Exercise3-数据科学与数据分析-20230314

1.安装所需要的包: RMySQL

```
# install.packages("RMySQL", repos = "https://mirrors.ustc.edu.cn/CRAN/")
```

2.载入包

```
any(grepl("RMySQL", installed.packages())) # 查看RMySQL是否安装成功
```

[1] TRUE

library (RMySQL)

3.连接MySQL数据库

查看链接信息

summary(mysqlconnection)

```
## <MySQLConnection:0,0>
## User: root
## Host: localhost
## Dbname: paper
## Connection type: localhost via TCP/IP
##
## Results:
```

查看该数据库内有哪些表dbListTables(mysqlconnection)

character(0)

4.使用R操作MySQL数据库的增删查改

```
# 在R环境中读取数据,为1000条论文文献数据
paper <- read.csv("paper.csv")

# 打印该数据的前6行
knitr::kable(head(paper))
```

id st title st abstract st year s10084 neutralization-1992 neutralization-sensitive merozoite surface antigens of babesia bovis encoded by members of a polymorphic sensitive gene family. Monospecific antibodies against native and recombinant versions of the major merozoite surface antigen (MSA-1) of Babesia bovis neutralize the infectivity of merozoites from Texas and Mexico merozoite surface antigens strains in vitro. Sequence analysis shows that MSA-1 and a related, co-expressed 44 kDa merozoite of babesia bovis surface protein (MSA-2) are encoded by members of a multigene family previously designated BabR. BabR encoded by genes, originally described in Australia strains of B. bovis, are notable because their marked polymorphism is apparently mediated by chromosomal rearrangements, but protein products of BabR genes have not members of a polymorphic previously been identified. The 3' terminal 173 nucleotides of the MSA-1 gene, including 60 nucleotides of untranslated sequence, are highly similar to the 3' terminal sequences of BabR 0.8 (84% identity) and MSAgene family. 2 (94% identity). Alignment of the predicted protein sequences demonstrates significant overall homology between MSA-1 and MSA-2, and between both proteins and the amino terminal BabR sequence. MSA-1 nucleic acid probes also hybridize weakly to genomic DNA from the Australia 'L' strain, even though this strain does not express merozoite surface epitopes cross-reactive with MSA-1 or MSA-2. Hybridization of these same probes to genomic DNA from the cloned Mexico strain reveals a pattern of bands compatible with two copies each of MSA-1 and MSA-2. Proteins encoded by this B. bovis gene family have been designated variable merozoite surface antigens (VMSA). The extent and mechanism of VMSA polymorphism among strains will be important when evaluating the role these surface proteins have in the host-parasite interaction, including immunity to blood stages. s10093 signalling 1992 signalling through the mhc class ii cytoplasmic domain is required for antigen presentation and induces b7 expression.. Class II major histocompatibility complex (MHC) molecules function as antigen-presenting through the mhc class ii elements as well as signal transducers on B lymphocytes. We previously reported that a B lymphoma cell cytoplasmic transfectant, 5C2, expressing genetically engineered I-Ak molecules with truncated cytoplasmic domains domain is was severely impaired in both antigen presentation and in anti-la-induced intracytoplasmic signalling. These two functions could be restored by preculturing 5C2 cells with cyclic AMP analogues. Here we demonstrate required for that impaired signal transduction by truncated class II molecules results in a deficiency in induction of the antigen presentation and newly defined B-cell accessory molecule B7 (ref. 8), which can be reversed by restoration of B7 expression. induces b7 These data imply that contact of the T-cell antigen receptor with MHC/antigen ligand results in signal transmission through the class II cytoplasmic domain. This signal, which can be mimicked by dibutyryl expression. cAMP, induces expression of B7, resulting in effective antigen presentation. The fact that crosslinking of surface class II MHC also induces B7 expression on normal resting human B cells supports this contention.

id	st_title	st_abstract	st_year
s10107	cloning and surface expression of pseudomonas aeruginosa o antigen in escherichia coli.	cloning and surface expression of pseudomonas aeruginosa o antigen in escherichia coli As a step toward developing recombinant oral vaccines, we have explored the feasibility of expression of O polysaccharide antigens from Pseudomonas aeruginosa by Escherichia coli. We cloned in E. coli HB101 a 26.2-kilobase DNA fragment from P. aeruginosa strain PA103 that specifies the production of the O polysaccharide of Fisher immunotype 2 (IT-2) strains. The recombinant organism incorporated the P. aeruginosa IT-2 O polysaccharide onto the core of the E. coli lipopolysaccharide (LPS). Transfer of the recombinant plasmid to three LPS-rough strains of P. aeruginosa resulted in synthesis of IT-2 O antigen, and two of these transconjugant strains also synthesized a second O polysaccharide, presumably representing expression of a repressed, or an incomplete, set of genes for an endogenous O polysaccharide. Rabbits injected with the purified recombinant LPS made antibody specific for P. aeruginosa IT-2 O side chains, as did mice fed the recombinant E. coli strain. Expression of P. aeruginosa O antigens by enteric bacteria makes it possible to study these recombinant strains as oral vaccines to prevent P. aeruginosa infections.	1992
s10198	murine b7 antigen provides a sufficient costimulatory signal for antigen-specific and mhc- restricted t cell activation.	murine b7 antigen provides a sufficient costimulatory signal for antigen-specific and mhc-restricted t cell activation We have previously shown that the murine B7 (mB7) molecule, when expressed in Chinese hamster ovary cells in stable fashion, can costimulate with anti-CD3 mAb or Con A to induce T cell activation. We have now derived, by gene transfection, Chinese hamster ovary cell lines that express the I-Ad molecule, either alone or in context with mB7. We have analyzed these transfectants for their capacity to present Ag to murine CD4+ T lymphocytes. I-Ad/mB7-double transfectants were able to stimulate mixed lymphocyte reactions and to present peptide Ag to specific T cells. Chinese hamster ovary cells that expressed only the I-Ad molecule were not able to stimulate T cell proliferation in these systems. Thus, the mB7 protein is a sufficient costimulatory molecule for the physiologic, Ag-dependent/MHC-restricted activation of murine CD4+ T cells. Stimulation of T cell bulk cultures resulted predominantly in the production of IL-2 and not of IL-4. The costimulatory activity of mB7 is not, however, restricted to the IL-2-secreting subset. We have identified one IL-4-secreting T cell clone, CDC35, which is responsive to mB7 triggering. Finally, we present experiments that suggest that mB7 and peptide/MHC complexes need to be expressed on the same cell for optimal induction of T cell activation.	1992

id st title st abstract st year 1992 s10289 the blood group the blood group antigen-related glycoepitopes: key structural determinants in immunogenesis and aids antigen-related pathogenesis.. This overview will focus on the functional and pathophysiological aspects of blood group glycoepitopes: antigen (BGA)-related glycodeterminants with regard to immunogenesis and AIDS pathogenesis. It has key structural been postulated that in a broad range of histogenetically different tissues and organs, BGA-related determinants in glycoepitopes are expressed on the cell surface at definite stages of cell differentiation. These glycoepitopes are expressed during embryogenesis, organogenesis, tissue repair, regeneration, immunogenesis remodelling and maturation when 'sorting-out' of one homotypic cell population from a heterotypic and aids pathogenesis. assemblage of cells occurs (1). In this event, the BGA-related glycoepitopes, if being expressed on the cell surface, play roles of key structural determinants in cell-cell recognition, association and aggregation. This mechanism will be discussed in relation to immunogenesis with regard to antigen presentation, self-non-self discrimination, and positive and negative selection during thymic education. It is postulated that the appearance of BGA-related glycoepitopes on the cell membrane is a consequence of the association of major histocompatibility complex antigens (MHC) and peptides, with the subsequent elimination of cells carrying a high density of BGA-related glycoepitopes on their surface. After human immunodeficiency virus (HIV) glycoproteins are glycosylated by host cell glycosyltransferases, the virus may use the BGA-related glycodeterminants as ligands and/or receptors for expansion to a spectrum of target cells during AIDS development and generalization of the infection throughout the body. We will review the experimental evidence that supports the concept that HIV uses an alternative to the gp120/CD4 ligand/receptor system, and that the alternative mechanism is probably carbohydrate-mediated in nature. 1992 s10545 proteasome proteasome subunits encoded by the major histocompatibility complex are not essential for antigen subunits presentation.. Major histocompatibility complex (MHC) class I molecules bind and deliver peptides derived encoded by the from endogenously synthesized proteins to the cell surface for survey by cytotoxic T lymphocytes. It is major believed that endogenous antigens are generally degraded in the cytosol, the resulting peptides being translocated into the endoplasmic reticulum where they bind to MHC class I molecules. Transporters histocompatibility containing an ATP-binding cassette encoded by the MHC class II region seem to be responsible for this complex are not essential for transport. Genes coding for two subunits of the '20S' proteasome (a multicatalytic proteinase) have been found in the vicinity of the two transporter genes in the MHC class II region, indicating that the proteasome antigen presentation. could be the unknown proteolytic entity in the cytosol involved in the generation of MHC class I-binding peptides. By introducing rat genes encoding the MHC-linked transporters into a human cell line lacking both transporter and proteasome subunit genes, we show here that the MHC-encoded proteasome subunit are not essential for stable MHC class I surface expression, or for processing and presentation of antigenic peptides from influenza virus and an intracellular protein.

- # 操作一: 在数据库中创建表
 # 利用R的数据框构建数据库的结构
 dbCreateTable(mysqlconnection, "paper", paper) # 新建一张名为paper的表
 dbListTables(mysqlconnection)
- ## [1] "paper"
- #操作二:向数据表中插入数据 dbWriteTable(mysqlconnection, "paper", paper, row. names=FALSE, append=TRUE)
- ## [1] TRUE
- # 注意:
- # 如果报错: could not run statement: Loading local data is disabled; this must be enabled on both the client and server sides
- # 需要调整MySQL全局参数,在MySQL中运行sql:set global local_infile=true;
- #操作三:查询,读取paper内的全部数据
 paper_from_mysql = dbGetQuery(mysqlconnection, "select * from paper") # 在双引号内写SQL查询式
 print(dim(paper_from_mysql))
- ## [1] 1000 4
- # 注意:
- # 如果报错: connection with pending rows, close resultSet before continuing
- # 执行dbClearResult(dbListResults(连接名)[[1]]),清理查询结果,重新跑
- # 操作四: 删除表dbRemoveTable(mysqlconnection,'paper')
- ## [1] TRUE

dbListTables(mysqlconnection)

character(0)

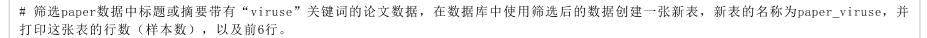
5.关闭连接

dbDisconnect(mysqlconnection)

[1] TRUE

6.练习

[1] TRUE



dbRemoveTable(mysqlconnection, 'paper')

[1] TRUE