Zhang-Lab 生信小课堂 第廿三期 Applied Bioinformatics Club (ABC)

和趣求真区兼实生信

(张建伟生物信息学课题组 https://zhang. hzau. edu. cn)

群体遗传学基础分析: π值与Tajima's D的计算与解读

2025.10.10 二综C102 15:00 欢迎大家交流学习!

主讲人: 丁悦 2025/10/10



目录

- 01 核心概念与理论基础
- 02 π值计算与解读
- 03 Tajima's D值计算与解读
- 04 综合评估



核心概念与理论基础

- 遗传变异:指一个物种群体内,不同个体在基因组特定位点上的DNA序列差异,是群体遗传分析的基本单元。最常见的类型是单核苷酸多态性。
- 中性理论:是分子进化的标准零假设。该理论认为,在分子水平上,绝大多数观察到的遗传变异是中性或近中性的,其命运主要由随机遗传漂变决定,而非自然选择。
- 自然选择:基于表型适应度的非随机过程,会定向改变等位基因频率,在基因组上留下特征性的统计模式,即选择信号。
- 选择信号:基因组中因自然选择作用而导致其遗传变异模式(如多样性水平、等位基因频率分布等)显著偏离中性理论 预期的区域。

群体遗传学通过量化种群内的遗传变异模式,来推断塑造这些模式的进化动力。

π **值计算与解读**-π值含义

PI值: (π值,核苷酸多样性)反映群体中任意两个随机选择的序列间平均核苷酸差异。

计算公式:

$$\pi = rac{\sum_{i < j} d_{ij}}{\binom{n}{2} imes L}$$

•dij: 第i条和第j条序列之间的核苷酸差异数。

•n: 群体中的序列数(样本量)。

•L: 分析的序列长度(排除缺失或不可比位点)。

•(n/2): 序列对的总数, 即n(n-1)/2

π值计算与解读-π值计算

PI值计算:

输入数据: 533 fixed.vcf(VCF文件,包含所有样本的基因型信息)

工具: `vcftools`

代码示例:

```
bsub -J vcftools -n 10 -R span[hosts=1] -o %J.out -e %J.err -q normal "vcftools --vcf 533_fixed.vcf --
window-pi 100000 --window-pi-step 100 --out 533
```

分亚群计算PI值

输入数据: cA.txt cB.txt GJ-tmp.txt GJ-trop.txt 533_fixed.vcf

```
#!/bin

sample_files=(../subgroup2/*.txt)

for sample in "${sample_files[@]}"; do

new_name1=$(echo "$sample" | awk -F '.' '{print $3}') # 添加前缀并保持扩展名不变

new_name2=$(echo "$new_name1" | awk -F '/' '{print $3}')

output="./subgroup2/pi_window_${new_name2}"

vcftools --vcf 533_fixed.vcf --window-pi 100000 --window-pi-step 100 --keep "$sample" --out "$output"

done
```

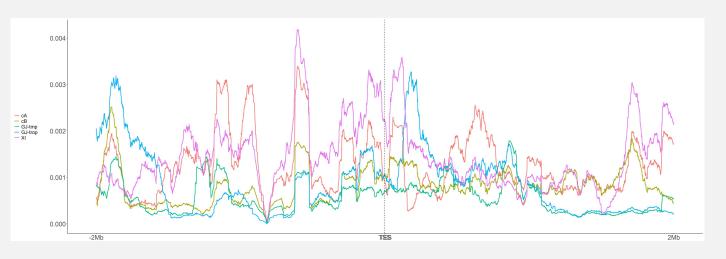
参数控制: Windows设置为100kb; 步长为100

比对结果: pi_window_*.windowed.pi

```
CHROM BIN_START BIN_END N_VARIANTS PI
1 1 100000 917 0.00255472
3 1 101 100100 918 0.00255666
4 1 201 100200 918 0.00255666
5 1 301 100300 918 0.00255666
6 1 401 100400 918 0.00255666
7 1 501 100500 918 0.00255666
8 1 601 100600 918 0.00255666
9 1 701 100700 918 0.00255666
10 1 801 100800 918 0.00255666
11 1 901 100900 918 0.00255666
```

π值计算与解读-可视化展示与结果分析

Rscript plot-gene-2Mb.r



*基因上下2Mb受选择情况

首先, 关注整个群体或染色体水平的平均π值

- ✓ 高 π 值 (例如 > 0.01): 表明群体**遗传多样性高**
- ✓ 低π值(例如 < 0.001): 表明群体遗传多样性低</p>

```
windowsFonts(A=windowsFont("Times New Roman"), B=windowsFont("Arial"))
4   # tajima_d_data <- read.table("rice_core_TajimaD.Tajima.D", header = TRUE, sep = "\t")
 # pi_data <- read.table("rice_core_pi.windowed.pi", header = TRUE, sep = "\t")
 # fst_data<-read.table("GJ_XI_fst.txt", header=TRUE, sep = "\t")</pre>
  Total data<-read.table("combined PI.txt", header=TRUE, sep = "\t")
  Total_data <- Total_data[Total_data$CHROM=='2']
 #Total data$BIN START <- Total data$BIN START / 100000</pre>
 Total data$BIN START <- Total data$BIN START / 100
  Total data <- Total data[Total data$BIN START>=323065.14 & Total data$BIN START <= 363133.19, ]
  ggplot(Total_data, aes(x = BIN_START, y = Value, color = TYPE)) +
    # geom rect(aes(xmin=100, xmax=300, ymin=0, ymax=60,col = "lightgrey",fill = "lightgrey")) + #
    # theme ridges(grid = FALSE)+
    xlab("LOC_0s01g59350") +
    scale\_color\_manual(values = c("#ff5722","#f8b500","#17b978","#07689f","#A079BF")) +
    #SCALE_COLOr_manual(Values = C( #FD/63F , #23BAC5 )) +
    # facet_grid(TYPE~.,scales = "free_y") +
    theme(legend.position = "left",
        axis.line = element line(color = "black", linewidth = 0.5),
          legend.title=element_blank(), ## 图例位置(可选:bottom/right/left/top)
          axis.text = element_text(size = 20),
          strip.text = element_text(size = 20),
          # axis.ticks.x = element_blank(),
          legend.text=element text(size = 15),
          strip.text.y = element_text(size = 12, face = "bold"),)+
          # panel.border = element_blank(),
    scale_x_continuous(breaks = c(323065.14,343065.14,343133.19,363133.19),
    #scale_x_continuous(breaks = c(34.306514), # 18414297 18415361
                       #labels = c("LOC_Os02g30850")) +
    # ylim(0,1)
    geom_vline(xintercept = 343065.14, col = "black", lty = 2)
    geom vline(xintercept = 343133.19, col = "black", lty = 2)
  ggsave("Chr2 2.5.png", width = 30, height = 10)
```

Tajima's D值计算与解读-Tajima's D值含义

Tajima's D: 一种用于检测群体遗传数据是否符合中性进化假说的统计量,通过比较两种基于突变率的遗传多样性估计值(π 值和Watterson's θ)的差异,揭示自然选择或群体历史事件(如扩张、瓶颈)的信号。

计算公式: Tajima D= θ π – θ s

Θπ:群体序列两两比较差异位点数累加/总两两比较对数,对中等频率变异敏感

θs:是群体总变异位点数/序列数的倒数累加,基于分离位点数,对低频变异敏感

Tajima's D值计算与解读-Tajima's D值计算

输入数据: 533_fixed.vcf(VCF文件,包含所有样本的基因型信息)

工具: `vcftools`

命令:

```
1 bsub -J vcftools -n 10 -R span[hosts=1] -o %J.out -e %J.err -q normal "vcftools --vcf 533_fixed.vcf --
TajimaD 100000 --out TajimaD_output"
```

分亚群计算Tajima's D值

输入数据: cA.txt cB.txt GJ-tmp.txt GJ-trop.txt 533_fixed.vcf

```
#!/bin

sample_files=(../subgroup2/*.txt)

for sample in "${sample_files[@]}"; do

new_name1=$(echo "$sample" | awk -F '.' '{print $3}') # 添加前缀并保持扩展名不变

new_name2=$(echo "$new_name1" | awk -F '/' '{print $3}')

output="./subgroup2/TajimaD_${new_name2}"

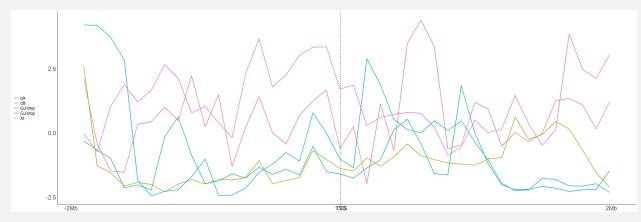
vcftools --vcf 533_fixed.vcf --TajimaD 100000 --keep "$sample" --out "$output" done
```

参数控制: Windows设置为100kb; 步长为100

比对结果: TajimaD_*.D

Tajima's D值计算与解读-可视化展示与结果解读

Rscript plot-gene-2Mb.r



基因上下2Mb受选择情况

计算值分析:

D=0:数据符合中性进化模型(π 与 θ W无显著差异)。

D>0: 可能表明平衡选择或群体收缩(低频等位基因较少)。

D<0: 可能提示定向选择(正选择)或群体扩张(近期种群增

长)。

```
Library(ggplot2)
  windowsFonts(A=windowsFont("Times New Roman"), B=windowsFont("Arial"))
4 # tajima_d_data <- read.table("rice_core_TajimaD.Tajima.D", header = TRUE, sep = "\t")
5  # pi data <- read.table("rice_core_pi.windowed.pi", header = TRUE, sep = "\t")
  # fst_data<-read.table("GJ_XI_fst.txt", header=TRUE, sep = "\t")</pre>
  Total data<-read.table("combined TajimaD.txt", header=TRUE, sep = "\t")
  Total_data <- Total_data[Total_data$CHROM=='2', ]</pre>
   Total data$BIN START <- Total data$BIN START / 100
  Total data <- Total data[Total data$BIN START>=164142.97 & Total data$BIN START <= 204153.61,
  ggplot(Total_data, aes(x = BIN_START, y = Value, color = TYPE)) +
    # geom_rect(aes(xmin=100, xmax=300, ymin=0, ymax=60,col = "lightgrey",fill = "lightgrey")) + #
     # theme ridges(grid = FALSE)+
     #xlab("LOC_Os01g59350") +
    xlab("LOC_Os02g30850") +
     scale color manual(values = c("#ff5722","#f8b500","#17b978","#07689f","#A079BF")) +
    #scale_color_manual(values = c("#FD/63F","#23BAC5")) +
    # facet_grid(TYPE~.,scales = "free_y") +
     theme(legend.position = "left",
         axis.line = element_line(color = "black", linewidth = 0.5), legend.title=element_blank(), ## 图例位置(可选:bottom/right/left/top)
           axis.title = element_text(size = 20),
           strip.text = element text(size = 20)
           # axis.ticks.x = element_blank(),
           legend.text=element text(size = 15),
           strip.text.y = element_text(size = 12, face = "bold"),)+
           # panel.border = element_blank(),
    scale x continuous(breaks = c(164142.97, 184142.97, 184153.61, 204153.61)
                        labels = c("-2Mb", "TSS", "TES", "2Mb")) +
                        #labels = c("LOC Os02g30850")) +
    geom vline(xintercept = 184142.97, col = "black", lty = 2)
    geom_vline(xintercept = 184153.61, col = "black", lty = 2)
   ggsave("Chr01.png", width = 30, height = 10)
```

综合评估

核苷酸多样性(π)定量评估了群体内的遗传变异数量,而 Tajima's D 则描述了其等位基因频率频谱的特征。将这两种统计量相结合,是识别偏离中性进化模式、从而推断自然选择作用的关键策略。

- ▶ 低 π 值 + 显著负的 Tajima's D: 该区域极有可能包含一个经历**强烈正选择的适应性基因**。
- ▶ 正常/高 π 值 + 显著正的 Tajima's D: **平衡选择**,该基因可能参与诸如病原体防御、自交不亲和等过程,其多样性本身具有适应意义。
- ▶ 低 π 值 + Tajima's D ≈ 0 或 轻微为正: 种群规模小或经历过瓶颈效应,无强烈局部选择。
- Arr 正常 π 值 + Tajima's D \approx 0: 中性进化——基因组的大部分常态。
- 低π值+显著正的 Tajima's D(罕见但存在): 长期平衡选择下的局部变异清除,在长期维持两个古老等位基因的同时,基因转换等机制清除了它们之间的重组区域变异,导致局部π值降低,但两个主效等位基因本身仍保持中等频率,从而产生正的D值。

Q&A

汇报人: 丁悦