 tn <- sum(target==test$cancer & target==-1)  
172 fp <- sum(test$cancer==1 & target==-1)  
173 fn <- sum(test$cancer==1 & target==1)  
174 acc <- (tp+tn)/(tp+tn+fp+fn)  
175 sens <- tp/(tp+fn)  
176 spec <- tn/(tn+fp)  
177 f.score <- 2*(acc*sens)/(acc+sens)  
178 cat(paste0(" | Real+ | Real- ", "\n",  
179 "predict+ | ", tp, " | ", fp, " | ", "\n",  
180 "predict- | ", fn, " | ", tn, " | ", "\n",  
181 "*****", "\n",  
182 "Accuracy = ", acc, "\n",  
183 "Sensitivity = ", sens, "\n",  
184 "Specificity = ", spec, "\n",  
185 "f-score = ", f.score, "\n",  
186 "*****"))|  
187 ...  
  
| Real+ | Real-  
predict+ | 0 | 0  
predict- | 2 | 10  
*****  
Accuracy = 0.8333333333333333  
Sensitivity = 0  
Specificity = 1  
f-score = 0  
*****  
  
188 <div class="header">  
189 <p>واعظ العالمين رب العالمين</p>  
190 </div>  
191  
6:58 R Markdown
```

مدلسازی

- Advanced Bioinformatics
- Special Packages in R



```
RStudio
Edit Code View Plots Session Build Debug Profile Tools Help
+ - + Go to file/function Addins
main.R x tT x note 1.Rmd x Paper I.R x Paper II.R x Paper V.R x
Source on Save | 
61 pcr <- data.frame(pc$x[,1:3], Group=colon$cancer)
62 pdf("results/PCA_samples.pdf")
63 ggplot(pcr, aes(PC1,PC2, color=Group)) + geom_point(size=3) + theme_bw()
64 dev.off()
65
66 library(limma)
67 design <- model.matrix(~cancer+0, colon)
68 colnames(design) <- c("neg", "pos")
69 fit <- lmFit(t(colon[,2:2001]), design)
70 cont.matrix <- makeContrasts(pos-neg, levels = design)
71 fit2 <- contrasts.fit(fit, cont.matrix)
72 fit2 <- eBayes(fit2, 0.01)
73 tT <- topTable(fit2, adjust="fdr", sort.by = "B", number = Inf)
74 write.table(tT, "results/Colon.txt", row.names = F, sep = "\t", quote = F)
75 #markGenes <- subset(tT, logFC>1 & adj.P.Val<0.05)
76 #aml <- sub("//.*","aml")
77 #library(caTools)
78 #set.seed(123)
79 #split <- sample.split(colon$cancer, splitRatio=0.8)
```

• اینار کار

• خواندن داده

• پیش پردازش

• مدلسازی

• ارزیابی

• ارائه

مدلسازی

- Advanced Bioinformatics
- Special Packages in R



The screenshot shows the RStudio interface with several files open in the session bar: main.R, tT, note 1.Rmd, Paper I.R, Paper II.R, and Paper V.R. The main window displays a data frame with 12 rows and 6 columns. The columns are labeled logFC, AveExpr, t, P.Value, adj.P.Val, and B. The data consists of gene identifiers and their statistical parameters. The first few rows are as follows:

	logFC	AveExpr	t	P.Value	adj.P.Val	B
Gene493	1.4970150	-3.387097e-07	5.758723	8.495853e-09	1.024875e-05	9.69269922
Gene377	1.4887588	-4.838710e-08	5.726963	1.024875e-08	1.024875e-05	9.51841067
Gene249	1.4307253	8.064516e-07	5.503720	3.726179e-08	2.484120e-05	8.32051824
Gene1635	1.4164825	1.935484e-07	5.448930	5.077294e-08	2.538647e-05	8.03380410
Gene1423	1.3426879	4.838710e-08	5.165057	2.407483e-07	9.208633e-05	6.59422304
Gene625	-1.3359808	-3.225806e-08	-5.139256	2.762590e-07	9.208633e-05	6.46719856
Gene245	1.2814361	-2.903226e-07	4.929434	8.257811e-07	2.198267e-04	5.45780927
Gene1771	-1.2782419	-9.677419e-08	-4.917146	8.793069e-07	2.198267e-04	5.40000290
Gene765	1.2614129	6.451613e-08	4.852408	1.221239e-06	2.368512e-04	5.09782248
Gene1772	-1.2609715	-2.903226e-07	-4.850710	1.231739e-06	2.368512e-04	5.08995071
Gene267	1.2580809	-3.548387e-07	4.839591	1.302682e-06	2.368512e-04	5.03846719

Showing 1 to 12 of 2,000 entries, 6 total columns

مدلسازی

- Advanced Bioinformatics
- Special Packages in R



• خواندن داده

• پیش پردازش

• مدلسازی

• ارزیابی

• ارائه

The screenshot shows the Enrichr web interface. In the top navigation bar, there are tabs for Analyze, What's new?, Libraries, Gene search, Term search, About, and Help. The main section is titled "Input data". It contains two input fields: one for uploading a BED file and another for pasting a list of Entrez gene symbols. Below these fields, a text box contains the gene symbols KRAS, NRAS, BRAF, and PIK3CA, separated by newlines. A message at the bottom of the text box says "0 gene(s) entered". There is also a field for entering a brief description of the list and a checkbox for contributing it to a public search index.

The screenshot shows the Enrichr web interface results page. At the top, there are tabs for Transcription, Pathways, Ontologies, Diseases/Drugs, Cell Types, Misc, Legacy, and Crowd. The "Description" tab is selected, showing the message "No description available (4 genes)". Below this, there is a grid of nine results cards. Each card has a title, a small icon, and a "Details" button. The titles include: ChEA 2016, ENCODE and ChEA Consensus TFs from CHIP-X, ARCHS4 TFs Coexp, TF Perturbations Followed by Expression, TRRUST Transcription Factors 2019, IncHUB lncRNA Co-Expression, Enrichr Submissions TF-Gene Cooccurrence, TRANSFAC and JASPAR PWMs, and Epigenomics Roadmap HM ChIP-seq.

مدلسازی

- Advanced Bioinformatics
- Special Packages in R

• خواندن داده

• پیش پردازش

• مدلسازی

• ارزیابی

• ارائه



The screenshot shows the RStudio interface with the following details:

- File Bar:** Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help.
- Toolbar:** Includes icons for file operations like Open, Save, and Print, along with Go to file/function and Addins dropdown.
- Project Explorer:** Shows files: main.R, test, tT, note 1.Rmd, Paper I.R, Paper II.R, Paper V.R.
- Code Editor:** Displays R code for logistic regression analysis, including variable assignment, model fitting, and performance metrics calculation.
- Console:** Shows the output of the R code, including accuracy, sensitivity, specificity, and f-score calculations.

```
85
86 ## Logistic Regression
87 ind <- sample(1:62, size = 50)
88 train <- rect[ind,]
89 test <- rect[-ind,]
90 classifier <- glm(formula = cancer ~ .,
91                      family = binomial,
92                      data = train)
93 prob_pred <- predict(classifier, type = 'response', newdata = test[-1])
94 y_pred <- ifelse(prob_pred>0.5, 1, -1)
95 tp <- sum(y_pred==test[[1]] & y_pred==1)
96 tn <- sum(y_pred==test[[1]] & y_pred==-1)
97 fn <- sum(y_pred!=test[[1]] & y_pred==1)
98 fp <- sum(y_pred!=test[[1]] & y_pred==1)
99 acc <- (tp+tn)/(tp+tn+fp+fn)
100 sens <- tp/(tp+fn)
101 spec <- tn/(tn+fp)
102 f.score <- 2*(acc*sens)/(acc+sens)
103 print(" | Real+ | Real- ")
104 print(paste0("predict+", tp, " | ", fp, " "))
105 print(paste0("predict-", fn, " | ", tn, " "))
106 print("*****")
107 print(paste0("Accuracy = ", acc))
108 print(paste0("Sensitivity = ", sens))
109 print(paste0("Specificity = ", spec))
110 print(paste0("f-score = ", f.score))
85:1 (Top Level) : 
```

```
C:/Users/shamsollah/Downloads/Data Mining/Project/
> print(paste0("Specificity = ", spec))
[1] "Specificity = 1"
> print(paste0("f-score = ", f.score))
[1] "f-score = 1"
> print("*****")
[1] "*****"
```

ارزیابی

• اینار کار

• خواندن داده

• پیش پردازش

• مدلسازی

• ارزیابی

• ارائه

Method	TP	FP	TN	FN	Accuracy	Sensitivity	Specificity	F-Score
RF	14	3	30	3	88	82.35	90.90	85
NB	0	0	10	2	83.33	0	100	0
SVM	8	0	4	0	100	100	100	100
Spec	5	0	7	0	100	100	100	100

ارائه

• اینار کار

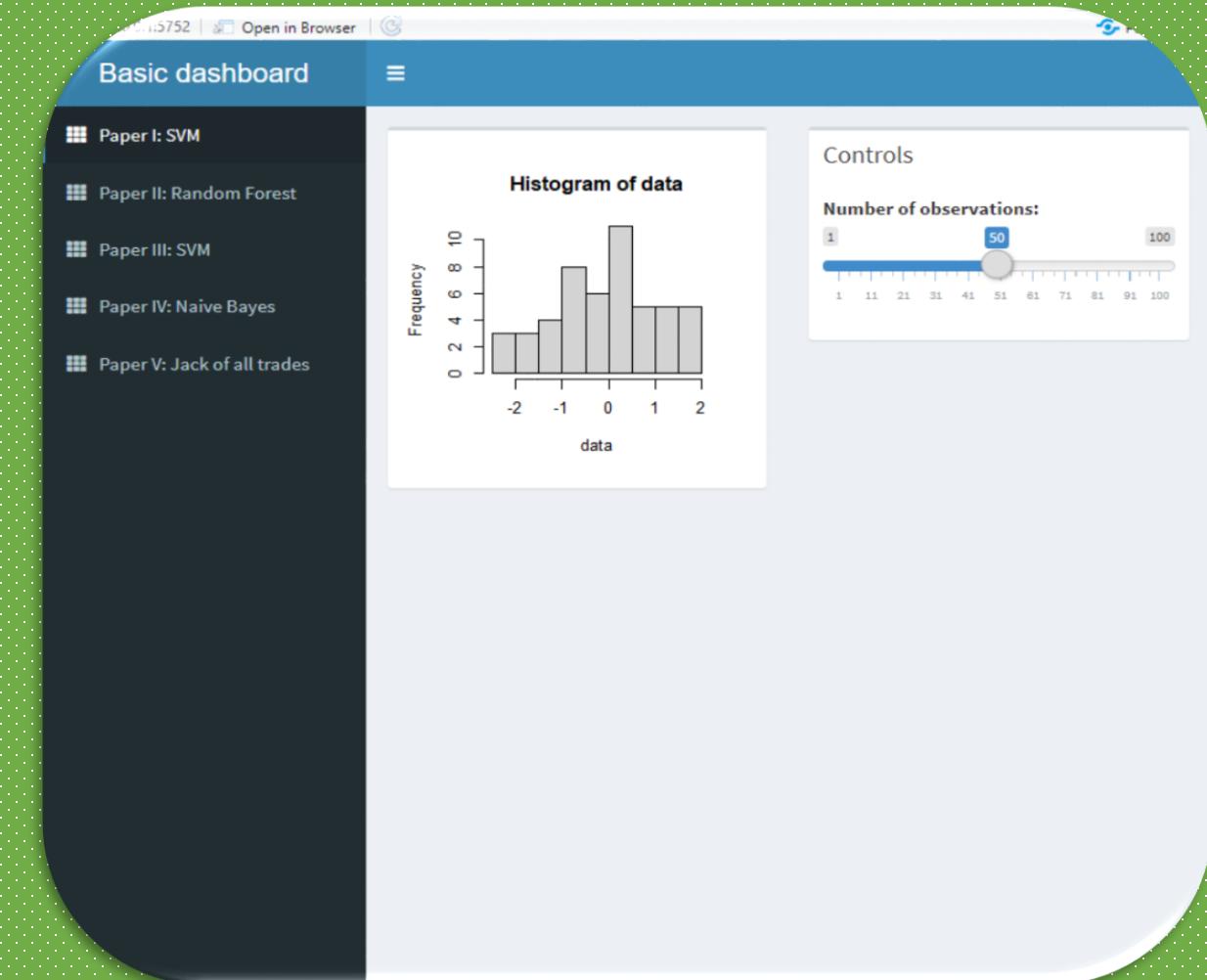
• خواندن داده

• پیش پردازش

• مدلسازی

• ارزیابی

• ارائه



مراجع

- <http://ce.sharif.edu/courses/97-98/2/ce550-1/>
- maktabkhooneh.org
- Foundations and Applications of Statistics An Introduction Using R [1,2,3]
- R for data science

با سپاس از توجه شما