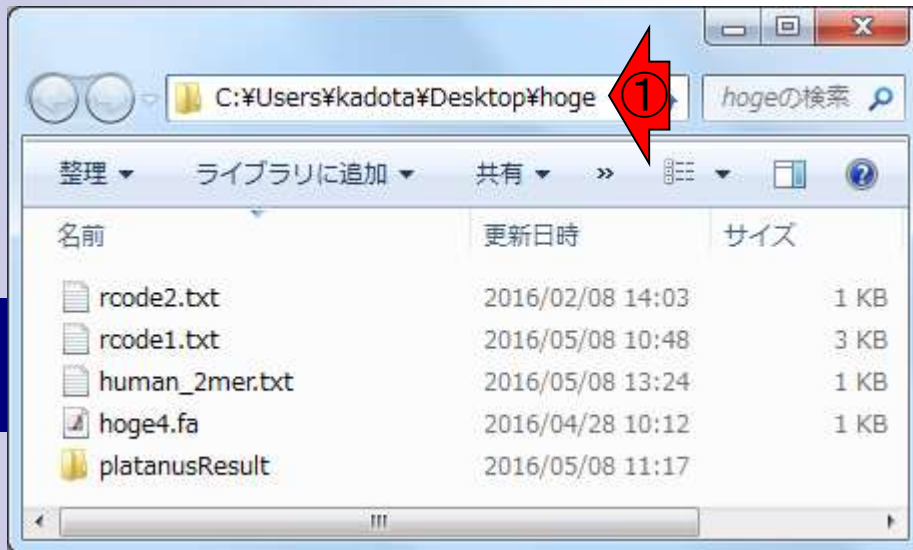


①hogeフォルダがあり、そこに解析用入力ファイルが存在するという前提で行います。印刷は2016.06.23版です。若干見栄えや文言を微修正しました



第1部：統計解析 ～ゲノム解析、塩基配列解析～



東京大学・大学院農学生命科学研究科
アグリバイオインフォマティクス教育研究プログラム

門田幸二(かどた こうじ)

kadota@iu.a.u-tokyo.ac.jp

<http://www.iu.a.u-tokyo.ac.jp/~kadota/>

実習用PC環境は、①の手順に従って「Rおよび必要なパッケージ」のインストールが完了している状態

講習会のPC環境

(Rで)塩基配列解析

～NGS、RNA-seq、ゲノム、トランスクリプトーム、正規化、発現変動、統計、モデル、バイオインフォマティクス～
(last modified 2016/01/29, since 2011)

What's new?

- このウェブページはインストール| | についての推奨手順 ([Windows2015.04.04版](#)と[Macintosh2015.04.03版](#))に従ってフリーソフトRと必要なパッケージをインストール済みであるという前提で記述しています。初心者の方は[基本的な利用法](#)([Windows2015.04.03版](#)と[Macintosh2015.04.03版](#))で自習してください。本ウェブページを体系的にまとめた書籍もあります。(2015/04/03)
- 多群間比較用の推奨ガイドライン提唱論文([Tang et al., BMC Bioinformatics, 2015](#))がpublishされました。論文概要については[門田](#)のページでも紹介しています。講習会でよく述べている「サンプル間クラスタリング結果からDEG検出結果のおおよその見積もりが可能である」という主張の根拠となる原著論文がこれになります。推奨ガイドライン周辺の関連項目もアップデートしました。(2015/11/05) **NEW**
- [日本乳酸菌学会誌](#)の[NGS関連連載](#)の第5回ウェブ資料を更新しました。2015年12月下旬に一気に全てやり直したので、若干プログラムのバージョンが上がっています。(2015/12/22)
- [解析|一般|アラインメント| | について](#)を追加しました。(2015/12/16)
- [日本乳酸菌学会誌](#)の[NGS関連連載](#)の第4回ウェブ資料を更新しました。2015年12月初旬に一気に全てやり直したので、若干プログラムのバージョンが上がっています。各回終了時点のovaファイル(約6GB)も提供可能です。(権利関係上無条件公開はできませんので...) 欲しい方は、メールのタイトルを「乳酸菌連載第x回終了時点のovaファイル希望」として私宛にメールしてください(本文は空でOK)。URLをお知らせします。(2015/12/11)

- [はじめに](#) (last modified 2015/03/31)
- [参考資料\(講義、講習会、本など\)](#) (last modified 2015/11/17)
- [過去のお知らせ](#) (last modified 2015/12/22)
- [インストール| | について](#) (last modified 2015/11/12)
- [インストール| | について](#) (last modified 2015/03/22) 推奨

[トップページへ](#)

予習事項の再確認

具体的な順番は、①R本体のインストール、②各種Rパッケージのインストールです。③基本的な利用法、および④H27年度講習会「データ解析環境R(7月29-30日分)」の自習

(Rで)塩基配列解析

～NGS、RNA-seq、ゲノム、トランスクリプトーム、正規化、発現変動、統計、モデル、バイオインフォマティクス～
(last modified 2016/01/29, since 2011)

What's new?

- このウェブページはインストール|についての推奨手順 (Windows2015.04.04版とMacintosh2015.04.03版)に従ってフリーソフトRと必要なパッケージをインストール済みであるという前提で記述しています。初心者の方は基本的な利用法(Windows2015.04.03版とMacintosh2015.04.03版)で自習してください。本ウェブページを体系的にまとめた書籍もあります。(2015/04/03)
 - 多群間比較用の推奨ガイドライン提唱論文(Tang et al.)については門田のページでも紹介しています。講習会結果のおおよその見積もりが可能である」という点と周辺の関連項目もアップデートしました。(2015/11/05)
 - 日本乳酸菌学会誌のNGS関連連載の第5回ウェブページなので、若干プログラムのバージョンが上がっています
 - 解析|一般|アラインメント|についてを追加しました
 - 日本乳酸菌学会誌のNGS関連連載の第4回ウェブページなので、若干プログラムのバージョンが上がっています(関係上無条件公開はできませんので...)欲しい方「希望」として私宛にメールしてください(本文は空で)
- はじめに (last modified 2015/03/31)
 - 参考資料(講義、講習会、本など) (last modified 2015/11/17)
 - 過去のお知らせ (last modified 2015/12/22)
 - インストール|について (last modified 2015/11/12)
 - インストール | R本体 | 最新版 | Win用 (last modified 2015/03/22)推奨
 - インストール | R本体 | 最新版 | Mac用 (last modified 2015/04/22)推奨
 - インストール | R本体 | 過去版 | Win用 (last modified 2015/03/22)
 - インストール | R本体 | 過去版 | Mac用 (last modified 2015/03/22)
 - インストール | Rパッケージ | ほぼ全て(20GB以上?!)(last modified 2015/05/25)
 - インストール | Rパッケージ | 必要最小限プラスアルファ(数GB?!)(last modified 2015/12/18)推奨
 - インストール | Rパッケージ | 必要最小限プラスアルファ(アグリバイオ|居室のみ)(last modified 2015/06/16)
 - インストール | Rパッケージ | 必要最小限(数GB?!)(last modified 2015/05/25)
 - インストール | Rパッケージ | 個別 (last modified 2015/06/10)
 - (削除予定)Rのインストールと起動 (last modified 2015/04/02)
 - (削除予定)個別パッケージのインストール (last modified 2015/02/20)
 - 基本的な利用法 (last modified 2015/04/03)
 - サンプルデータ (last modified 2015/06/15)
 - バイオインフォマティクス人材育成カリキュラム(次世代シーケンサ)|NGSハンズオン講習会2015 (last modified 2015/06/15)



基本スタンスは、優先順位とエフォート。基本独裁、一匹狼。受益者が金をかけずに効率的に学べる教材整備が最優先。

私の主な活動

- 東大アグリバイオの大学院講義(バイオインフォ全般)
 - Rを中心としたハンズオン講義(平成16年度～)
 - 受講人数が多い(最大130名)ので、クラウド(ウェブツール)系実習は実質的に不可能
 - 講義補助員(TA)が数名のみなので、Linux系実習も困難
- NBDC/東大アグリバイオのNGSハンズオン講義(NGSに特化)
 - Linuxを中心としたハンズオン講義(平成26年度～)
 - 受講人数は多い(H27年度は最大71名;おそらくアグリバイオ本体に次ぐ規模)が、受講生の意識レベルが高く(きっちり予習をやるヒトが多数派)、環境構築済みノートPC数、TA数が充実しているため、本格的なLinux実習が成立しうる。
- 日本乳酸菌学会誌のNGS連載
 - Linuxを中心とした自習用教材(平成26年度～)
 - バクテリア(乳酸菌)データを、主にBio-Linux上で解析するノウハウを提供。
 - 第6回(2016年3月)分以降は、DDBJ Pipeline(ウェブツール)の利用法も紹介。
 - データ取得・インストール・実行に時間がかかるものも、自習なので時間を気にせずにできる。ハンズオン講義よりも心穏やか。
- その他
 - 研究(発現変動解析精度向上のためのアルゴリズム開発や評価)
 - ...


Contents

- NGS解析手段、ウェブツール(DDBJ Pipeline)との連携
- DDBJ PipelineでPlatanusを実行
- k-mer解析(k個の連続塩基に基づく各種解析)の基礎
 - 塩基ごとの出現頻度解析(k=1)、一気に計算
 - 2連続塩基の出現頻度解析(k=2)、基本スキルの復習や作図
- de novoアセンブリ時のエラー補正やゲノムサイズ推定の基本的な考え方
 - ランダムな塩基配列(仮想ゲノムおよび仮想NGSデータ)の作成
 - k-mer解析の応用、ゲノムサイズ推定の基礎
 - ゲノムサイズ推定(1,000 bpの仮想ゲノムの場合)
 - ゲノムサイズ推定(1,000 bpの仮想ゲノム; 4X → 10X coverageの場合)
 - k-mer出現頻度分布
 - シークエンスエラーを含む場合
 - 最終確認



NGSデータ解析戦略

自分の置かれている環境・予算・ポリシーによっても異なる。どの選択肢でも基本正解。Rは、主に統計解析部分で使われる

- 解析受託企業に外注：Linuxコマンドを知らなくてもよい
 -  ngs 受託解析
- クラウド（ウェブツール）：Linuxコマンドを知らなくてもよい
 - DDBJ Pipeline (Nagasaki et al., *DNA Res.*, 20: 383–390, 2013)
 - Illumina BaseSpace
 - Galaxy (Goecks et al., *Genome Biol.*, 11: R86, 2010)
 - ...
- Linuxコマンドを駆使（旧来型）
 - なるべく自力で解析
 - LinuxコマンドやNGS解析用プログラムのインストールなどを練習し、スパコンを使いこなす
 - NBDC/東大アグリバイオの「NGSハンズオン講習会」の方向性

DDBJ PipelineとR

①DDBJ Pipelineだけで全てのNGS解析ができるわけではない。②Rもまた然り。特にRでは、(門田の知る限り) *de novo*アセンブリは不可能。現実を知り、うまく使い分けるべし。DDBJ Pipeline上で*de novo*アセンブリを行った結果の解釈や確認をRで行うところからスタート

■ 解析受託企業に外注: Linuxコマンド

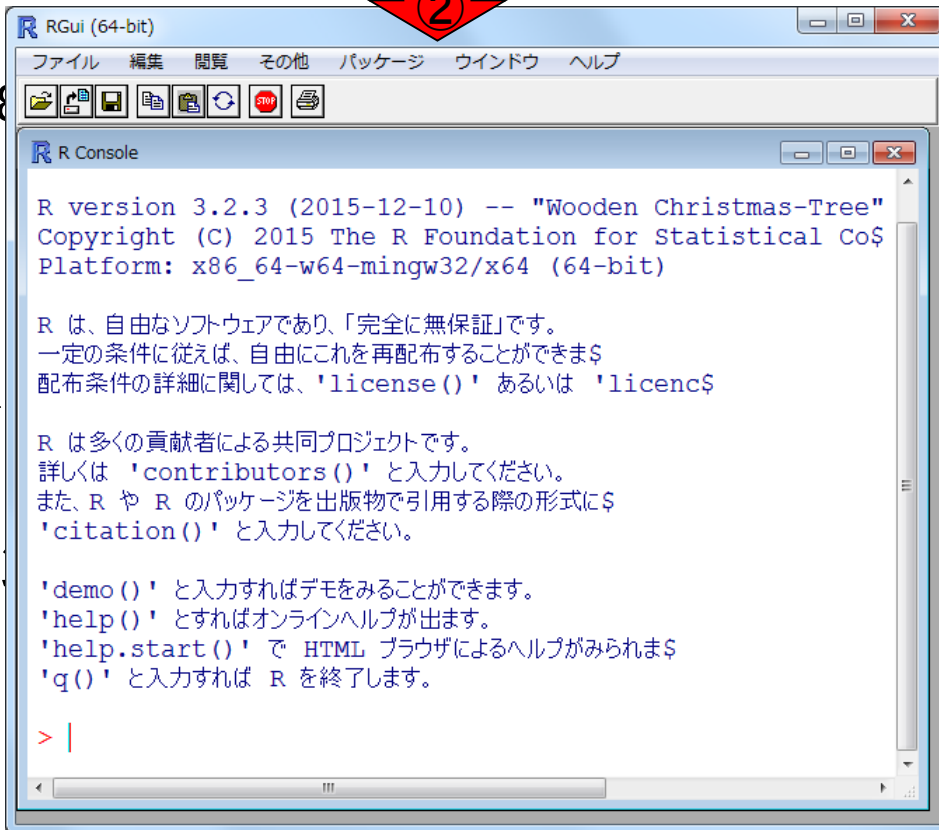
-  ngs 受託解析

■ クラウド(ウェブツール): Linuxコマンドを知らなくてもよい

- ① □ DDBJ Pipeline (Nagasaki et al., *DNA Res.*, 20: 383-390, 2013)
- Illumina BaseSpace
- Galaxy (Goecks et al., *Genome Biol.*, 11: R119, 2010)
- ...

■ Linuxコマンドを駆使(旧来型)

- なるべく自力で解析
- LinuxコマンドやNGS解析用プログラムのインストール/アップデートをこなし
- NBDC/東大アグリバイオの「NGSハンズオン」



```
R GUI (64-bit)
ファイル 編集 閲覧 その他 パッケージ ウィンドウ ヘルプ

R Console
R version 3.2.3 (2015-12-10) -- "Wooden Christmas-Tree"
Copyright (C) 2015 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R は、自由なソフトウェアであり、「完全に無保証」です。
一定の条件に従えば、自由にこれを再配布することができます。
配布条件の詳細に関しては、'license()' あるいは 'licenc$'
を参照してください。

R は多くの貢献者による共同プロジェクトです。
詳しくは 'contributors()' と入力してください。
また、R や R のパッケージを出版物で引用する際の形式に '$
'citation()' と入力してください。

'demo()' と入力すればデモをみることができます。
'help()' とすればオンラインヘルプが出ます。
'help.start()' で HTML ブラウザによるヘルプがみられます。
'q()' と入力すれば R を終了します。

> |
```

DDBJ Pipeline

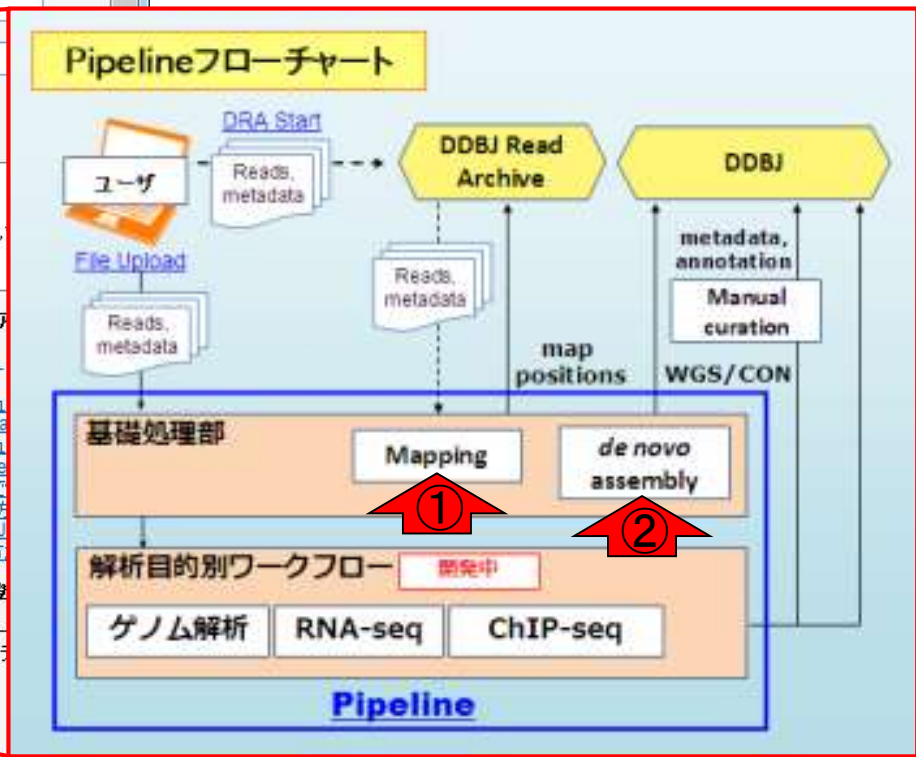
DDBJ Pipeline(の基礎処理部)では、主に①マッピングや② *de novo*アセンブリができる。特に後者ができるのは非常に有難い。③新規アカウント作成から*de novo*アセンブリまでの詳細手順は、乳酸菌NGS連載第6回ウェブ資料を参照

The screenshot shows the DDBJ Read Annotation Pipeline website. At the top, there is a navigation bar with the DDBJ logo and the text "DDBJ Read Annotation Pipeline". Below this, there are language selection buttons for "English" and "Japanese". A description states: "DDBJ Read Annotation Pipelineは、次世代シーケンサ配列のクラウド型データ解析プラットフォームです。" (DDBJ Read Annotation Pipeline is a cloud-based data analysis platform for next-generation sequencing data.)

The main content area is titled "LOGIN" and includes a "新規アカウント作成" (New Account Creation) link. There are input fields for "User ID:" and "Password:" and a "Login" button. Below the login form, there is a section for "動作中JOBの確認" (Check running jobs) with a note: "PipelineのIDをお持ちでないことができます。" (You can check jobs without the Pipeline ID.)

On the left side, there is a "Pipelineフローチャート" (Pipeline Flowchart) section. It shows a flow from "ユーザ" (User) through "DRA Start" (Reads, metadata) to "DDBJ Read Archive" (Reads, metadata) and finally to "DDBJ" (metadata, annotation, Manual curation). Below this, it details the "基礎処理部" (Basic Processing Unit) which includes "Mapping" and "de novo assembly". A "解析目的別ワークフロー" (Workflow by analysis purpose) is also shown, including "ゲノム解析" (Genome analysis), "RNA-seq", and "ChIP-seq".

At the bottom left, there is a "Tweets" section with a tweet from @pipeline_info dated 22 Dec, stating: "DDBJ Pipeline e-mail support will not open from Dec 25 - Jan 3. Thank you for your cooperation."



NGSデータ

乳酸菌(*Lactobacillus hokkaidonensis* LOOC260^T) ゲノム解読論文(PMID: 25879859)。Illumina MiSeqデータ(DRR24501)とPacBioデータ(DRR024500*)を併用することでcomplete genomeを得ることができた、という内容。講習会では、MiSeqデータをDDBJ Pipelineで*de novo*アセンブリした結果をRで解析し、アセンブリ過程の理解を深める。これは*k*-mer解析のイントロでもある

The screenshot shows the PubMed interface for the article. The URL is <http://www.ncbi.nlm.nih.gov/pubmed/25879859>. The article title is "Complete genome sequence and analysis of *Lactobacillus hokkaidonensis* LOOC260(T), a psychrotrophic lactic acid bacterium isolated from silage." The authors listed are Tanizawa Y^{1,2}, Tohno M³, Kaminuma E⁴, Nakamura Y⁵, and Arita M^{6,7}. The abstract text is visible, starting with "BACKGROUND: *Lactobacillus hokkaidonensis* is an obligate heterofermentative lactic acid bacterium, which is isolated from Timothy grass silage in Hokkaido, a subarctic region of Japan. This bacterium is expected to be useful as a silage starter culture in cold regions because of its remarkable psychrotolerance; it can grow at temperatures as low as 4°C. To elucidate its genetic background, particularly in relation to the source of psychrotolerance, we constructed the complete genome sequence of *L. hokkaidonensis* LOOC260(T) using PacBio single-molecule real-time sequencing technology." The results and conclusions sections are also partially visible. On the right side of the page, there are links for "Full text links" (BioMed Central, PMC Full text), "Save items" (Add to Favorites), "Similar articles" (listing related papers), and "Related information".

NGSデータ

①Full textリンク先で全文を見られる。②Availability of supporting dataという項目をよく眺めると、生データがDDBJ Sequence Read Archive (DDBJ SRA; 略してDRA)に DRR024500とDRR24501というIDで登録されていることがわかる。尚、DRR024500は登録内容の問題が判明し、2016年1月末に削除されDRR054113-054116に差し替えられている

Abstract

BMC Genomics. 2015 Mar 25;16:240. doi: 10.1186/s12864-015-1435-2.

Complete genome sequence and analysis of *Lactobacillus hokkaidonensis* LOOC260(T), a psychrotrophic lactic acid bacterium isolated from silage.

Tanizawa Y^{1,2}, Tohno M³, Kaminuma E⁴, Nakamura Y⁵, Arita M^{6,7}.

Author information

Abstract

BACKGROUND: *Lactobacillus hokkaidonensis* is an obligate heterofermentative lactic acid bacterium, which is isolated from Timothy grass silage in Hokkaido, a subarctic region of Japan. This bacterium is expected to be useful as a silage starter culture in cold regions because of its remarkable psychrotolerance; it can grow at temperatures as low as 4°C. To elucidate its genetic background, particularly in relation to the source of psychrotolerance, we constructed the complete genome sequence of *L. hokkaidonensis* LOOC260(T) using PacBio single-molecule real-time sequencing technology.

RESULTS: The genome of LOOC260(T) comprises one circular chromosome (2.28 Mbp) and two circular plasmids: pLOOC260-1 (81.6 kbp) and pLOOC260-2 (41.0 kbp). We identified diverse genetic elements, such as prophages, integrated and conjugative elements, and conjugative plasmids, which may reflect adaptation to plant-associated niches. Comparative genome analysis also revealed unique genomic features, such as genes involved in pentose assimilation and NADPH-dependent reactions.

CONCLUSIONS: This is the first complete genome in the *L. vaccinostercus* group, which is characterized, so the genomic information obtained in this study provides insight into the genetic evolution of this group. We also found several factors that may contribute to the ability of *L. hokkaidonensis* to grow at cold temperatures. The results of this study will facilitate further research on the cold-tolerance mechanism of *L. hokkaidonensis*.

PMID: 25879859 [PubMed - in process] PMID: PMC4377027 [Free PMC Article](#)

Full text links

Read free full text at BioMed Central

PMC Full text

Save items

Add to Favorites

Similar articles

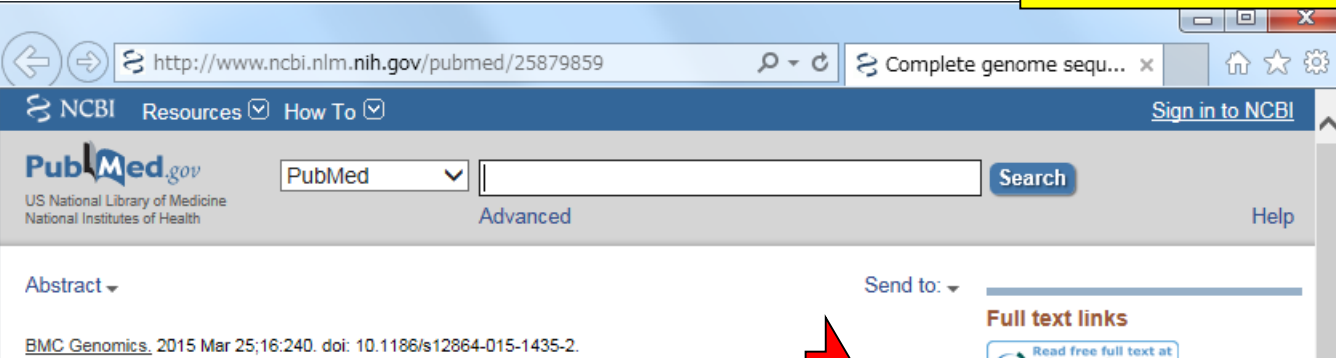
Availability of supporting data

The complete genome sequence of *L. hokkaidonensis* LOOC260^T and its annotations were deposited at DDBJ/ENA/GenBank under accession numbers AP014680 (chromosome), AP014681 (plasmid pLOOC260-1), and AP014682 (plasmid pLOOC260-2). All of the sequencing data were deposited in the DDBJ Sequence Read Archive under accession numbers DRR024500 and DRR024501. The phylogenetic tree and associated data matrix for in Additional file 1: Figure S2 are available in TreeBASE database (Accession URL: <http://purl.org/phylo/treebase/phylovs/study/TB2:S17206>).

Related information

NGSデータ

① Genome sequencing and *de novo* assemblyという項目を見ると、② paired-endで5,942,620リードと書いてある。一応公共DB(DRA)上で確認する



Complete genome sequence and analysis of *Lactobacillus hokkaidensis* LOOC260(T), a psychrotrophic lactic acid bacterium isolated from Timothy grass silage in Hokkaido, Japan

Tanizawa Y^{1,2}, Tohno M³, Kaminuma E⁴, Nakamura Y⁵, Arita M^{6,7}.

Author information

Abstract

BACKGROUND: *Lactobacillus hokkaidensis* is an obligate heterofermentative lactic acid bacterium which is isolated from Timothy grass silage in Hokkaido, a subarctic region of Japan. This bacterium is expected to be useful as a silage starter culture in cold regions because of its remarkable psychrotolerance; it can grow at temperatures as low as 4°C. To elucidate its genetic background, particularly in relation to the source of psychrotolerance, we constructed the complete genome sequence of *L. hokkaidensis* LOOC260(T) using PacBio single-molecule real-time sequencing technology.

RESULTS: The genome of LOOC260(T) comprises one circular chromosome (2.28 Mbp) and two circular plasmids: pLOOC260-1 (81.6 kbp) and pLOOC260-2 (41.0 kbp). We identified diverse genetic elements, such as prophages, integrated and conjugative elements, and conjugative plasmids, which may reflect adaptation to plant-associated niches. Comparative genome analysis also revealed unique genomic features, such as genes involved in pentose assimilation and NADPH-dependent reactions.

CONCLUSIONS: This is the first complete genome in the *L. vaccinostercus* group, which is characterized, so the genomic information obtained in this study provides insight into the genetic evolution of this group. We also found several factors that may contribute to the ability of *L. hokkaidensis* to grow at cold temperatures. The results of this study will facilitate further research on the cold-tolerance mechanism of *L. hokkaidensis*.

PMID: 25879859 [PubMed - in process] PMCID: PMC4377027 Free PMC Article



① Genome sequencing and *de novo* assembly

The cells of *L. hokkaidensis* LOOC260^T were cultured in MRS (de Man, Rogosa, and Sharpe) broth (Difco) and were harvested in the mid-logarithmic phase. The genomic DNA was extracted and purified using Qiagen Genomic-tip 500/G and Qiagen Genomic DNA Buffer Set with lysozyme (Sigma) and proteinase K (Qiagen) according to the manufacturer's instruction. PacBio SMRT whole-genome sequencing was performed using a PacBio RSII sequencer with P4-C2 chemistry. Four SMRT cells were used for sequencing, thereby yielding 163,376 adapter-trimmed reads (subreads) with an average read length of approximately 4 kbp, which corresponded to approximately 250-fold coverage. *De novo* assembly was conducted using the HGAP method based on the SMRT Analysis package 2.0, which yielded seven contigs. Independent genome sequencing using the 250-bp paired-end Illumina MiSeq system generated 5,942,620 reads, which were assembled into contigs using Platanus assembler v1.2.2 with the default settings [40]. The initial contigs

Related information

Assembly

NGSデータ

①Genome sequencing and *de novo* assemblyという項目を見ると、②paired-endで5,942,620リードと書いてある。③MiSeqデータに相当するDRR024501は、④2,971,310リードと書いてある。5,942,620/2 = 2,971,310である。ウェブサイト上の数値は、single-endとしてのリード数と考えれば妥当

DRASearch

DRR024501 FASTQ SRA

Run Detail	
Alias	DRR024501
Instrument model	
Date of run	
Run center	
Number of spots	2,971,310
Number of bases	1,491,597,620

Navigation

- Submission DRA002643 FTP
- Study DRP002401

READS (joined) quality show 10 rows << < 1

```
>DRR024501.1
ATGNATCGAAACAGTATTTACAAGATTTGCATACTGAAATTTGAAGCTGATCAACACGAAACCATTC
AATCTAAACCCACCCATTAGCTGTTATTGAAGCTTTGCAGCAACGAGTTGATGATAAAATGACCGTTT
GAGCCATTATATTTGGATGGCCCGGCACCTCCGAAGTTATGAGCCTCGCCATTTATTGTTTAGTAA
TTGGAGTGGCGATGAACCGTATTAAGCCCTAAACGAACGGCTGTCTCCAGTCTTGTCCAGTAAATA
CCAGAAACAGAGACTGATTTAGCATTGGGCCGAACCTAACCGCAGCCGAAATTTGACCAAGGTAGCGC
GCATCCCACTTAACAATAAATGGCGAGGCTCATAACTTCGGAAGTGCCGGGCCATCCAAATATA
TCAAACGCCAAGTTCATATTAT

>DRR024501.2
GTCNGAACACATGAATGGTGAACGGCGCTGAACTTTTACGGACGCGGCACGAGGATCCACAGGGC
AAACAGCTCAACGACTTGTATCACCGCATTATTACGTGAGTGGATTTCTGATAAAACAGTGTAA
AATTGATTTCTAGCCAATCAAAGACAGATTTAACGAAATCACAGATGACCAAAATTTGGCATTGAGT
TAAAAAAGTGATGGACCCCTCTTTAACCCTAAGTTGTCCCGAATAACATTTCGAAACTCTCTGCTTT
ACCTCAAATGATTTGCCCAATCAATTTGACTTGTTCCTTGTGTCATGATCTGATTTCACTGTTTAT
ACTCAATGCCAATTTGGTCATCTGTGATTTCTGTTAACTCTGTCTTTGATTGGCTAGACATCGAATT
TTCACACTGTTTTATACGAAAT

>DRR024501.3
CAANGATACAATCATTATCATGAACCTAATGCCGGTTCCTGGTATGCACTTCTACTGTTTGGCTTGAAGCC
CTGGTGACACGCACTCAGTTCCTTAACAAACTAGGCGATTAT
```

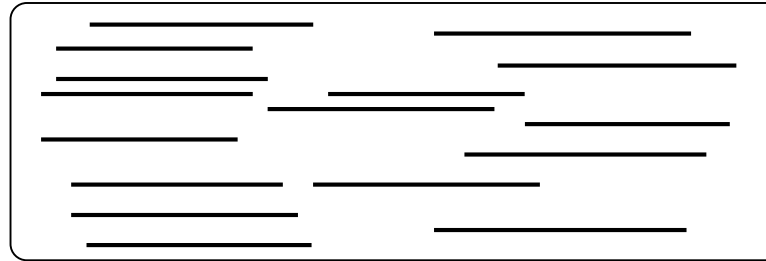
Genome sequencing and *de novo* assembly

The cells of *L. hokkaidonensis* LOOC260^T were cultured in MRS (de Man, Rogosa, and Sharpe) broth (Difco) and were harvested in the mid-logarithmic phase. The genomic DNA was extracted and purified using Qiagen Genomic-tip 500/G and Qiagen Genomic DNA Buffer Set with lysozyme (Sigma) and proteinase K (Qiagen) according to the manufacturer's instruction. PacBio SMRT whole-genome sequencing was performed using a PacBio RSII sequencer with P4-C2 chemistry. Four SMRT cells were used for sequencing, thereby yielding 163,376 adapter-trimmed reads (subreads) with an average read length of approximately 4 kbp, which corresponded to approximately 250-fold coverage. *De novo* assembly was conducted using the HGAP method based on the SMRT Analysis package 2.0, which yielded seven contigs. Independent genome sequencing using the 250-bp paired-end Illumina MiSeq system generated 5,942,620 reads, which were assembled into contigs using Platanus assembler v1.2.2 with the default settings [40]. The initial contigs

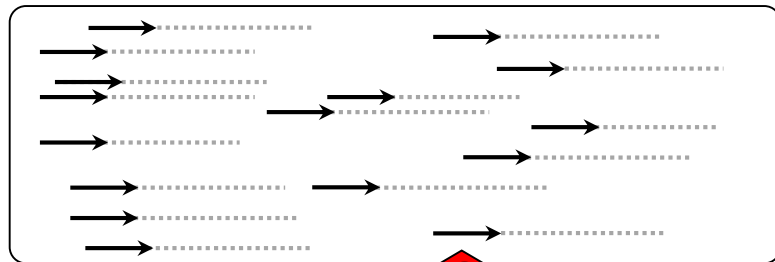
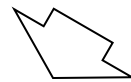
リードとは、Sequencerで読んだ塩基配列のこと。①黒矢印の一本一本がリードに相当する

用語：リード

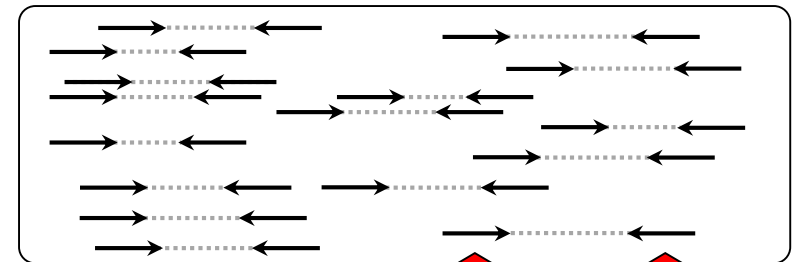
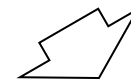
断片化されたゲノム配列



single-endの場合



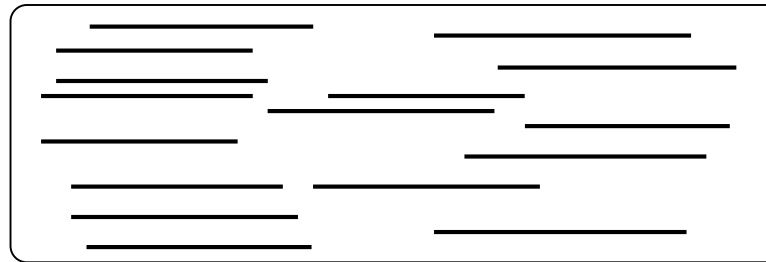
paired-endの場合



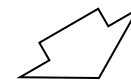
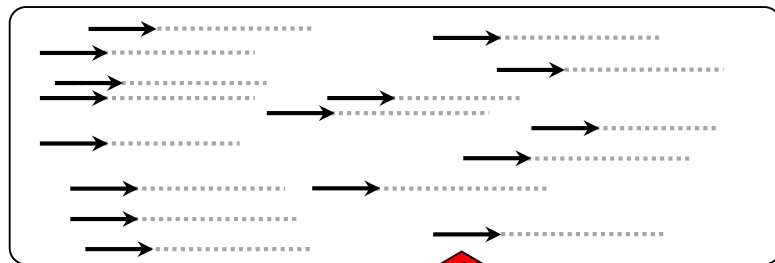
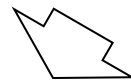
①断片化された配列の片側のみを読む場合を single-end という。この場合は右向き矢印のみ

用語 : single-end

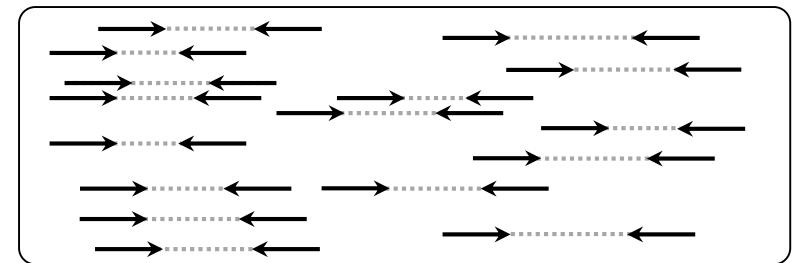
断片化されたゲノム配列



single-endの場合



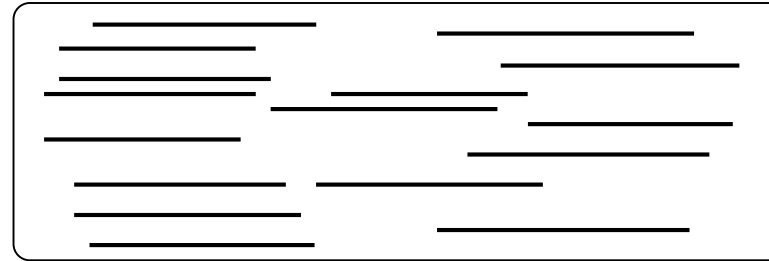
paired-endの場合



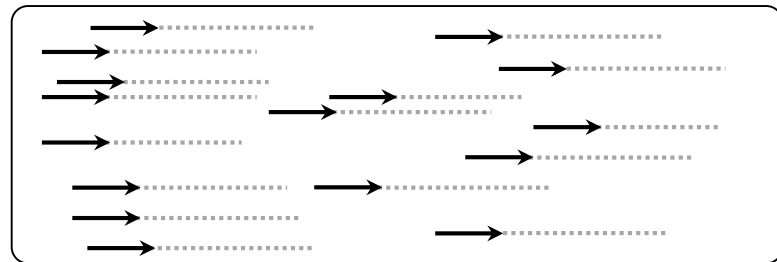
用語: paired-end

①断片化された配列の両側から読む場合を paired-end という。②右向き矢印と③左向き矢印のリードが読まれることになる。それぞれを forward側リード、reverse側リードなどと呼ぶ

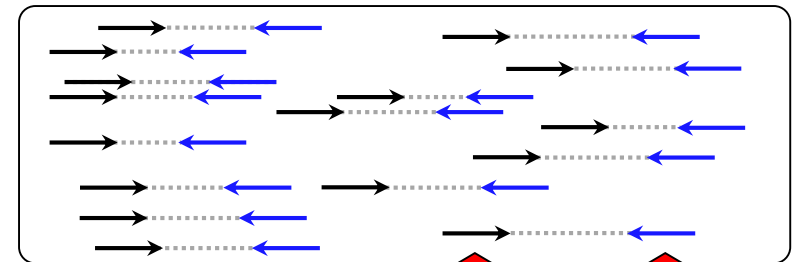
断片化されたゲノム配列



single-endの場合



paired-endの場合



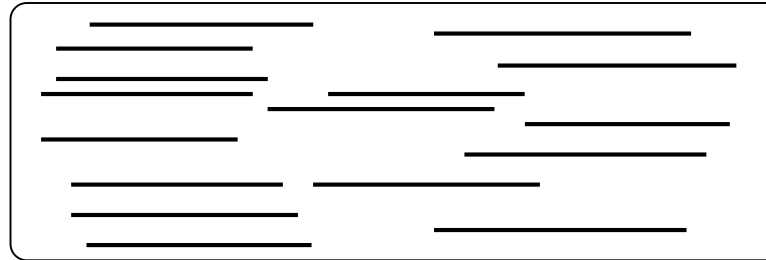
forward側

reverse側

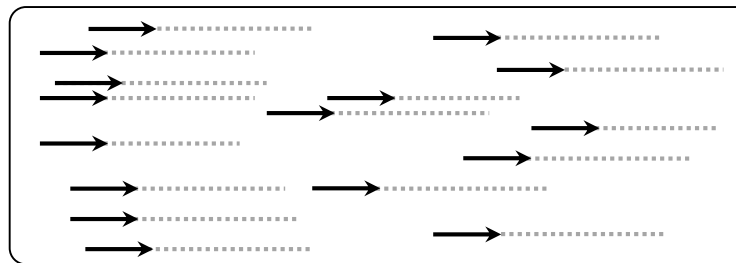
用語 : paired-end

ILLUMINA MiSeqデータ(DRR24501)の場合、① forward側、②reverse側ともに、矢印の長さが250 bp、矢印の本数(リード数)が計5,942,620個(約594万;片側のみで約297万)に相当

断片化されたゲノム配列

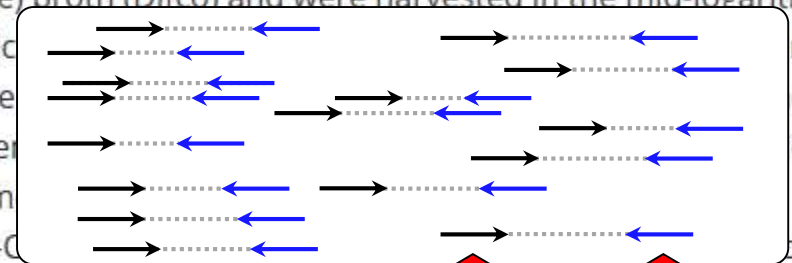


single-endの場合



Genome sequencing and *de novo* assembly

The cells of *L. hokkaidonensis* LOOC260^T were cultured in MRS (de Man, Rogosa, and Sharpe) broth (Difco) and were harvested in the mid-logarithmic phase. The genomic DNA was extracted using the DNeasy spin 500/G and Qiagen proteinase K (Qiagen) and sequenced using the SMRT whole-genome sequencer with P4-C100 chemistry, thereby yielding 163,376 adapter-trimmed reads with an average read length of approximately 4 kbp, which corresponded to approximately 250-fold coverage. *De novo* assembly was conducted using the HGAP method based on the SMRT Analysis package 2.0, which yielded seven contigs. Independent genome sequencing using the 250-bp paired-end Illumina MiSeq system generated 5,942,620 reads, which were assembled into contigs using Platanus assembler ver 1.2 with the default settings [40]. The initial contigs



① ②

DDBJ SRA (DRA)

①forward側 : DRR024501_1.fastq.bz2、
 ②reverse側 : DRR024501_2.fastq.bz2、
 のような感じ。DRAの場合は、bzip2圧縮FASTQファイルをダウンロード可能。乳酸菌ゲノム配列決定論文では、このデータを入力として *de novo* アセンブリが行われた

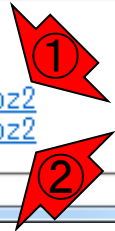
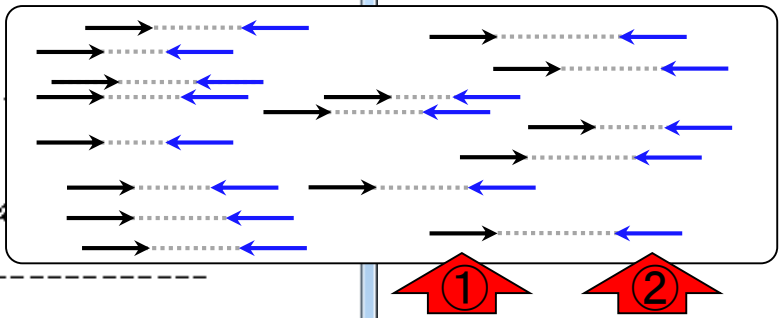
FTP ディレクトリ /ddbj_database/dra/fastq/DRA002/DRA002643/D
 ftp.ddbj.nig.ac.jp

エクスプローラーでこの FTP サイトを表示するには、Alt キーを押して、[表示] をクリックし、[エクスプローラーで FTP サイトを開く] をクリックしてください。

-Welcome to DDBJ FTP Archive, running on ftp.ddbj.nig.ac.jp!
 -
 -Please contact ddbj@ddbj.nig.ac.jp when you have any problem
 -access to this archive, downloading the data, and etc.
 -
 -For details on the directory structure and file contents, please refer to the README.TXT placed in the "/ddbj_database/".

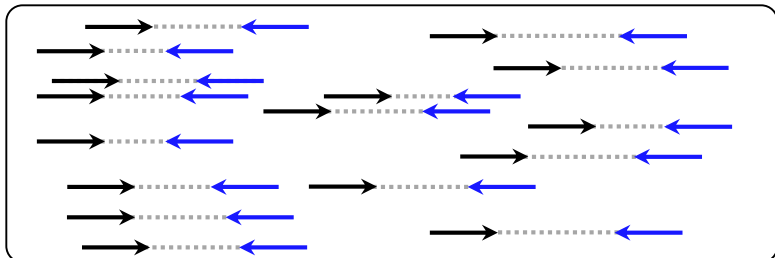
1 階層上のディレクトリへ

11/12/2014 12:00午前	470,724,676	DRR024501_1.fastq.bz2
11/12/2014 12:00午前	528,430,059	DRR024501_2.fastq.bz2

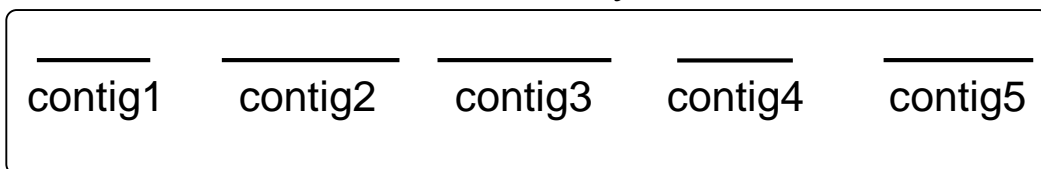


用語：コンティグ

入力：paired-end FASTQファイル



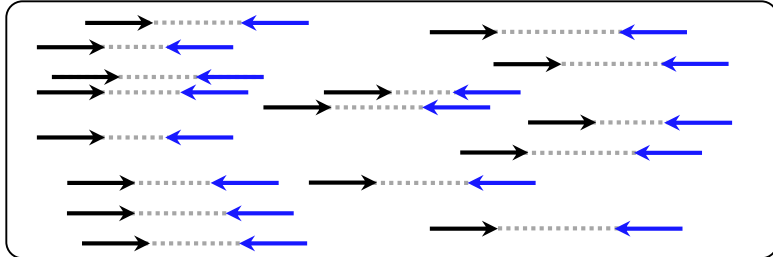
↓ Assembly (コンティグの作成)



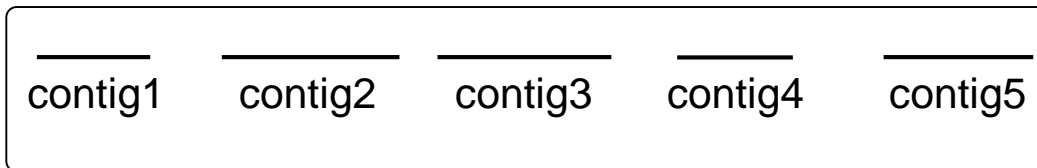
(通常は) paired-endのリードファイルを入力として、*de novo*アセンブリプログラムを実行した結果として得られる、異なる複数のリードが(ACGTの切れ目なく)つなげられたもの。contiguous sequence (連続的な配列) という意味。通常、元のリード長よりも長くなる

用語 : scaffold

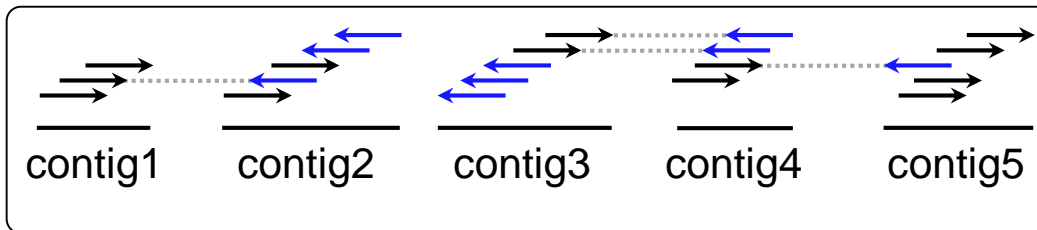
入力 : paired-end FASTQファイル



↓ Assembly (コンティグの作成)

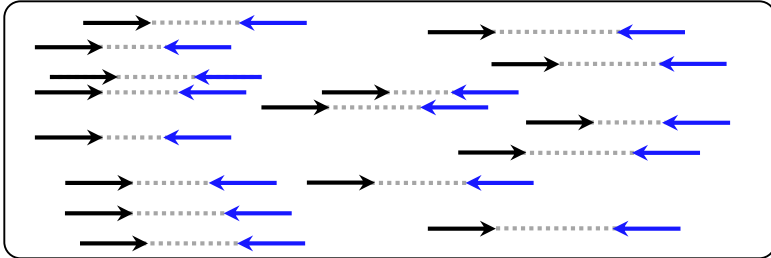


↓ Scaffold

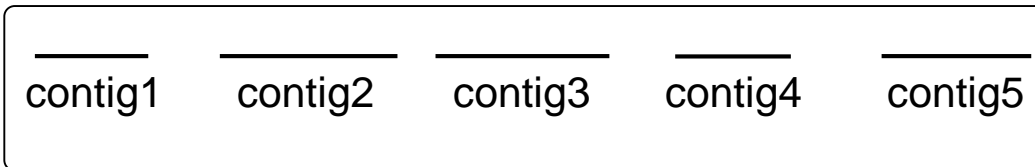


用語 : scaffold

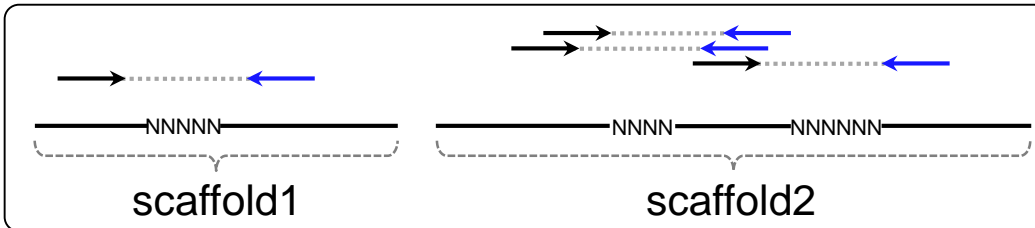
入力 : paired-end FASTQファイル



↓ Assembly (コンティグの作成)



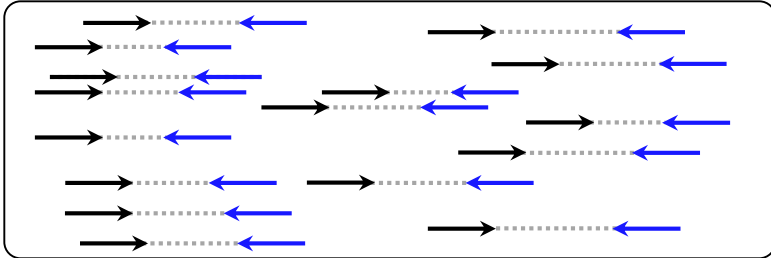
↓ Scaffold



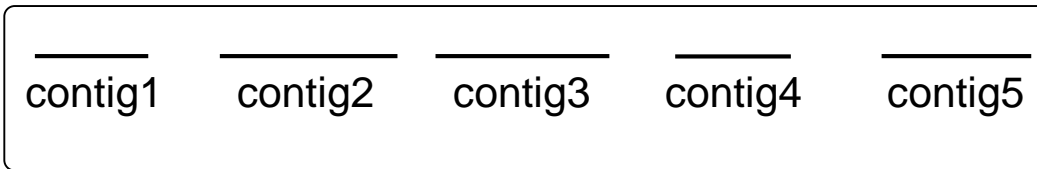
得られたコンティグにリードをマップし…ペアの情報を頼りにコンティグ間にNを入れて連結したもの。supercontigともいう。scaffoldの数はcontigの数よりも少なくなる。尚、Nを入れた部分をgapという

用語: gap close

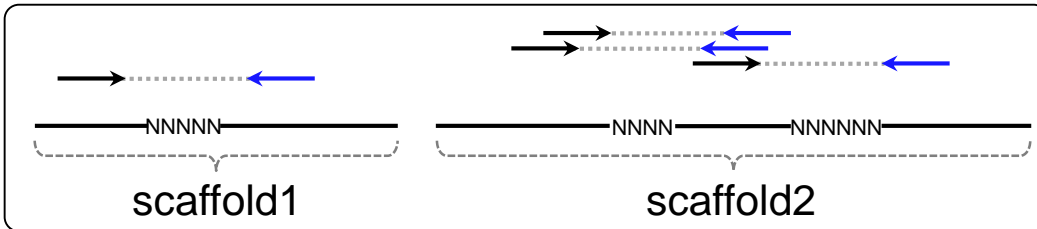
入力: paired-end FASTQファイル



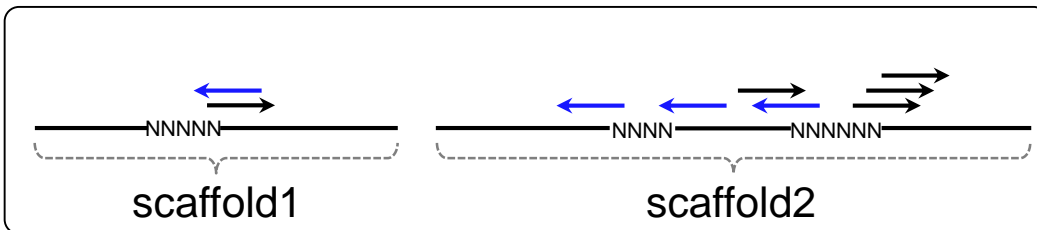
Assembly (コンティグの作成)



Scaffold

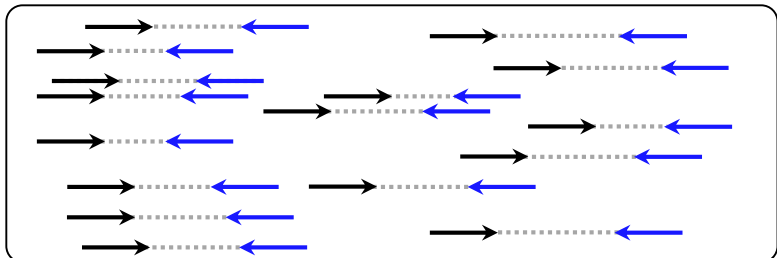


Gap close



用語: gap close

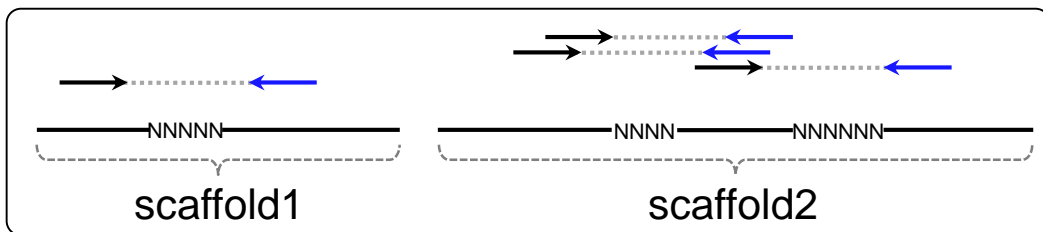
入力: paired-end FASTQファイル



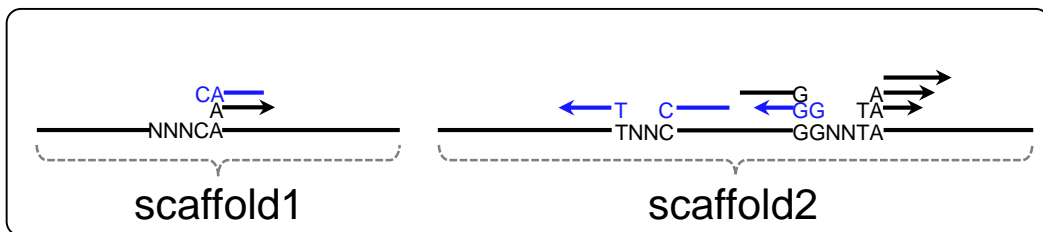
Assembly (コンティグの作成)



Scaffold



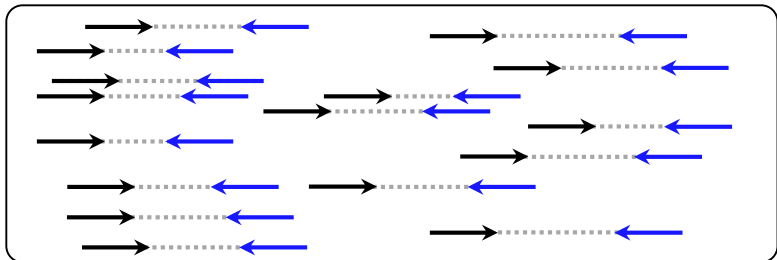
Gap close



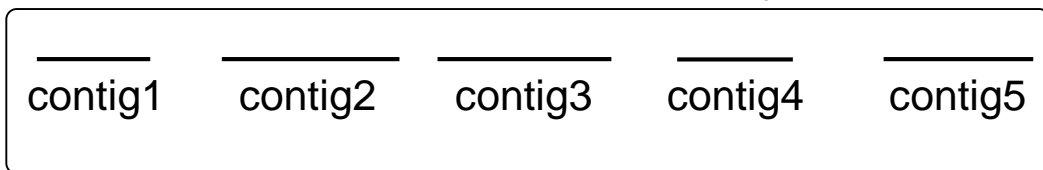
得られたscaffoldsにリードをマップし…gap
周辺にマップされたリードの塩基でNを置
換。gapのNがなくなり、閉じていく(close)の
でgap closeという(おそらく)

de novoアセンブリ

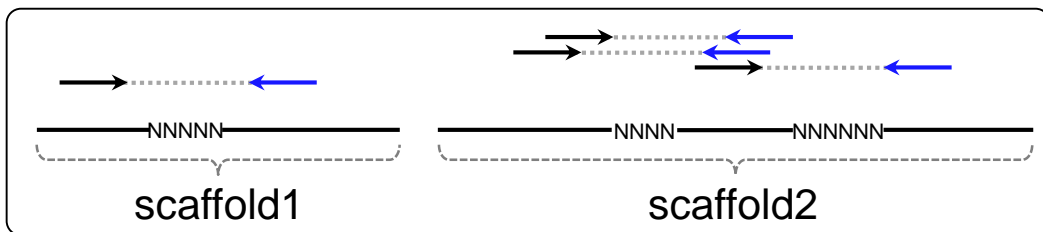
入力: paired-end FASTQファイル



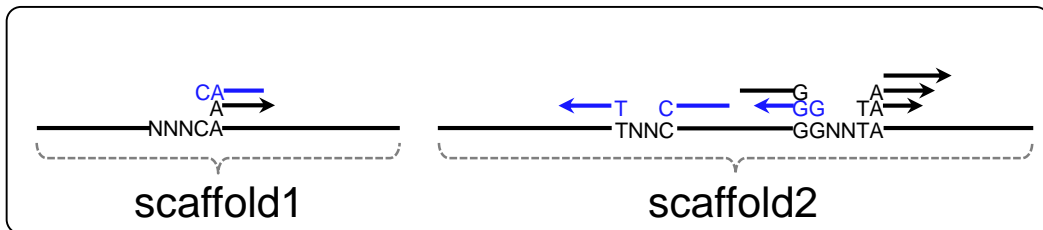
Step1: Assembly



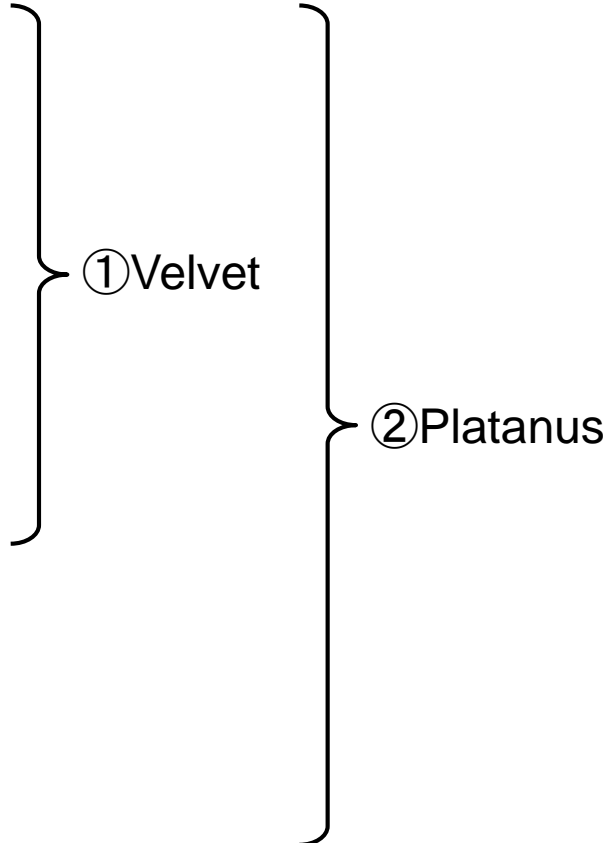
Step2: Scaffold



Step3: Gap close



①最も有名なNGSデータ用 *de novo* ゲノムアセンブリプログラムであるVelvet (Zerbino and Birney, *Genome Res.*, 2008)は、Step2までを実行。②比較的最近開発されたPlatanus (Kajitani et al., *Genome Res.*, 2014)は、Step3までを実行してくれる



乳酸菌論文は...

乳酸菌(*Lactobacillus hokkaidonensis* LOOC260^T) ゲノム解読論文では、Illumina MiSeqデータ(DRR24501)の*de novo*アセンブリに①Platanus (ver. 1.2)を利用

DRASearch

DRR024501 FASTQ SRA

Run Detail	
Alias	DRR024501
Instrument model	
Date of run	
Run center	
Number of spots	2,971,310
Number of bases	1,491,597,620

Navigation

- Submission DRA002643 FTP
- Study DRP002401

READS (joined) quality show 10 rows << < 1

```
>DRR024501.1
ATGNATCGAAACAGTATTTACAAGATTTGCATACTGAAATTTGAAGCTGATCAACACGAAACCATTC
AATCTAAACACCCATTAGCTGTTATTGAAGCTTTGCAGCAACGAGTTGATGATAAAATGACCGTTT
GAGCCATTATATTTGGATGGCCCGGCACCTCCGAAGTTATGAGCCTCGCCATTTATTGTTTAGTAA
TTGGAGTGGCGATGAACCGTATTAAGCCCTAAACGAACGGCTGCTCCAGTTCTTGTCCAGTAAAT
CCAGAAACAGAGACTGATTTAGCATTGGGCCGAACCTAACCGCAGCCGAAATTTGACCAAGGTAGCGC
GCATCCCACTTAACAATAAATGGCGAGGCTCATAACTTCGGAAGTGCCGGGCCATCCAAATATA
TCAAACGCAACGTTTCATATTAT

>DRR024501.2
GTCNGAACACATGAATGGTGAACGGCGCTGAACTTTTACGGACGCGGCACGAGGATCCACAGGGC
AAACACGTCACAGCCTTGTATCACCGCATTATTACGTGAGTGGATTTCTGATAAAACAGTGTAA
AATTCGATTTCTAGCCAATCAAAGACAGATTTAACGAAATCACAGATGACCAAAATTTGGCATTGAGT
TAAAAAAGTGATGGACCCCTCTTTAACCCTAAGTTGTCCCGAATAACATTCGAAACTCTCTGCTTT
ACCTCAAATGATTTGCCCAATCAATTTGACTTGTTCCTTGTGCATGATCTGATTTCACTGTTTAT
ACTCAATGCCAATTTGGTCATCTGTGATTTCTGTTAACTCTGTCTTTGATTGGCTAGACATCGAA
TTCACACTGTTTTATACGAAAT

>DRR024501.3
CAANGATACAATCATTATCATGAACCTAATGCCGGTCTCGGTGATGACGTTGCTAGTGTGCGTGAAGCC
CTGGTGACACGCACTCAGTTCTTAACAACTAGGCGATTA
```

Genome sequencing and *de novo* assembly

The cells of *L. hokkaidonensis* LOOC260^T were cultured in MRS (de Man, Rogosa, and Sharpe) broth (Difco) and were harvested in the mid-logarithmic phase. The genomic DNA was extracted and purified using Qiagen Genomic-tip 500/G and Qiagen Genomic DNA Buffer Set with lysozyme (Sigma) and proteinase K (Qiagen) according to the manufacturer's instruction. PacBio SMRT whole-genome sequencing was performed using a PacBio RSII sequencer with P4-C2 chemistry. Four SMRT cells were used for sequencing, thereby yielding 163,376 adapter-trimmed reads (subreads) with an average read length of approximately 4 kbp, which corresponded to approximately 250-fold coverage. *De novo* assembly was conducted using the HGAP method based on the SMRT Analysis package 2.0, which yielded seven contigs. Independent genome sequencing using the 250-bp paired-end Illumina MiSeq system generated 5,942,620 reads, which were assembled into contigs using Platanus assembler ver 1.2 with the default settings [40]. The initial contigs

Contents

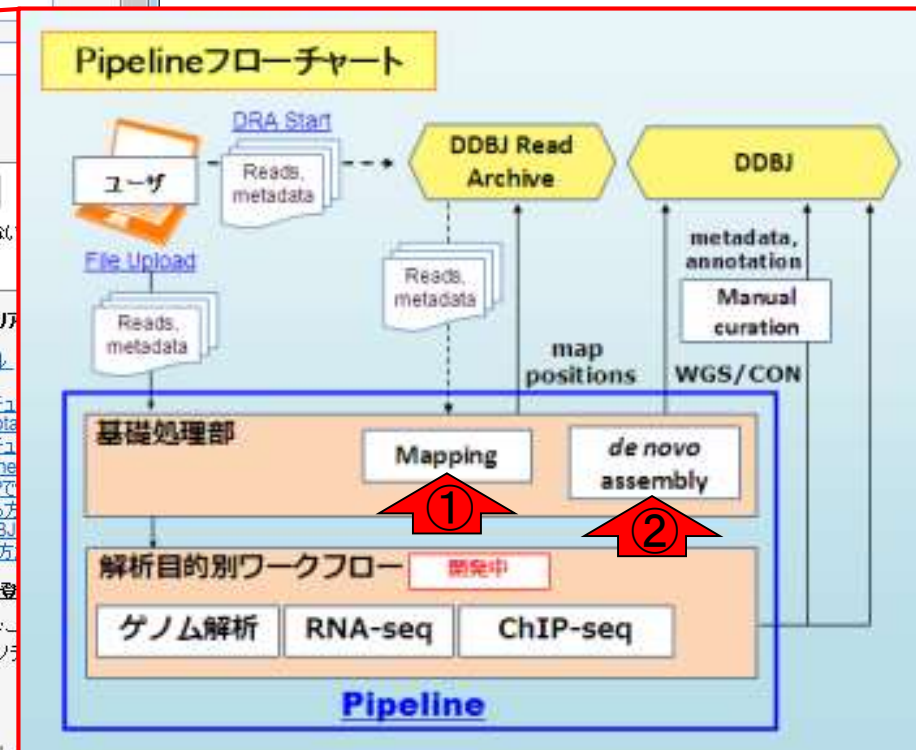
- NGS解析手段、ウェブツール(DDBJ Pipeline)との連携
- DDBJ PipelineでPlatanusを実行
- k-mer解析(k個の連続塩基に基づく各種解析)の基礎
 - 塩基ごとの出現頻度解析(k=1)、一気に計算
 - 2連続塩基の出現頻度解析(k=2)、基本スキルの復習や作図
- de novoアセンブリ時のエラー補正やゲノムサイズ推定の基本的な考え方
 - ランダムな塩基配列(仮想ゲノムおよび仮想NGSデータ)の作成
 - k-mer解析の応用、ゲノムサイズ推定の基礎
 - ゲノムサイズ推定(1,000 bpの仮想ゲノムの場合)
 - ゲノムサイズ推定(1,000 bpの仮想ゲノム; 4X → 10X coverageの場合)
 - k-mer出現頻度分布
 - シークエンスエラーを含む場合
 - 最終確認



DDBJ Pipeline

DDBJ Pipelineでは、主に①マッピングや②*de novo*アセンブリができる。特に後者ができるのは非常に有難い。③新規アカウント作成から*de novo*アセンブリまでの詳細については、乳酸菌連載第6回ウェブ資料を参照。ここでは、説明は必要最小限にして、Rのハンズオンへと移行する

The screenshot shows the DDBJ Pipeline website interface. At the top, there's a navigation bar with the DDBJ logo and the text "DDBJ Read Annotation Pipeline". Below this, there's a language selector for "English" and "Japanese". A description states: "DDBJ Read Annotation Pipelineは、次世代シーケンサ配列のクラウド型データ解析プラットフォームです。" The main section is titled "LOGIN" and includes fields for "User ID:" and "Password:" with a "Login" button. To the right of the login fields is a link for "新規アカウント作成" (New Account Creation). Below the login section, there's a "Pipelineフローチャート" (Pipeline Flowchart) which is a smaller version of the diagram shown in the red-bordered inset. At the bottom, there's a "Tweets" section with a tweet from @pipeline_info dated 22 Dec.



動作中JOBの確認
PipelineのIDをお持ちでない
ことができます。

マニュアルおよびチュートリアル

- 日本語チュートリアル
- 英語マニュアル
- DBCLS 統合TVチュートリアル
- DBCLS 統合TVチュートリアル
- チュートリアル: FTPでPipelineへ登録する方法
- チュートリアル: DDBJのアセンブリを行う方法

塩基配列・解析結果の登録

- DRA: NGS出力データの登録
- DDBJ-INSDC: Annotation

Citation

- Nagasaki, H. et al.,

DDBJ Pipelineのプログラム選択画面。
①Velvetや②Platanusを選択可能

DDBJ PipelineでPlatanus

The screenshot shows the 'Selecting Tools for Basic...' interface. On the left, there are navigation tabs for 'Preprocessing Start', 'step-1 Preprocessing', 'step-2 Workflow', 'JOB STATUS', and 'HELP'. The main area displays a table of tools for 'de novo Assembly' with a total limit of 22 Gbp. The 'Platanus' tool is highlighted in yellow and marked with a red arrow and the number '2'. The 'Velvet' tool is also marked with a red arrow and the number '1'. Below the table, there is a checkbox for 'Mapping Contigs by de novo Assemble to Reference Sequences.' and a small table showing 'BLAT' as a selected tool for 'Single-end analysis only'.

Tool	Help	Version	Base space	Color space	Paired-end	MSS (WGS)	Comment
<input type="checkbox"/> SOAPdenovo	Help	1.05	✓		✓		
<input type="checkbox"/> ABySS	Help	1.3.2	✓		✓		Maximum K-mer value is 64.
<input type="checkbox"/> Velvet	Help	1.2.10	✓		✓	✓	We severe recommend when performing Velvet, total length of those reads is up to 22G bp.Maximum K-mer value is 64.
<input type="checkbox"/> Trinity	Help	r2013-02-25	✓		✓		RNA-Seq De novo Assembly
<input checked="" type="checkbox"/> Platanus	Help	1.2.2	✓		✓		
<input type="checkbox"/> HGAP	Help	Protocol3 (v 2.2.0)					HGAP Pipeline for PacBio Sequence based on SMRT Analysis v2.2.0. For bax.h5 file only. (Beta version)

DDBJ Pipeline: Nagasaki et al., *DNA Res.*, **20**: 383-390, 2013

Platanus: Kajitani et al., *Genome Res.*, **24**: 1384-1395, 2014

DDBJ PipelineでPlatanus

de novoアセンブリの一般的な手順がわか
ていれば、赤枠内のStep1-3の説明の意
味がなんとなくわかる。①DDBJ Pipelineは
基本的にボタンをポチポチ押していただ

ACCOUNT
login ID [agribio]
Logout
Change password

ANALYSIS
Data setup
DRA Start
FTP upload
HTTP upload
DRA Import
Preprocessing Start
step-1
Preprocessing
Mapping /
de novo Assembly
step-2
Workflow
Genome (SNP/Short
Indel)
RNA-seq (Tag count)
ChIP-seq

JOB STATUS
step1.
Preprocessing
step1.
Mapping
step1.
de novo Assembly
step2-All status

HELP
HELP
TUTORIAL
Contact Us.
DDBJ Read Annotation

Setting for De Novo Assembly

platanus

Set optional parameters of the paired-end analysis

Memory Usage : Low (recommended) High

If you request "High" memory usage during the time Nig super computer system is congested, you might be kept waiting long before job starts running,

Step1) Assembly : Construct contigs using the algorithm based on the de Bruijn graph.

platanus assemble -t 15 -m 120 -o out [input] -f PE1.fastq PE2.fastq

Step2) Scaffold : Map paired-end (mate-pair) reads on contigs and construct scaffolds

platanus scaffold -t 8 -o out [input] -c out_contig.fa -b out_contigBubble.fa -IP1 PE1.fastq PE2.fastq (-OP2 MP1.fastq MP2.fastq)

Step3) Gap Close : Map paired-end (mate-pair) reads on scaffolds and assemble reads on gaps and close gaps

platanus gap_close -t 8 -o out [input] -c out_scaffold.fa -IP1 PE1.fastq PE2.fastq (-OP2 MP1.fastq MP2.fastq)

Step5) Create assembled sequences in FASTA file from pileupped reads to [submit WGS division of DDBJ.](#)

Set filtered length for contigs
 perl lengthfilter.pl pileupFile [input] out_WGS.txt

BACK NEXT

DDBJ PipelineでPlatanus

アセンブリ終了後の画面。①Platanus
実行結果ファイル(platanusResult.zip)
をダウンロードして解凍したのが…

ANALYSIS

- Data setup
 - DRA Start
 - FTP upload
 - HTTP upload
 - DRA Import
 - Preprocessing Start
- step-1
 - Preprocessing
 - Mapping / *de novo* Assembly
- step-2
 - Workflow**
 - Genome (SNP/Short Indel)
 - RNA-seq (Tag count)
 - ChIP-seq

JOB STATUS

- step1. **Preprocessing**
- step1. **Mapping**
- step1. ***de novo* Assembly**
- step2-All status

HELP

- HELP ☞
- TUTORIAL
- Contact Us.
 - DDBJ Read Annotation Pipeline.
 - Development Team.

Job info

ID: 21211

Tool (Version): Platanus (1.2.2)

RunAccession or Filename	Download	Read length	Alias
QC.1.trimmed.fastq.gz	QC.1.trimmed.fastq.gz	N.A. bp	L.hokkaidonensis_MiSeq_denovo

Download modified queries

- [QC.1.trimmed.fastq.gz](#) (Original size 189.4 MB)
- [QC.2.trimmed.fastq.gz](#) (Original size 189.6 MB)

Download wgs file

- [out_WGS.fasta.gz](#) (Original size 2.3 MB)

Assembly statistics

Contig # : 117
Total contig size : 2,356,019
Maximum contig size : 257,728
Minimum contig size : 101
N50 contig size : 92,304

Time

Wait time	Start time	End time
0: 0:9	2016-01-20 18:33:36	2016-01-21 10:10:06

Command	Start time	End time	Log1	Log2	Result	MD5
platanus assemble -m 120 -f QC.1.trimmed.fastq QC.2.trimmed.fastq	2016-01-20 18:33:37	2016-01-21 10:08:51		View	Download(2.2 MB)	MD5
platanus scaffold -c out_contig.fa -b out_contigBubble.fa -IP1 QC.1.trimmed.fastq QC.2.trimmed.fastq	2016-01-21 10:09:02	2016-01-21 10:09:12		View	Download(2.2 MB)	MD5
platanus gap_close -c out_scaffold.fa -IP1 QC.1.trimmed.fastq QC.2.trimmed.fastq	2016-01-21 10:09:23	2016-01-21 10:09:34		View	Download(2.2 MB)	MD5

①

BACK

Top of page

DDBJ PipelineでPlatanus

アセンブリ終了後の画面。①Platanus
実行結果ファイル(platanusResult.zip)
をダウンロードして解凍したのが…②
platanusResultというフォルダ

The screenshot displays the DDBJ Pipeline web interface for a job with ID 21211. The interface is divided into several sections: ANALYSIS, Job info, Assembly statistics, Time, and a command log. A file explorer window is overlaid on the right, showing the contents of the 'platanusResult' folder.

Job info

ID	21211		
Tool (Version)	Platanus (1.2.2)		
RunAccession or Filename	Download	Read length	Alias
QC.1.trimmed.fastq.gz	QC.1.trimmed.fastq.gz	N.A. bp	L.hokkaidonensis_MiSeq_denovo

Download modified queries

- [QC.1.trimmed.fastq.gz \(Original size 189.4 MB\)](#)
- [QC.2.trimmed.fastq.gz \(Original size 189.6 MB\)](#)

Download wgs file

- [out_WGS.fasta.gz \(Original size 2.3 MB\)](#)

Assembly statistics

Contig #	
Total contig size	: 2,356
Maximum contig size	: 257
Minimum contig size	: 10
N50 contig size	: 92

Time

Wait time	Start time	End time
0: 0:9	2016-01-20 18:33:36	2016-01-21 10:10:06

Command Log

Command	Start time	End time	Log1	Log2	Result	MD5
platanus assemble -m 120 -f QC.1.trimmed.fastq QC.2.trimmed.fastq	2016-01-20 18:33:37	2016-01-21 10:08:51		View	Download(2.2 MB)	MD5
platanus scaffold -c out_contig.fa -b out_contigBubble.fa -IP1 QC.1.trimmed.fastq QC.2.trimmed.fastq	2016-01-21 10:09:02	2016-01-21 10:09:12		View	Download(2.2 MB)	MD5
platanus gap_close -c out_scaffold.fa -IP1 QC.1.trimmed.fastq QC.2.trimmed.fastq	2016-01-21 10:09:23	2016-01-21 10:09:34		View	Download(2.2 MB)	MD5

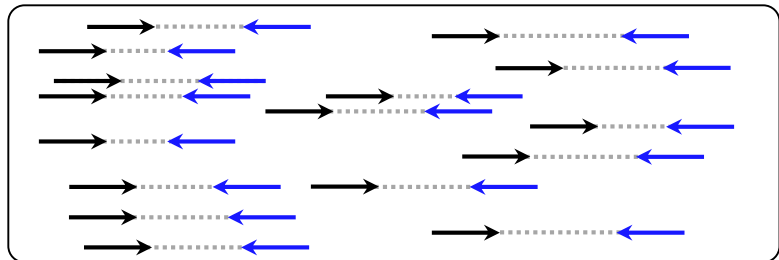
File Explorer (platanusResult)

名前	サイズ
out_32merFrq.tsv	12 KB
out_contig.fa	2,380 KB
out_contigBubble.fa	1 KB
out_gapClosed.fa	2,332 KB
out_lib1_insFreq.tsv	234 KB
out_scaffold.fa	2,334 KB
out_scaffoldBubble.fa	0 KB
out_scaffoldComponent.tsv	5 KB

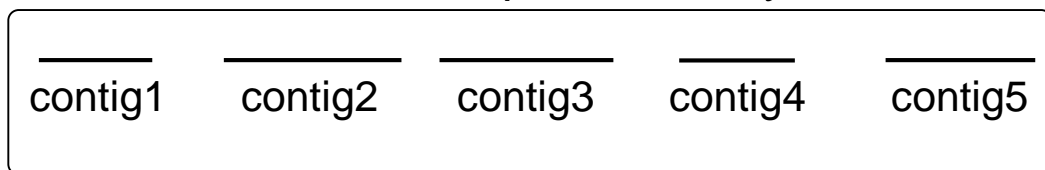
DDBJ PipelineでPlatanu

一般的な *de novo* アセンブリの手順を知っておけば、ファイル名から最終的な結果が③out_gapClosed.faだと認識できる

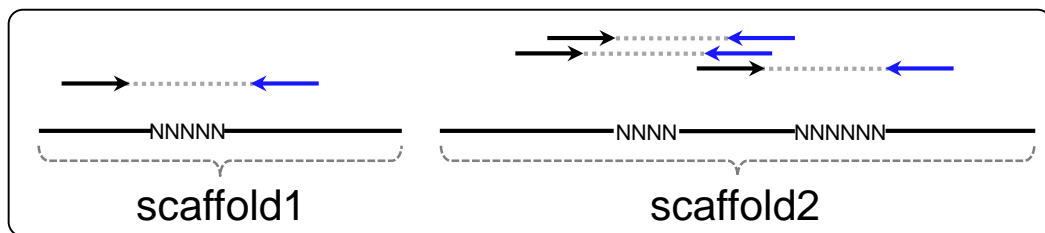
入力: paired-end FASTQファイル



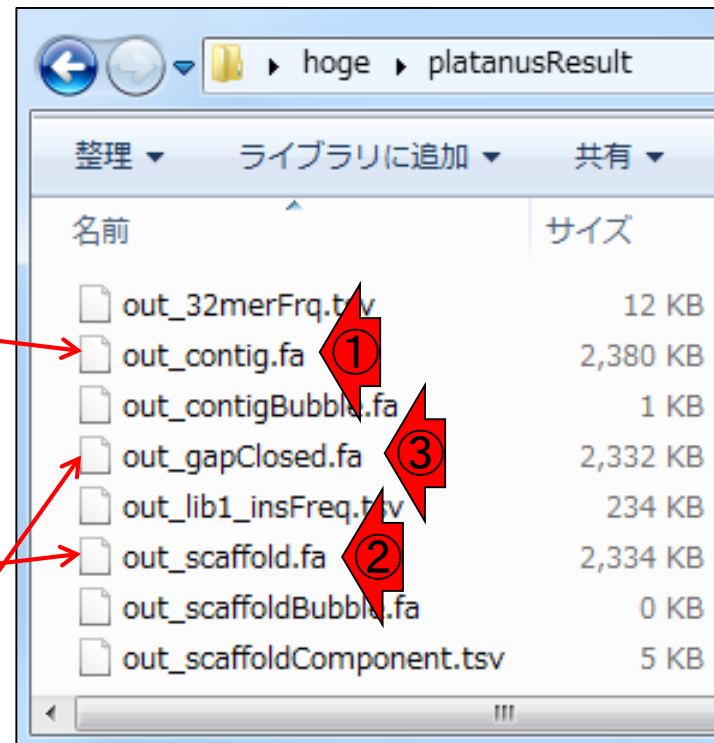
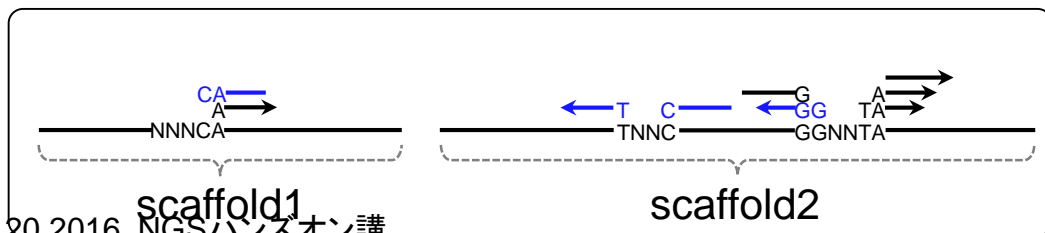
Step1: Assembly



Step2: Scaffold



Step3: Gap close



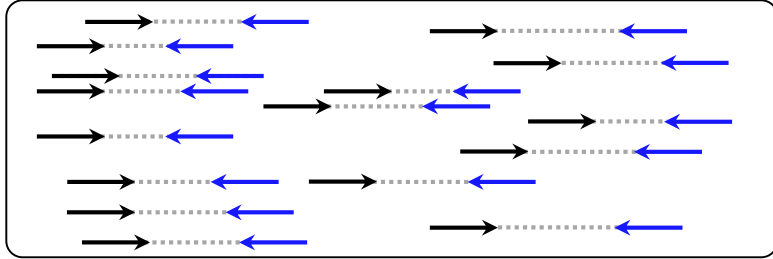
Contents

- NGS解析手段、ウェブツール(DDBJ Pipeline)との連携
- DDBJ PipelineでPlatanusを実行
- k-mer解析(k個の連続塩基に基づく各種解析)の基礎
 - 塩基ごとの出現頻度解析(k=1)、一気に計算
 - 2連続塩基の出現頻度解析(k=2)、基本スキルの復習や作図
- de novoアセンブリ時のエラー補正やゲノムサイズ推定の基本的な考え方
 - ランダムな塩基配列(仮想ゲノムおよび仮想NGSデータ)の作成
 - k-mer解析の応用、ゲノムサイズ推定の基礎
 - ゲノムサイズ推定(1,000 bpの仮想ゲノムの場合)
 - ゲノムサイズ推定(1,000 bpの仮想ゲノム; 4X → 10X coverageの場合)
 - k-mer出現頻度分布
 - シークエンスエラーを含む場合
 - 最終確認

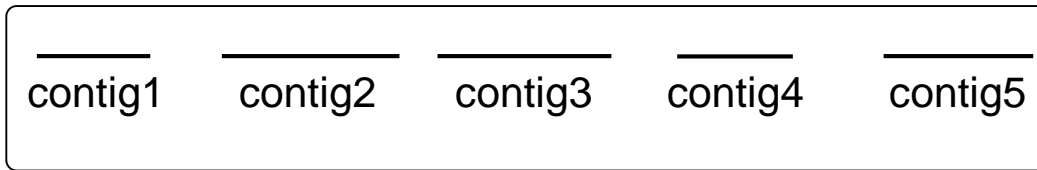


k-mer解析(k=1)

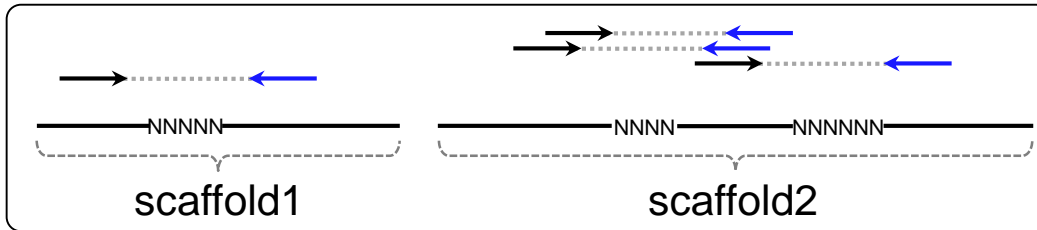
入力: paired-end FASTQファイル



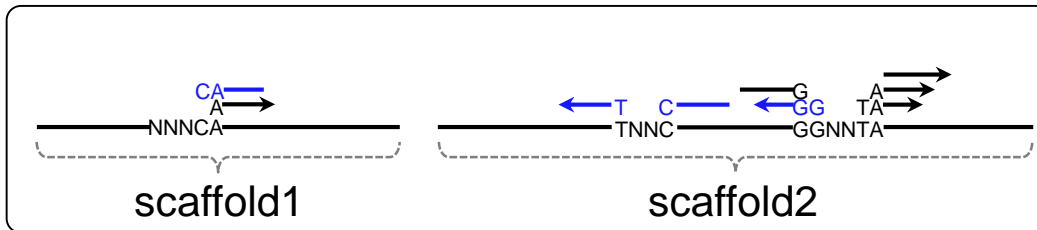
Step1: Assembly



Step2: Scaffold



Step3: Gap close



(アセンブリ実行結果の) multi-FASTAファイルを読み込んで、塩基ごとの出現頻度解析(k=1のk-mer解析に相当)ができる。①Step1実行後(out_contig.fa)はNがなく、②Step2実行後(out_scaffold.fa)にNができて、③Step3実行後(out_gapClosed.fa)にNが減るのだろうと妄想できる。それを自力で確認することで、アルゴリズムの理解を深められる

名前	サイズ
out_32merFrq.tsv	12 KB
out_contig.fa	2,380 KB
out_contigBubbles.fa	1 KB
out_gapClosed.fa	2,332 KB
out_lib1_insFreq.tsv	234 KB
out_scaffold.fa	2,334 KB
out_scaffoldBubbles.fa	0 KB
out_scaffoldComponent.tsv	5 KB

① (アセンブリ実行結果の) multi-FASTAファイルを読み込んで、塩基ごとの出現頻度解析を行う項目

k-mer解析(k=1)

(Rで)塩基配列解析

~NGS, RNA-seq, ゲノム, トランスクリプトーム, 正規化, 発現変動, 統計, モデル, バイオインフォマティクス~
(last modified 2016/02/03, since 2011)

What's new?

- このウェブページはフリーソフト Rと必要 [利用法\(Windows20書籍\)](#) があります。(2016/02/03)
- 多群間比較用の推奨については [門田のページ](#) についてはおおよそ検索結果のおおよそイン周辺の関連項目
- Erratum**, 2014.06.22
ここでは、kの値を小さくはいいです。(見てきたm(_)_m(2016/02/03))

- ・ [イントロ](#) | [一般](#) | [翻訳配列\(translate\)を取得\(応用\)](#) | [seqinr\(Charif 2005\)](#) (last modified 2015/03/09)
- ・ [イントロ](#) | [一般](#) | [相補鎖\(complement\)を取得](#) (last modified 2013/06/14)
- ・ [イントロ](#) | [一般](#) | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- ・ [イントロ](#) | [一般](#) | [逆鎖\(reverse\)を取得](#) (last modified 2013/06/14)
- ・ [イントロ](#) | [一般](#) | [k-mer解析 | k=1\(塩基ごとの出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/02/03) **NEW**
- ・ [イントロ](#) | [一般](#) | [k-mer解析 | k=2\(2連続塩基の出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/01/28) **NEW**
- ・ [イントロ](#) | [一般](#) | [k-mer解析 | k=3\(3連続塩基の出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/01/28) **NEW**
- ・ [イントロ](#) | [一般](#) | [k-mer解析 | k=n\(n連続塩基の出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/01/28) **NEW**
- ・ (削除予定) [イントロ](#) | [一般](#) | [2連続塩基の出現頻度情報を取得](#) (last modified 2015/04/20)
- ・ (削除予定) [イントロ](#) | [一般](#) | [3連続塩基の出現頻度情報を取得](#) (last modified 2015/02/19)
- ・ (削除予定) [イントロ](#) | [一般](#) | [任意の長さの連続塩基の出現頻度情報を取得](#) (last modified 2015/02/19)

イントロ | 一般 | k-mer解析 | k=1(塩基ごとの出現頻度解析) | Biostrings **NEW**

Biostringsパッケージを用いて、multi-FASTA形式ファイルを読み込んで、「A」、「C」、「G」、「T」、...、「N」、...など塩基ごとの出現頻度を調べるやり方を示します。k-mer解析のk=1の場合に相当します。

「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. [イントロ](#) | [一般](#) | [ランダムな塩基配列を作成](#)の4を実行して得られたmulti-FASTAファイル([hoge4.fa](#))の場合:
配列ごとに出現頻度をカウントした結果を返すやり方です。

```
in_f <- "hoge4.fa" #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.txt" #出力ファイル名を指定してout_fに格納

#必要なパッケージをロード
library(Biostrings) #パッケージの読み込み

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta") #in_fで指定したファイルの読み込み

#本番
out <- alphabetFrequency(fasta) #A,C,G,T,..の数を各配列ごとにカウントした結果をoutに格納
#outの中身を表示
```

①例題7が、PlatanusのStep3実行後のファイル(②out_gapClosed.fa)を入力とするものなので、そのままコピーできて便利。これを実行します

k-mer解析(k=1)

Biostringsパッケージを用いて、multi-FASTA形式ファイルを読み込んで、"A", "C", "G", "T", ..., "N", ...など塩基ごとの出現頻度を調べるやり方を示します。k-mer解析のk=1の場合に相当します。

「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. イントロ | 一般 | ランダムな塩基配列を作成の4.を実行して得られたmulti-FASTAファイル(hoge4.fa)の場合:

配列ごとに出現頻度をカウント



① 7. FASTA形式ファイル(out_gapClosed.fa)の場合:

[DDBJ Pipeline \(Nagasaki et al., DNA Res., 2013\)](#)上で de novoゲノムアセンブリプログラム [Platanus \(Kajitani et al., Genome Res., 2014\)](#) を実行して得られたmulti-FASTA形式ファイル(out_gapClosed.fa; 約2.4MB)です。



```

in_f <- "hoge4.fa"
out_f <- "hoge1.txt"

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNASTringSet

#本番
out <- alphabetFrequency(
out

#ファイルに保存
tmp <- cbind(names(fasta)
write.table(tmp, out_f, s

```

```

in_f <- "out_gapClosed.fa"
out_f <- "hoge7.txt"
param_base <- c("A", "C", "G", "T", "N")

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")

#本番
hoge <- alphabetFrequency(fasta)
obj <- is.element(colnames(hoge), param_base)
#out <- colSums(hoge[, obj])
out <- apply(as.matrix(hoge[, obj]), 2, sum)

#ファイルに保存
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=T, col.names=F)

```

k-mer解析(k=1)

イントロ | 一般 | k-mer解析 | k=1(塩基ごとの出現頻度解析) | Biostrings

つまり、Platanus実行結果ファイル(platanusResult.zip)をダウンロードし、解凍して得られた①platanusResultフォルダ中の②out_gapClosed.faを入力として、塩基ごとの出現頻度解析を行う

Biostringsパッケージを用いて、multi-FASTA形式ファイルを読み込んで、"A", "C", "G", "T", ..., "N", ...など塩基の出現頻度を調べるやり方を示します。k-mer解析のk=1の場合に相当します。

「ファイル」→「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. イントロ | 一般 | ランダムな塩基配列を作成の4.を実行して得られたmulti-FASTAファイル(hoge4.fa)の場合:

配列ごとに出現頻度をカウントした

```
in_f <- "hoge4.fa"
out_f <- "hoge1.txt"

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f)

#本番
out <- alphabetFrequency(fasta)

#ファイルに保存
tmp <- cbind(names(fasta), out)
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=T, col.names=F)
```

7. FASTA形式ファイル(out_gapClosed.fa)の場合:

[DDBJ Pipeline \(Nagasaki et al., DNA Res., 2013\)](#)上で de novoゲノムアセンブリを実行して得られたmulti-FASTA形式ファイル(out_gapClosed.fa)を入力として、塩基ごとの出現頻度解析を行う。

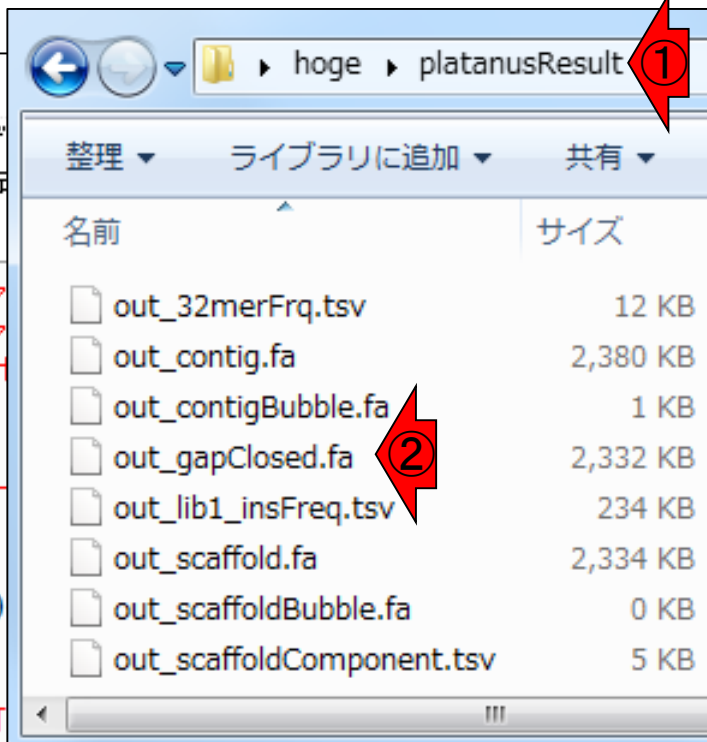
```
in_f <- "out_gapClosed.fa" #入力ファイル
out_f <- "hoge7.txt" #出力ファイル
param_base <- c("A", "C", "G", "T", "N") #出力する塩基

#必要なパッケージをロード
library(Biostrings) #パッケージをロード

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")

#本番
hoge <- alphabetFrequency(fasta) #A,C,G,T
obj <- is.element(colnames(hoge), param_base) #条件を満たすかどうかを判定した結果をobjに格納
#out <- colSums(hoge[, obj]) #列ごとの総和をoutに格納
out <- apply(as.matrix(hoge[, obj]), 2, sum) #列ごとの総和をoutに格納

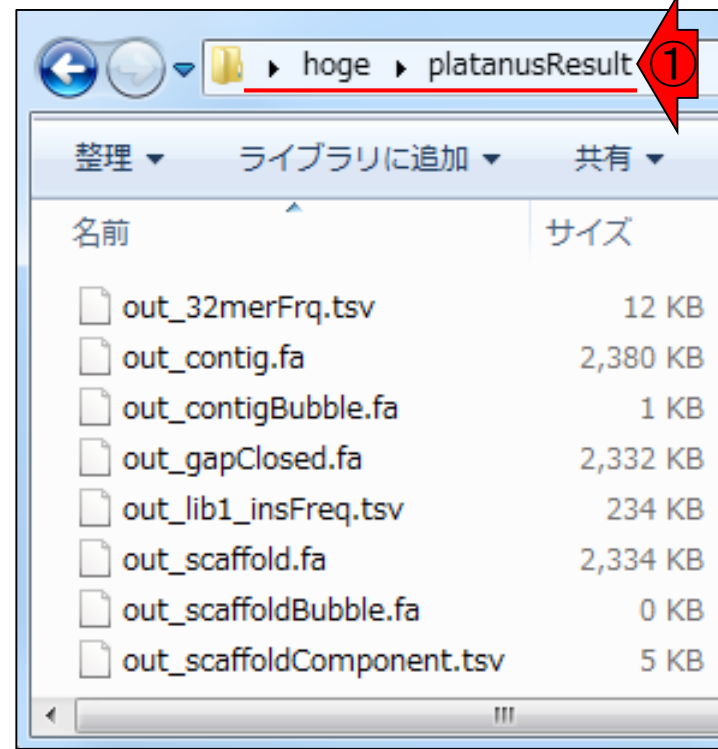
#ファイルに保存
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=T, col.names=F)
```



getwd()とlist.files()

```
R Console
'demo()' と入力すればデモをみることができます。
'help()' とすればオンラインヘルプが出ます。
'help.start()' で HTML ブラウザによるヘルプがみら$
'q()' と入力すれば R を終了します。

> getwd()
[1] "C:/Users/kadota/Desktop/hoge/platanusResult"
> list.files()
[1] "out_32merFrq.tsv"
[2] "out_contig.fa"
[3] "out_contigBubble.fa"
[4] "out_gapClosed.fa"
[5] "out_lib1_insFreq.tsv"
[6] "out_scaffold.fa"
[7] "out_scaffoldBubble.fa"
[8] "out_scaffoldComponent.tsv"
> |
```



コピペ実行結果

①解析結果(塩基ごとの出現頻度情報)はhoge7.txtに保存される。②list.files()で、自分が出力として指定した③hoge7.txtが存在することを確認。④EXCELで眺めるとこんな感じ

7. FASTA形式ファイル(out_gapClosed.fa)の場合:

DDBJ Pipeline (Nagasaki et al., DNA Res., 2013)上で de novoゲノムアセンブリプログラムPlatanus (Kajitani et al., Genome Res., 2014)を実行して得られたmulti-FASTA形式ファイル(out_gapClosed.fa; 約2.4MB)です。

```

in_f <- "out_gapClosed.fa" #入力ファイル名を指定してin_fに格納
out_f <- "hoge7.txt" #出力ファイル名を指定してout_fに格納
param_base <- c("A", "C", "G", "T", "N") #出力

#必要なパッケージをロード
library(Biostrings) #パッケージ

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")

#本番
hoge <- alphabetFrequency(fasta) #A,C,G,T,N
obj <- is.element(colnames(hoge), param_base) #A,C,G,T,N
#out <- colSums(hoge[, obj]) #列ごとの出現頻度
out <- apply(as.matrix(hoge[, obj]), 2, sum)

#ファイルに保存
write.table(out, out_f, sep="\t", append=F, q

```

```

R Console
> hoge <- alphabetFrequency(fasta)
> obj <- is.element(colnames(hoge), param_base)
> #out <- colSums(hoge[, obj])
> out <- apply(as.matrix(hoge[, obj]), 2, sum)
>
> #ファイルに保存
> write.table(out, out_f, sep="\t", append=F, q
> list.files()
[1] "hoge7.txt"
[2] "out_32merFrq.tsv"
[3] "out_contig.fa"
[4] "out_contigBubble.fa"
[5] "out_gapClosed.fa"
[6] "out_lib1_insFreq.tsv"
[7] "out_scaffold.fa"
[8] "out_scaffoldBubble.fa"
[9] "out_scaffoldComponent.tsv"
> |

```

④

	A	B
1	A	729772
2	C	458211
3	G	440585
4	T	727451
5	N	0

総塩基数

①outオブジェクトは数値ベクトル。②和を計算するsum関数実行結果(2,356,019)は、入力ファイル(out_gapClosed.fa)中の総塩基数に相当

7. FASTA形式ファイル(out_gapClosed.fa)の場合:

DDBJ Pipeline (Nagasaki et al., DNA Res., 2013)上で de novoゲノムアセンブリプログラムPlatanus (Kajitani et al., Genome Res., 2014)を実行して得られたmulti-FASTA形式ファイル(out_gapClosed.fa; 約2.4MB)です。

```

in_f <- "out_gapClosed.fa"      #入力ファイル名を指定してin_fに格納
out_f <- "hoge7.txt"           #出力ファイル名を指定してout_fに格納
param_base <- c("A", "C", "G", "T", "N")#出力

```

```

#必要なパッケージをロード
library(Biostrings)           #バック

```

```

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")

```

```

#本番
hoge <- alphabetFrequency(fasta) #A,C,G
obj <- is.element(colnames(hoge), param_base) #列ごと
#out <- colSums(hoge[, obj])
out <- apply(as.matrix(hoge[, obj]), 2, sum)#

```

```

#ファイルに保存
write.table(out, out_f, sep="\t", append=F,

```

```

R Console
[4] "out_contigBubble.fa"
[5] "out_gapClosed.fa"
[6] "out_lib1_insFreq.tsv"
[7] "out_scaffold.fa"
[8] "out_scaffoldBubble.fa"
[9] "out_scaffoldComponent.tsv"
> getwd()
[1] "C:/Users/kadota/Desktop/hog
> date()
[1] "Thu Feb 04 14:27:38 2016"
> out_f
[1] "hoge7.txt"
> out
      A      C      G      T      N
729772 458211 440585 727451      0
> sum(out)
[1] 2356019
> |

```

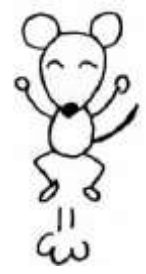
	A	B
1	A	729772
2	C	458211
3	G	440585
4	T	727451
5	N	0



総塩基数

①DDBJ Pipeline実行結果画面上の数値と同じ。
 ②入力ファイル(out_gapClosed.fa)は、DDBJ Pipeline上でPlatanusという*de novo*アセンブリプログラムを実行した結果だったことを思い出そう

ID		21211				
Tool (Version)		Platanus (1.2.2)				
RunAccession or Filename	Download	Read length	Alias			
QC.1.trimmed.fastq.gz	QC.1.trimmed.fastq.gz	N.A. bp	L.hokkaidonensis_MiSeq_denovo			
Download modified queries						
<ul style="list-style-type: none"> • QC.1.trimmed.fastq.gz (Original size 189.4 MB) • QC.2.trimmed.fastq.gz (Original size 189.6 MB) 						
Download wgs file						
<ul style="list-style-type: none"> • out_WGS.fasta.gz (Original size 2.3 MB) 						
Assembly statistics						
		Contig # : 117 Total contig size : 2,356,019 Maximum contig size : 257,728 Minimum contig size : 101 N50 contig size : 92,304				
Time						
Wait time	Start time	End time				
0: 0:9	2016-01-20 18:33:36	2016-01-21 10:10:06				
Command	Start time	End time	Log1	Log2	Result	MD5
platanus assemble -m 120 -f QC.1.trimmed.fastq QC.2.trimmed.fastq	2016-01-20 18:33:37	2016-01-21 10:08:51		View	Download(2.2 MB)	MD5
platanus scaffold -c out_contig.fa -b out_contigBubble.fa -IP1 QC.1.trimmed.fastq QC.2.trimmed.fastq	2016-01-21 10:09:02	2016-01-21 10:09:12		View	Download(2.2 MB)	MD5
platanus gap_close -c out_scaffold.fa -IP1 QC.1.trimmed.fastq QC.2.trimmed.fastq	2016-01-21 10:09:23	2016-01-21 10:09:34		View	Download(2.2 MB)	MD5



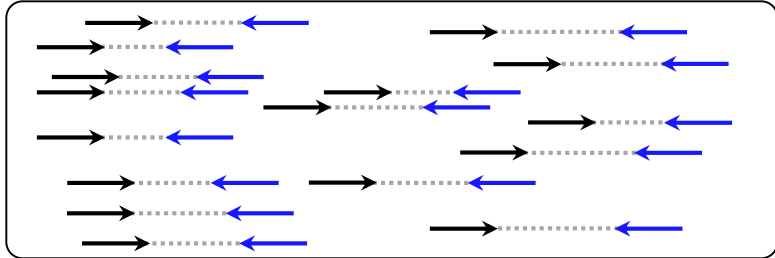
Contents

- NGS解析手段、ウェブツール(DDBJ Pipeline)との連携
- DDBJ PipelineでPlatanusを実行
- k-mer解析(k個の連続塩基に基づく各種解析)の基礎
 - 塩基ごとの出現頻度解析(k=1)、一気に計算
 - 2連続塩基の出現頻度解析(k=2)、基本スキルの復習や作図
- de novoアセンブリ時のエラー補正やゲノムサイズ推定の基本的な考え方
 - ランダムな塩基配列(仮想ゲノムおよび仮想NGSデータ)の作成
 - k-mer解析の応用、ゲノムサイズ推定の基礎
 - ゲノムサイズ推定(1,000 bpの仮想ゲノムの場合)
 - ゲノムサイズ推定(1,000 bpの仮想ゲノム; 4X → 10X coverageの場合)
 - k-mer出現頻度分布
 - シークエンスエラーを含む場合
 - 最終確認

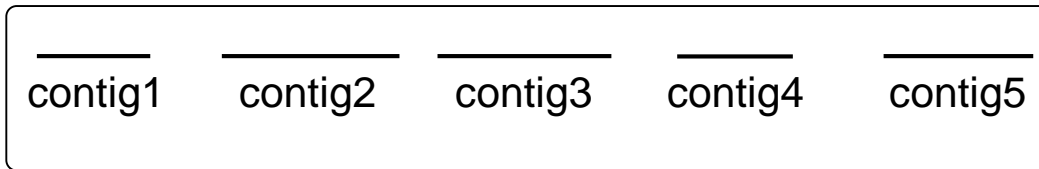


目的をおさらい

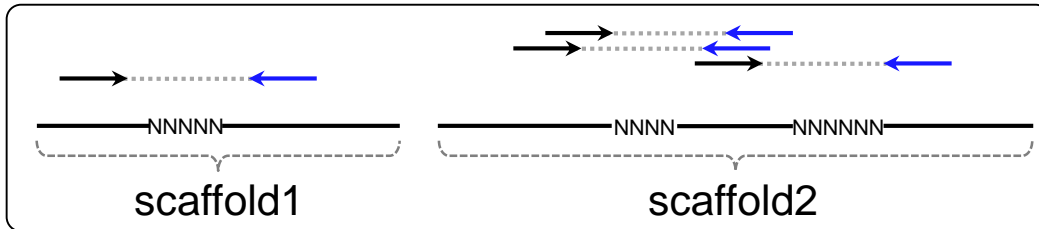
入力: paired-end FASTQファイル



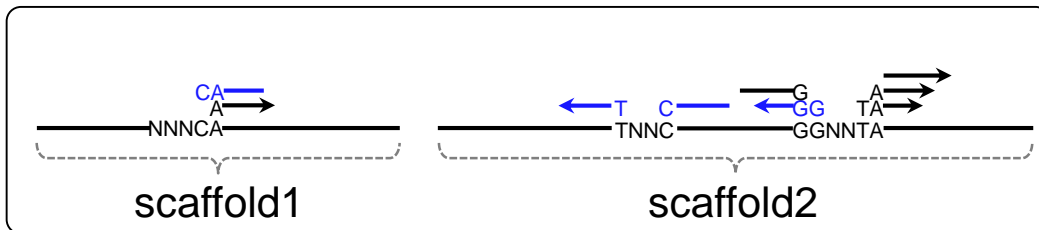
Step1: Assembly



Step2: Scaffold



Step3: Gap close

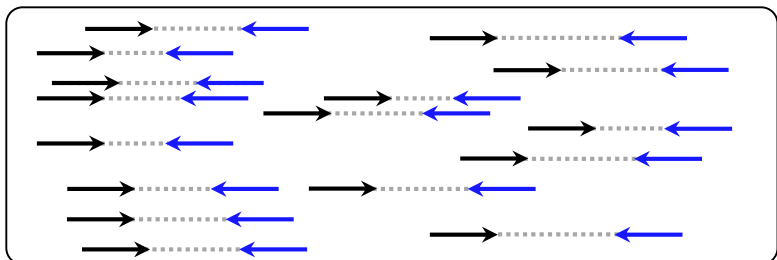


(アセンブリ実行結果の) multi-FASTAファイルを読み込んで、塩基ごとの出現頻度解析ができる。①Step1実行後 (out_contig.fa) はNがなく、②Step2実行後 (out_scaffold.fa) にNができて、③Step3実行後 (out_gapClosed.fa) にNが減るので、それを自力で確認することで、アルゴリズムの理解を深めることができる

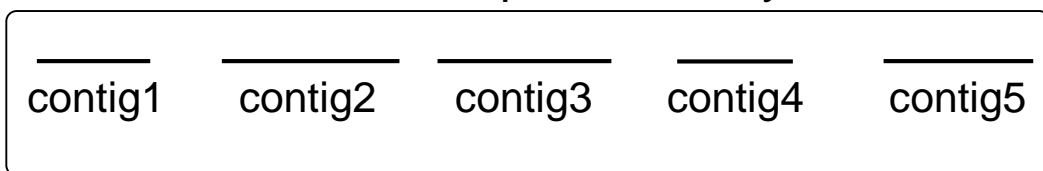
名前	サイズ
out_32merFrq.tsv	12 KB
out_contig.fa	2,380 KB
out_contigBubbles.fa	1 KB
out_gapClosed.fa	2,332 KB
out_lib1_insFreq.tsv	234 KB
out_scaffold.fa	2,334 KB
out_scaffoldBubbles.fa	0 KB
out_scaffoldComponent.tsv	5 KB

一気に結果を得る

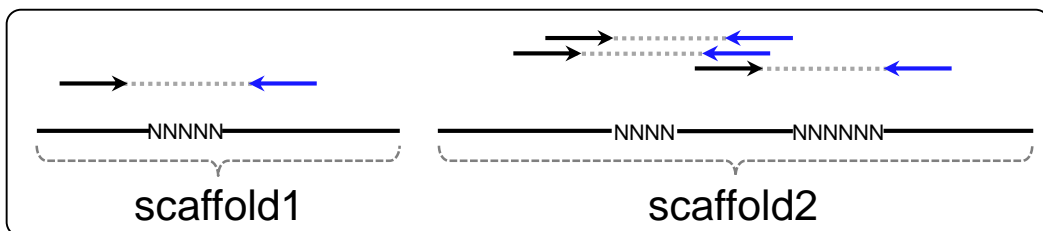
入力: paired-end FASTQファイル



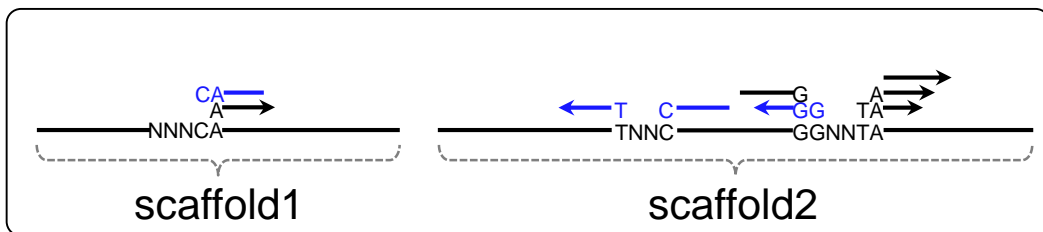
Step1: Assembly



Step2: Scaffold



Step3: Gap close



- ① Step1実行後 (out_contig.fa) はNがなく、
- ② Step2実行後 (out_scaffold.fa) にNができて、
- ③ Step3実行後 (out_gapClosed.fa) にNが減るのだろうという妄想を一気に確認したい!

名前	サイズ
out_32merFrq.tsv	12 KB
out_contig.fa	2,380 KB
out_contigBubbles.fa	1 KB
out_gapClosed.fa	2,332 KB
out_lib1_insFreq.tsv	234 KB
out_scaffold.fa	2,334 KB
out_scaffoldBubbles.fa	0 KB
out_scaffoldComponent.tsv	5 KB

実際の利用時は、hogeフォルダ直下にある①
rcode1.txtのように、無駄なコメントを除いてスリムに
した一連のスクリプトを作成しておき、一気にコピー

一気に結果を得る

```
rcode1.txt
param_base <- c("A", "C", "G", "T", "N") #出力させたい塩基を指定↓
library(Biostrings) #パッケージの読み込み↓
#####↓
### Step 1↓
#####↓
in_f <- "out_contig.fa" #入力ファイル名を指定してin_fに格納↓
out_f <- "result_step1.txt" #出力ファイル名を指定してout_fに格納↓
↓
fasta <- readDNASTringSet(in_f, format="fasta") #in_fで指定したファイルの読み込み
hoge <- alphabetFrequency(fasta) #A,C,G,T,..の数を各配列ごとにカウントした
obj <- is.element(colnames(hoge), param_base) #条件を満たすかどうかを判定した結果
out <- apply(as.matrix(hoge[, obj]), 2, sum) #列ごとの総和をoutに格納↓
write.table(out, out_f, sep="¥t", append=F, quote=F, row.names=T, col.names=F) #t
↓
#####↓
### Step 2↓
#####↓
in_f <- "out_scaffold.fa" #入力ファイル名を指定してin_fに格納↓
out_f <- "result_step2.txt" #出力ファイル名を指定してout_fに格納↓
↓
fasta <- readDNASTringSet(in_f, format="fasta") #in_fで指定したファイルの読み込み
hoge <- alphabetFrequency(fasta) #A,C,G,T,..の数を各配列ごとにカウントした
obj <- is.element(colnames(hoge), param_base) #条件を満たすかどうかを判定した結果
out <- apply(as.matrix(hoge[, obj]), 2, sum) #列ごとの総和をoutに格納↓
write.table(out, out_f, sep="¥t", append=F, quote=F, row.names=T, col.names=F) #t
↓
#####↓
### Step 3↓
#####↓
```

一気に結果を得る

hogeフォルダ直下にある、①rcode1.txtのように無駄なコメントを除いてスリムにした一連の скриプトを作成しておき、一気にコピペ。②コピペ後に自分が指定した出力ファイルができていることを確認

```
rcode1.txt
param_base <- c("A", "C", "G", "T", "N") #出力させたい塩基を指定↓
library(Biostrings) #パッケージの読み込み↓
#####↓
### Step 1↓
#####↓
in_f <- "out_contig.fa" #入力ファイル名を指定してin_fに格納↓
out_f <- "result_step1.txt" #出力ファイル名を指定してout_fに格納↓
↓
fasta <- readDNASTringSet(in_f, format="fasta") #in
hoge <- alphabetFrequency(fasta) #A,C,G,T,...
obj <- is.element(colnames(hoge), param_base) #条件
out <- apply(as.matrix(hoge[, obj]), 2, sum) #列ごと
write.table(out, out_f, sep="¥t", append=F, quote=
↓
#####↓
### Step 2↓
#####↓
in_f <- "out_scaffold.fa" #入力ファイル
out_f <- "result_step2.txt" #出力ファイル
↓
fasta <- readDNASTringSet(in_f, format="fasta") #in
hoge <- alphabetFrequency(fasta) #A,C,G,T,...
obj <- is.element(colnames(hoge), param_base) #条件
out <- apply(as.matrix(hoge[, obj]), 2, sum) #列ごと
write.table(out, out_f, sep="¥t", append=F, quote=
↓
#####↓
### Step 3↓
#####↓
```

```
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/platanusResult"
> list.files()
[1] "hoge7.txt"
[2] "out_32merFrq.tsv"
[3] "out_contig.fa"
[4] "out_contigBubble.fa"
[5] "out_gapClosed.fa"
[6] "out_lib1_insFreq.tsv"
[7] "out_scaffold.fa"
[8] "out_scaffoldBubble.fa"
[9] "out_scaffoldComponent.tsv"
[10] "result_step1.txt"
[11] "result_step2.txt"
[12] "result_step3.txt" }
```

結果のまとめ

```

rcode1.txt x
param_base <- c("A", "C", "G", "T", "N")#出力させたい塩基を指定↓
library(Biostrings) #パッケージの読み込み↓
#####↓
### Step 1↓
#####↓
in_f <- "out_contig.fa" #入力ファイル名を指定してin_fに格納↓
out_f <- "result_step1.txt" #出力ファイル名を指定してout_fに格納↓
↓
fasta <- readDNASTringSet(in_f, format="fasta")#in
hoge <- alphabetFrequency(fasta) #A,C,G,T,..
obj <- is.element(colnames(hoge), param_base)#条件
out <- apply(as.matrix(hoge[, obj]), 2, sum)#列ごと
write.table(out, out_f, sep="¥t", append=F, quote=
↓
#####↓
### Step 2↓
#####↓
in_f <- "out_scaffold.fa" #入力ファイ
out_f <- "result_step2.txt" #出力ファイ
↓
fasta <- readDNASTringSet(in_f, format="fasta")#in
hoge <- alphabetFrequency(fasta) #A,C,G,T,..
obj <- is.element(colnames(hoge), param_base)#条件
out <- apply(as.matrix(hoge[, obj]), 2, sum)#列ごと
write.table(out, out_f, sep="¥t", append=F, quote=
↓
#####↓
### Step 3↓
#####↓

```

```

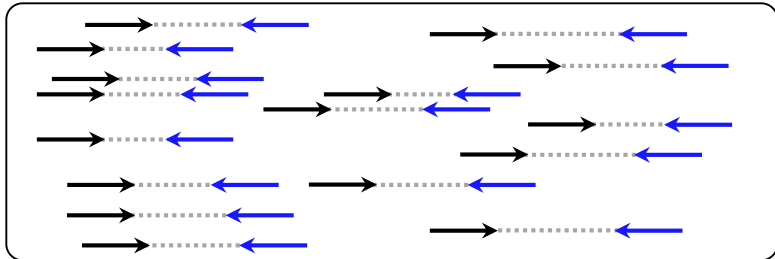
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/platanusResult"
> list.files()
[1] "hoge7.txt"
[2] "out_32merFrq.tsv"
[3] "out_contig.fa"
[4] "out_contigBubbler"
[5] "out_gapClosed"
[6] "out_lib1_insFle"
[7] "out_scaffold.fa"
[8] "out_scaffoldBubb"
[9] "out_scaffoldComponent.tsv"
[10] "result_step1.txt"
[11] "result_step2.txt"
[12] "result_step3.txt"
> |

```

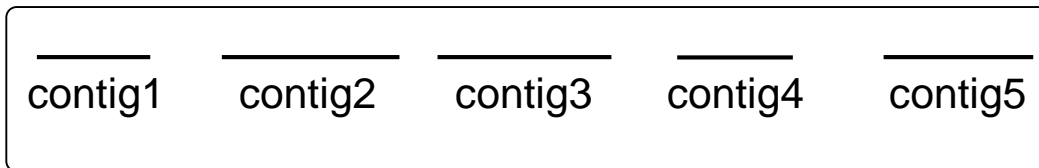
base	Step1	Step2	Step3
A	739836	729635	729772
C	469377	458119	458211
G	446806	440510	440585
T	742714	727306	727451
N	0	491	0

結果の解釈

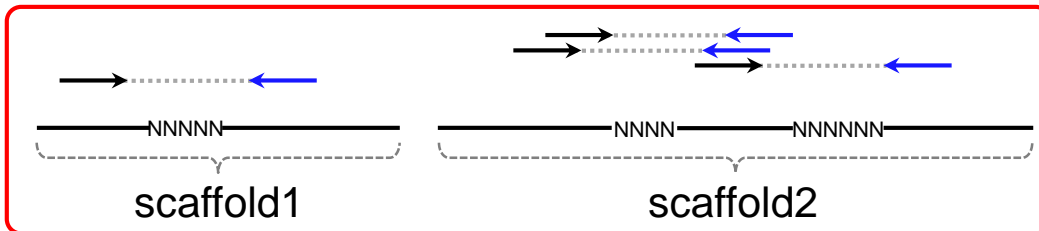
入力: paired-end FASTQファイル



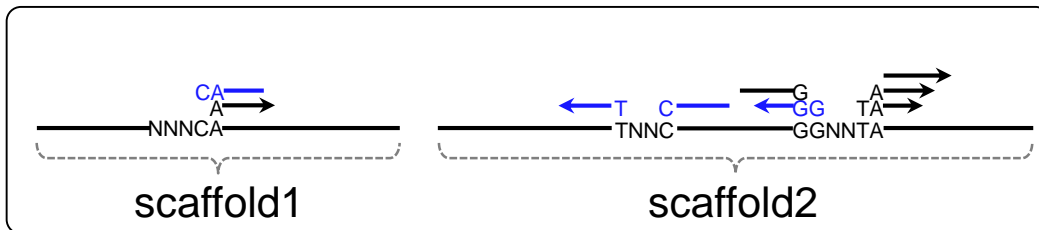
Step1: Assembly



Step2: Scaffold



Step3: Gap close



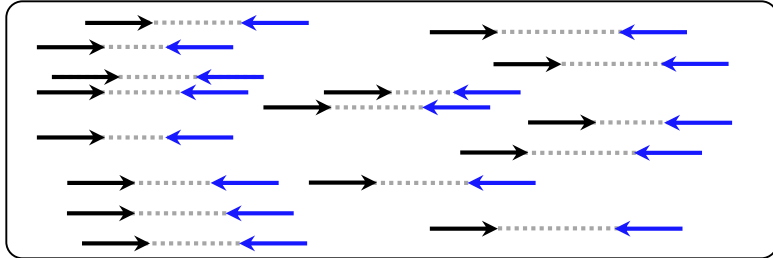
①Step1実行後はNが0。②Step2実行後にNが491個生成されたということは、いくつかのcontigsがまとめられてscaffoldsになったのだろう。③Step3でNが0個になったのはおそらくたまたまうまくいっただけ。491個よりも減ったということが重要で、gap closeがうまく機能したと判断できる



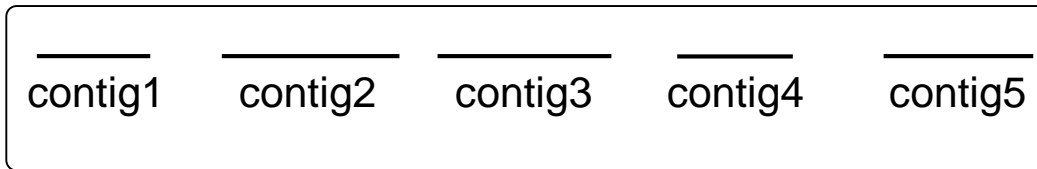
base	Step1	Step2	Step3
A	739836	729635	729772
C	469377	458119	458211
G	446806	440510	440585
T	742714	727306	727451
N	0	491	0

配列数

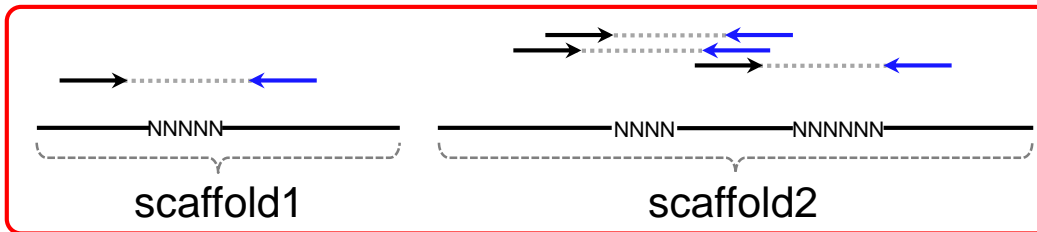
入力: paired-end FASTQファイル



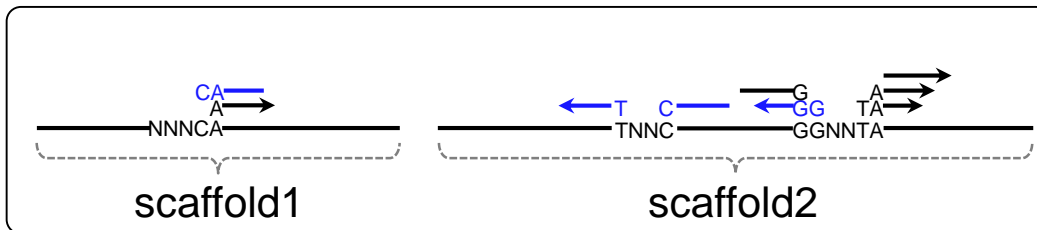
Step1: Assembly



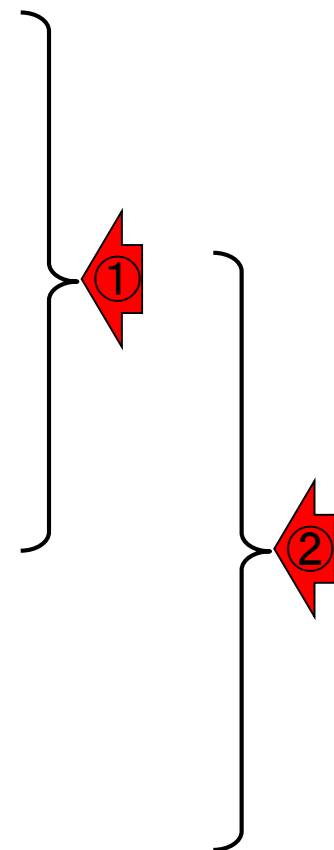
Step2: Scaffold



Step3: Gap close



Rを使うことで、アセンブリプログラムの内部挙動の把握や理解ができる。他の例は、配列数 (contig数やscaffold数と書くと説明しづらいので配列数に統一)。配列数は、①Step1 → Step2で減り、②Step2 → Step3では不変だろうと予想



配列数

配列数は、①Step1 → Step2で減り、②Step2 → Step3では不変だろうと予想。③(hogeフォルダ直下の)rcode2.txtは、配列数をカウントする必要最小限のコード。349 → 117 → 117で予想通りの結果

```
rcode2.txt ③
library(Biostrings) #パッケージの読み込み↓
### Step 1 ###↓
in_f <- "out_contig.fa" #入力ファイル
fasta <- readDNAStringSet(in_f, format="fasta")#ir
length(fasta) #配列数を表
### Step 2 ###↓
in_f <- "out_scaffold.fa" #入力ファイル
fasta <- readDNAStringSet(in_f, format="fasta")#ir
length(fasta) #配列数を表
### Step 3 ###↓
in_f <- "out_gapClosed.fa" #入力ファイル
fasta <- readDNAStringSet(in_f, format="fasta")#ir
length(fasta) #配列数を表
```

```
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/platanusResult"
> list.files(pattern="*.fa")
[1] "out_contig.fa" "out_contigBubble.fa"
[3] "out_gapClosed.fa" "out_scaffold.fa"
[5] "out_scaffoldBubble.fa"
> ### Step 1 ###
> in_f <- "out_contig.fa" #入力ファイル$
> fasta <- readDNAStringSet(in_f, format="fasta")
> length(fasta) #配列数$
[1] 349
>
> ### Step 2 ###
> in_f <- "out_scaffold.fa" #入力ファイル$
> fasta <- readDNAStringSet(in_f, format="fasta")
> length(fasta) #配列数$
[1] 117
>
> ### Step 3 ###
> in_f <- "out_gapClosed.fa" #入力ファイル$
> fasta <- readDNAStringSet(in_f, format="fasta")
> length(fasta) #配列数$
[1] 117
```

配列数

①DDBJ Pipeline最終結果の数値(117個)と同じことがわかります。ここまでの話は、2016.03.03の講習会資料のダイジェスト版でした。「2016.03.03」でページ内検索

ID											
21211											
Tool (Version)											
Platanus (1.2.2)											
RunAccession or Filename	Download										
QC.1.trimmed.fastq.gz	QC.1.trimmed.fastq.gz										
Read length	Alias										
N.A. bp	L.hokkaidonensis_MiSeq_denovo										
Download modified queries											
<ul style="list-style-type: none"> • QC.1.trimmed.fastq.gz (Original size 189.4 MB) • QC.2.trimmed.fastq.gz (Original size 189.6 MB) 											
Download wgs file											
<ul style="list-style-type: none"> • out_WGS.fasta.gz (Original size 2.3 MB) 											
Assembly statistics											
<table border="0"> <tr> <td>Contig #</td> <td>: 117</td> </tr> <tr> <td>Total contig size</td> <td>: 2,356,019</td> </tr> <tr> <td>Maximum contig size</td> <td>: 257,728</td> </tr> <tr> <td>Minimum contig size</td> <td>: 101</td> </tr> <tr> <td>N50 contig size</td> <td>: 92,304</td> </tr> </table>		Contig #	: 117	Total contig size	: 2,356,019	Maximum contig size	: 257,728	Minimum contig size	: 101	N50 contig size	: 92,304
Contig #	: 117										
Total contig size	: 2,356,019										
Maximum contig size	: 257,728										
Minimum contig size	: 101										
N50 contig size	: 92,304										
Time											
Wait time	Start time										
0: 0:9	2016-01-20 18:33:36										
	End time										
	2016-01-21 10:10:06										
Command	Start time	End time	Log1	Log2	Result	MD5					
platanus assemble -m 120 -f QC.1.trimmed.fastq QC.2.trimmed.fastq	2016-01-20 18:33:37	2016-01-21 10:08:51		View	Download(2.2 MB)	MD5					
platanus scaffold -c out_contig.fa -b out_contigBubble.fa -IP1 QC.1.trimmed.fastq QC.2.trimmed.fastq	2016-01-21 10:09:02	2016-01-21 10:09:12		View	Download(2.2 MB)	MD5					
platanus gap_close -c out_scaffold.fa -IP1 QC.1.trimmed.fastq QC.2.trimmed.fastq	2016-01-21 10:09:23	2016-01-21 10:09:34		View	Download(2.2 MB)	MD5					



Contents

- NGS解析手段、ウェブツール(DDBJ Pipeline)との連携
- DDBJ PipelineでPlatanusを実行
- k-mer解析(k個の連続塩基に基づく各種解析)の基礎
 - 塩基ごとの出現頻度解析(k=1)、一気に計算
 - 2連続塩基の出現頻度解析(k=2)、基本スキルの復習や作図
- de novoアセンブリ時のエラー補正やゲノムサイズ推定の基本的な考え方
 - ランダムな塩基配列(仮想ゲノムおよび仮想NGSデータ)の作成
 - k-mer解析の応用、ゲノムサイズ推定の基礎
 - ゲノムサイズ推定(1,000 bpの仮想ゲノムの場合)
 - ゲノムサイズ推定(1,000 bpの仮想ゲノム; 4X → 10X coverageの場合)
 - k-mer出現頻度分布
 - シークエンスエラーを含む場合
 - 最終確認



ヒトゲノムを調べると、①CGという連続塩基の出現確率が他(特にCC, GC, GG)に比べて少ないという有名な話を調べる。②染色体番号。実際に行うことは、2連続塩基の出現頻度解析(k=2のときのk-mer解析に相当)。2016.04.25の講義資料(スライド35~)をベースに作成

k-mer解析(k=2)

②

①

	AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	9.5%	5.0%	7.1%	7.4%	7.3%	5.4%	1.0%	7.1%	6.0%	4.4%	5.4%	5.0%	6.3%	6.0%	7.3%	9.6%
2	10.0%	5.0%	7.0%	7.9%	7.2%	5.0%	0.9%	7.0%	5.9%	4.1%	5.0%	5.0%	6.7%	5.9%	7.2%	10.0%
3	10.1%	5.0%	6.9%	8.0%	7.2%	4.9%	0.8%	6.9%	5.9%	4.0%	4.9%	5.0%	6.9%	5.9%	7.2%	10.2%
4	10.6%	5.0%	6.7%	8.5%	7.1%	4.5%	0.8%	6.7%	5.9%	3.8%	4.5%	5.0%	7.3%	5.8%	7.1%	10.6%
5	10.2%	5.0%	6.9%	8.1%	7.2%	4.8%	0.9%	6.9%	5.9%	4.0%	4.8%	5.1%	6.9%	5.9%	7.2%	10.3%
6	10.2%	5.0%	6.9%	8.1%	7.2%	4.8%	0.9%	6.9%	5.9%	4.0%	4.9%	5.0%	6.9%	5.9%	7.2%	10.2%
7	9.8%	5.0%	7.0%	7.7%	7.3%	5.1%	1.0%	7.0%	6.0%	4.2%	5.1%	5.1%	6.5%	5.9%	7.3%	10.0%
8	10.0%	5.1%	6.9%	7.9%	7.2%	5.0%	0.9%	6.9%	6.0%	4.1%	5.0%	5.0%	6.7%	5.9%	7.2%	10.0%
9	9.7%	5.1%	7.0%	7.6%	7.3%	5.3%	1.0%	7.0%	6.0%	4.3%	5.3%	5.0%	6.4%	6.0%	7.3%	9.7%
10	9.6%	5.0%	7.1%	7.5%	7.3%	5.3%	1.0%	7.1%	6.0%	4.4%	5.3%	5.1%	6.3%	6.0%	7.4%	9.7%
11	9.5%	5.1%	7.1%	7.5%	7.3%	5.3%	1.0%	7.1%	6.1%	4.3%	5.4%	5.0%	6.3%	6.0%	7.3%	9.6%
12	9.8%	5.0%	7.0%	7.7%	7.2%	5.1%	1.0%	7.0%	6.0%	4.2%	5.2%	5.1%	6.6%	6.0%	7.2%	9.9%
13	10.5%	5.0%	6.8%	8.4%	7.1%	4.5%	0.9%	6.7%	5.9%	3.8%	4.6%	5.0%	7.2%	5.8%	7.1%	10.6%
14	9.7%	5.0%	7.0%	7.7%	7.2%	5.1%	1.0%	7.0%	6.0%	4.2%	5.2%	5.1%	6.6%	5.9%	7.3%	9.9%
15	9.4%	5.1%	7.1%	7.3%	7.3%	5.4%	1.1%	7.1%	6.0%	4.5%	5.5%	5.1%	6.1%	6.0%	7.4%	9.5%
16	8.6%	5.1%	7.3%	6.7%	7.5%	6.1%	1.4%	7.2%	6.1%	5.0%	6.1%	5.1%	5.4%	6.1%	7.6%	8.8%
17	8.5%	5.1%	7.3%	6.4%	7.4%	6.3%	1.5%	7.4%	6.2%	5.1%	6.4%	5.0%	5.2%	6.1%	7.5%	8.6%
18	10.1%	5.1%	7.0%	7.7%	7.3%	5.3%	1.0%	7.1%	6.0%	4.4%	5.3%	5.1%	6.3%	6.0%	7.4%	9.7%

Lander et al., *Nature*, 409: 860-921, 2001

k-mer解析(k=2)

②平成27年度NGSハンズオン講習会(2015.07.29の slides 94~)で用いたhoge4.faを入力として、配列ごと(この場合コンティグごと)に16種類の2連続塩基の出現頻度解析の概念を説明します

- ・ イントロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- ・ イントロ | 一般 | [逆鎖\(reverse\)を取得](#) (last modified 2013/06/14)
- ・ イントロ | 一般 | [k-mer解析 | k=1\(塩基ごとの出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/04/27) **NEW**
- ・ イントロ | 一般 | [k-mer解析 | k=2\(2連続塩基の出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/01/28)
- ・ イントロ | 一般 | [k-mer解析 | k=3\(3連続塩基の出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/01/28)
- ・ イントロ | 一般 | [k-mer解析 | k=n\(n連続塩基の出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/01/28)



イントロ | 一般 | k-mer解析 | k=2(2連続塩基の出現頻度解析) | Biostrings

Biostringsパッケージを用いて、multi-FASTA形式ファイルを読み込んで、"AA", "AC", "AG", "AT", "CA", "CC", "CG", "CT", "GA", "GC", "GG", "GT", "TA", "TC", "TG", "TT"の計 $4^2 = 16$ 通りの2連続塩基の出現頻度を調べるやり方を示します。k-mer解析のk=2の場合に相当します。ヒトゲノムで"CG"の割合が期待値よりも低い(Lander et al., 2001; Saxonov et al., 2006)ですが、それを簡単に検証できます。

「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピペ。

1. イントロ | 一般 | ランダムな塩基配列を作成の4.を実行して得られたmulti-FASTAファイル(hoge4.fa)の場合:

タイトル通りの出現頻度です。

```
in_f <- "hoge4.fa"
out_f <- "hoge1.txt"

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")

#本番
out <- dinucleotideFrequency(fasta)

#ファイルに保存
tmp <- cbind(names(fasta), out)
write.table(tmp, out_f, sep="\t", append=TRUE)
```



#入力ファイル名を指定してin_fに格納

```
# hoge4.fa - メモ帳
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCAGAGACCCAGCACACCCT
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG
```

①例題1。②出力ファイルは、配列ごと(この場合コンティグごと)に16種類の2連続塩基の出現頻度をカウントしたもの。(作業ディレクトリにhoge4.faがないというエラーに何人か遭遇するだろうが、基本自力で解決せよ)

k-mer解析(k=2)

1. イントロ | 一般 | ランダムな塩基配列を作成の4.を実行して得られた

タイトル通りの出現頻度です。

```

in_f <- "hoge4.fa"           #入力ファイル
out_f <- "hoge1.txt"        #出力ファイル

#必要なパッケージをロード
library(Biostrings)        #パッケージ

#入力ファイルの読み込み
fasta <- readDNAStringSet(in_f, format="fasta") #確認して
fasta

#本番
out <- dinucleotideFrequency(fasta) #連続塩基

#ファイルに保存
tmp <- cbind(names(fasta), out) #保存した
write.table(tmp, out_f, sep="\t", append=F, qu

```

```

hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCAGAGACCCAGCACACCCT
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG

```

出力:hoge1.txt

	AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
contig_1	0	1	1	2	2	2	3	2	2	2	3	0	0	3	0	0
contig_2	4	6	9	1	11	11	5	6	4	9	10	8	1	8	6	3
contig_3	2	4	5	4	4	2	5	2	4	3	7	6	6	4	3	3
contig_4	3	6	2	3	5	3	3	4	3	3	1	2	3	2	4	1

①例題2。②出力ファイルは、配列ごとに2連続塩基の出現確率をカウントしたもの。
③as.probオプションをTRUEにしているだけ

k-mer解析(k=2)

2. イントロ | 一般 | ランダムな塩基配列を作成の4.を実行して得られたmulti-FASTAファイル(hoge4.fa)の場合:

出現頻度ではなく、出現確率を得るやり方です。

```

in_f <- "hoge4.fa" #入力ファイル
out_f <- "hoge2.txt" #出力ファイル

#必要なパッケージをロード
library(Biostrings) #パッケージ

#入力ファイルの読み込み
fasta <- readDNAStringSet(in_f, format="fasta") #確認して
fasta

#本番
out <- dinucleotideFrequency(fasta, as.prob=T) #

#ファイルに保存
tmp <- cbind(names(fasta), out) #保存した
write.table(tmp, out_f, sep="\t", append=F, qu
    
```

```

hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCAGAGACCCAGCACACCCT
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG
    
```

出力:hoge2.txt

	AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
contig_1	0.0%	4.3%	4.3%	8.7%	8.7%	8.7%	13.0%	8.7%	8.7%	8.7%	13.0%	0.0%	0.0%	13.0%	0.0%	0.0%
contig_2	3.9%	5.9%	8.8%	1.0%	10.8%	10.8%	4.9%	5.9%	3.9%	8.8%	9.8%	7.8%	1.0%	7.8%	5.9%	2.9%
contig_3	3.1%	6.3%	7.8%	6.3%	6.3%	3.1%	7.8%	3.1%	6.3%	4.7%	10.9%	9.4%	9.4%	6.3%	4.7%	4.7%
contig_4	6.3%	12.5%	4.2%	6.3%	10.4%	6.3%	6.3%	8.3%	6.3%	6.3%	2.1%	4.2%	6.3%	4.2%	8.3%	2.1%

①例題7。②ヒトゲノムRパッケージを入力とすることもできます。③fastaオブジェクト以降は同じ。平成27年度講習会(2015.07.30のスライド40あたり)にもあり

k-mer解析(k=2)

2. イントロ | 一般 | ランダムな塩基配列を作成の4.を実行して得られたmulti-FASTAファイル(hoge4.fa)の場合:

出現頻度ではなく、出現確率を得るやり方です。

```

in_f <- "hoge4.fa"           #入力ファイル名を指定してin_fに格納
out_f <- "hoge2.txt"        #出力ファイル名を指定してout_fに格納

#必要なパッケージをロード
library(Biostrings)        #パッケージの読み込み

```

```

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, as.fastq=T)
fasta

#本番
out <- dinucleotideFrequency(fasta)

#ファイルに保存
tmp <- cbind(names(fasta), out)
write.table(tmp, out_f, sep="\t")

```



7. ヒトゲノム配列パッケージ(BSgenome.Hsapiens.NCBI.GRCh38)の場合:

2013年12月にリリースされたGenome Reference Consortium GRCh38です。出力は出現確率です。

```

out_f <- "hoge7.txt"        #出力ファイル名を指定してout_fに格納
param_bsgenome <- "BSgenome.Hsapiens.NCBI.GRCh38" #パッケージ名を指定(BSgenome系のゲノム)

#必要なパッケージをロード
library(Biostrings)        #パッケージの読み込み
library(param_bsgenome, character.only=T) #指定したパッケージの読み込み

#前処理(指定したパッケージ中のオブジェクト名をgenomeに統一)
tmp <- ls(paste("package", param_bsgenome, sep=":")) #指定したパッケージで利用可能なオブジェクト
genome <- eval(parse(text=tmp)) #文字列tmpをRオブジェクトとしてgenomeに格納(パッケージ)
fasta <- getSeq(genome)    #ゲノム塩基配列情報を抽出した結果をfastaに格納
names(fasta) <- seqnames(genome) #description情報を追加している
fasta                      #確認してるだけです

```



```

#本番
out <- dinucleotideFrequency(fasta, as.prob=T) #連続塩基の出現確率情報をoutに格納

#ファイルに保存
tmp <- cbind(names(fasta), out) #保存したい情報をtmpに格納
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F) #tmpの中身を指定したファイルに保存

```



例題7実行結果ファイル。約3分。CGの連続塩基が他に比べて確かに低いことがわかる。出現確率は染色体ごとにはばらつきがある。私ならbox plotを採用

k-mer解析(k=2)

出力:hoge7.txt

	AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	9.5%	5.0%	7.1%	7.4%	7.3%	5.4%	1.0%	7.1%	6.0%	4.4%	5.4%	5.0%	6.3%	6.0%	7.3%	9.6%
2	10.0%	5.0%	7.0%	7.9%	7.2%	5.0%	0.9%	7.0%	5.9%	4.1%	5.0%	5.0%	6.7%	5.9%	7.2%	10.0%
3	10.1%	5.0%	6.9%	8.0%	7.2%	4.9%	0.8%	6.9%	5.9%	4.0%	4.9%	5.0%	6.9%	5.9%	7.2%	10.2%
4	10.6%	5.0%	6.7%	8.5%	7.1%	4.5%	0.8%	6.7%	5.9%	3.8%	4.5%	5.0%	7.3%	5.8%	7.1%	10.6%
5	10.2%	5.0%	6.9%	8.1%	7.2%	4.8%	0.9%	6.9%	5.9%	4.0%	4.8%	5.1%	6.9%	5.9%	7.2%	10.3%
6	10.2%	5.0%	6.9%	8.1%	7.2%	4.8%	0.9%	6.9%	5.9%	4.0%	4.9%	5.0%	6.9%	5.9%	7.2%	10.2%
7	9.8%	5.0%	7.0%	7.7%	7.3%	5.1%	1.0%	7.0%	6.0%	4.2%	5.1%	5.1%	6.5%	5.9%	7.3%	10.0%
8	10.0%	5.1%	6.9%	7.9%	7.2%	5.0%	0.9%	6.9%	6.0%	4.1%	5.0%	5.0%	6.7%	5.9%	7.2%	10.0%
9	9.7%	5.1%	7.0%	7.6%	7.3%	5.3%	1.0%	7.0%	6.0%	4.3%	5.3%	5.0%	6.4%	6.0%	7.3%	9.7%
10	9.6%	5.0%	7.1%	7.5%	7.3%	5.3%	1.0%	7.1%	6.0%	4.4%	5.3%	5.1%	6.3%	6.0%	7.4%	9.7%
11	9.5%	5.1%	7.1%	7.5%	7.3%	5.3%	1.0%	7.1%	6.1%	4.3%	5.4%	5.0%	6.3%	6.0%	7.3%	9.6%
12	9.8%	5.0%	7.0%	7.7%	7.2%	5.1%	1.0%	7.0%	6.0%	4.2%	5.2%	5.1%	6.6%	6.0%	7.2%	9.9%
13	10.5%	5.0%	6.8%	8.4%	7.1%	4.5%	0.9%	6.7%	5.9%	3.8%	4.6%	5.0%	7.2%	5.8%	7.1%	10.6%
14	9.7%	5.0%	7.0%	7.7%	7.2%	5.1%	1.0%	7.0%	6.0%	4.2%	5.2%	5.1%	6.6%	5.9%	7.3%	9.9%
15	9.4%	5.1%	7.1%	7.3%	7.3%	5.4%	1.1%	7.1%	6.0%	4.5%	5.5%	5.1%	6.1%	6.0%	7.4%	9.5%
16	8.6%	5.1%	7.3%	6.7%	7.5%	6.1%	1.4%	7.2%	6.1%	5.0%	6.1%	5.1%	5.4%	6.1%	7.6%	8.8%
17	8.5%	5.1%	7.3%	6.4%	7.4%	6.3%	1.5%	7.4%	6.2%	5.1%	6.4%	5.0%	5.2%	6.1%	7.5%	8.6%
18	10.1%	5.1%	7.0%	7.9%	7.2%	4.7%	0.9%	6.9%	6.1%	4.0%	4.9%	5.1%	6.7%	5.9%	7.3%	10.3%

作図(box plot): 基本形

②例題10。box plotをPNG形式ファイルで出力するやり方の基本形。エアーハンズオン(やったつもり)でよい

- ・ イントロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- ・ イントロ | 一般 | [逆鎖\(reverse\)を取得](#) (last modified 2013/06/14)
- ・ イントロ | 一般 | k-mer解析 | k=1(塩基ごとの出現頻度解析) | [Biostrings](#) (last modified 2016/02/07)
- ・ イントロ | 一般 | k-mer解析 | k=2(2連続塩基の出現頻度解析) | [Biostrings](#) (last modified 2016/01/28)
- ・ イントロ | 一般 | k-mer解析 | k=3(3連続塩基の出現頻度解析) | [Biostrings](#) (last modified 2016/01/28)
- ・ イントロ | 一般 | k-mer解析 | k=n(n連続塩基の出現頻度解析) | [Biostrings](#) (last modified 2016/01/28)



イントロ | 一般 | k-mer解析 | k=2(2連続塩基の出現頻度解析) | Biostrings

[Biostrings](#)パッケージを用いて、multi-FASTA形式ファイルを読み込んで、"AA", "AC", "AG", "AT", "CA", "CC", "CG", "CT", "GA", "GC", "GG", "GT", "TA", "TC", "TG", "TT"の計 $4^2 = 16$ 通りの2連続塩基の出現頻度を調べるやり方を示します。k-mer解析のk=2の場合に相当します。ヒトゲノムで"CG"の割合が期待値よりも低い([Lander et al., 2001](#); [Saxonov et al., 2006](#))ですが、それを簡単に検証できます。

「ファイル」→「ディレクトリ」→ **②** 10. ヒトゲノム配列パッケージ([BSgenome.Hsapiens.NCBI.GRCh38](#))の場合:

1. イントロ | 一般 | ランク

タイトル通りの出現頻度

```
in_f <- "hoge4.fa"
out_f <- "hoge1.tx"
```

```
#必要なパッケージをロード
library(Biostrings)
```

```
#入力ファイルの読み込み
fasta <- readDNAST
fasta
```

7.と基本的に同じですが、box plotのPNGファイルも出力しています。

```
out_f1 <- "hoge10.txt" #出力ファイル名を指定してout_f1に格納
out_f2 <- "hoge10.png" #出力ファイル名を指定してout_f2に格納
param_bsgenome <- "BSgenome.Hsapiens.NCBI.GRCh38" #パッケージ名を指定(BSgenome系のゲノム)
param_fig <- c(700, 400) #ファイル出力時の横幅と縦幅を指定(単位はピクセル)
```

```
#必要なパッケージをロード
library(Biostrings) #パッケージの読み込み
library(param_bsgenome, character.only=T) #指定したパッケージの読み込み
```

```
#前処理(指定したパッケージ中のオブジェクト名をgenomeに統一)
tmp <- ls(paste("package", param_bsgenome, sep=":")) #指定したパッケージで利用可能なオブジェクト
genome <- eval(parse(text=tmp)) #文字列tmpをRオブジェクトとしてgenomeに格納(パッケージ)
fasta <- getSeq(genome) #ゲノム塩基配列情報を抽出した結果をfastaに格納
names(fasta) <- seqnames(genome) #description情報を追加している
fasta #確認してるだけです
```

```
#本番
out <- dinucleotideFrequency(fasta, as.prob=T) #連続塩基の出現確率情報をoutに格納
```

①例題10実行結果。②CGの連続塩基は、当然縦軸上で0.01付近に位置する。染色体全体の分布をbox plotで眺めても、確かに低い

作図(box plot): 基本形

10. ヒトゲノム配列パッケージ([BSgenome.Hsapiens.NCBI.GRCh38](#))の場合:

7. と基本的に同じですが、box plotのPNGファイルも出力しています。

```

out_f1 <- "hoge10.txt"           #出力ファイル名を指定してout_f1に格納
out_f2 <- "hoge10.png"         #出力ファイル名を指定してout_f2に格納
param_bsgenome <- "BSgenome.Hsapiens.NCBI.GRCh38" #パッケージ名を指定(BSgenome系のゲノム)
param_fig <- c(700, 400)       #ファイル出力時の横幅と縦幅を指定(単位はピクセル)

```

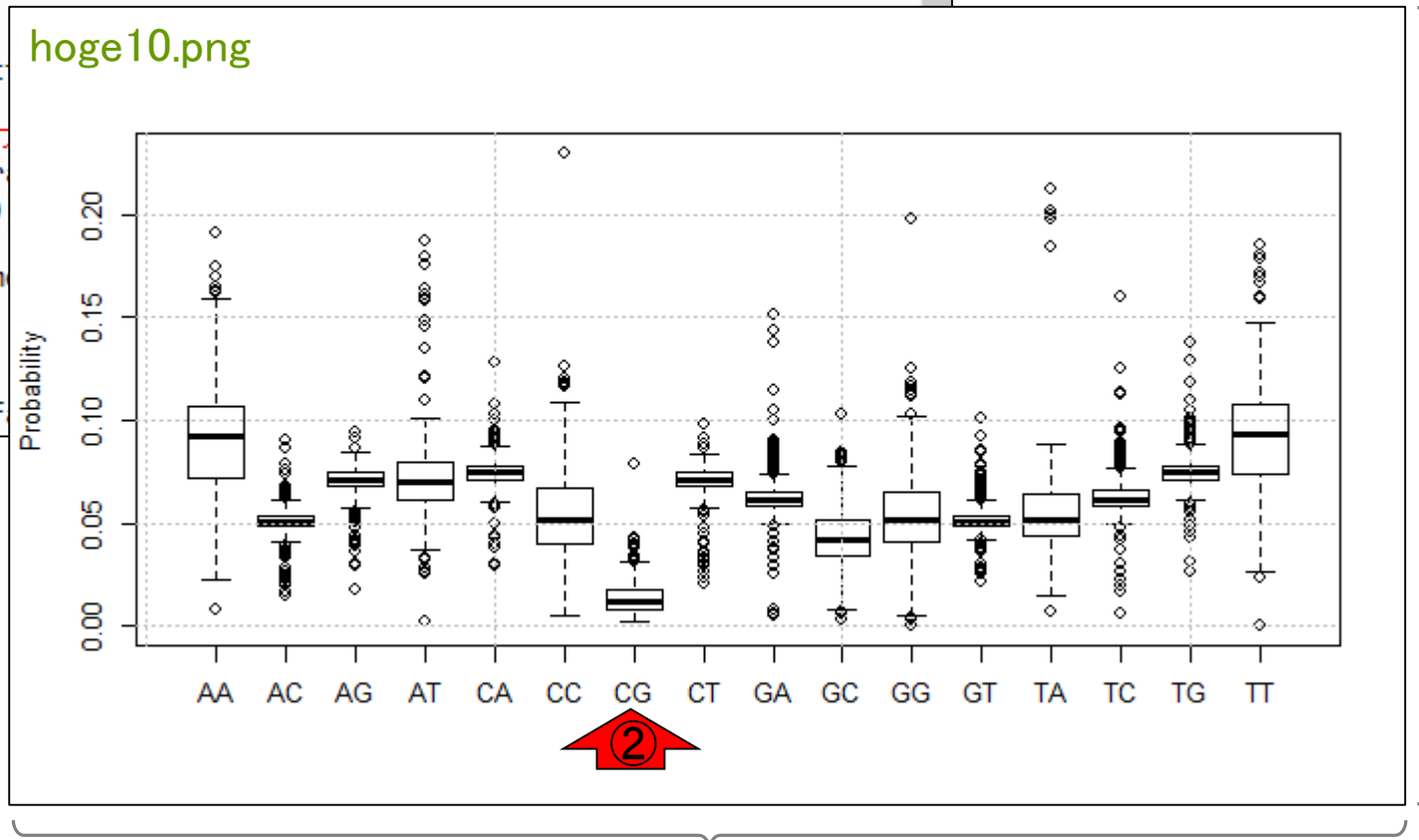
```

#必要なパッケージをロード
library(Biostrings)
library(param_bsgenome, character.only=TRUE)

#前処理(指定したパッケージ中のオブジェクトをロード)
tmp <- ls(paste("package", param_bsgenome))
genome <- eval(parse(text=tmp))
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)
fasta

#本番
out <- dinucleotideFrequency(fasta)

```



700 pixels

400 pixels

作図(box plot): 基本形

10. ヒトゲノム配列パッケージ(BSgenome.Hsapiens.NCBI.GRCh38)の場合:

7.と基本的に同じですが、box plotのPNGファイルも出力しています。

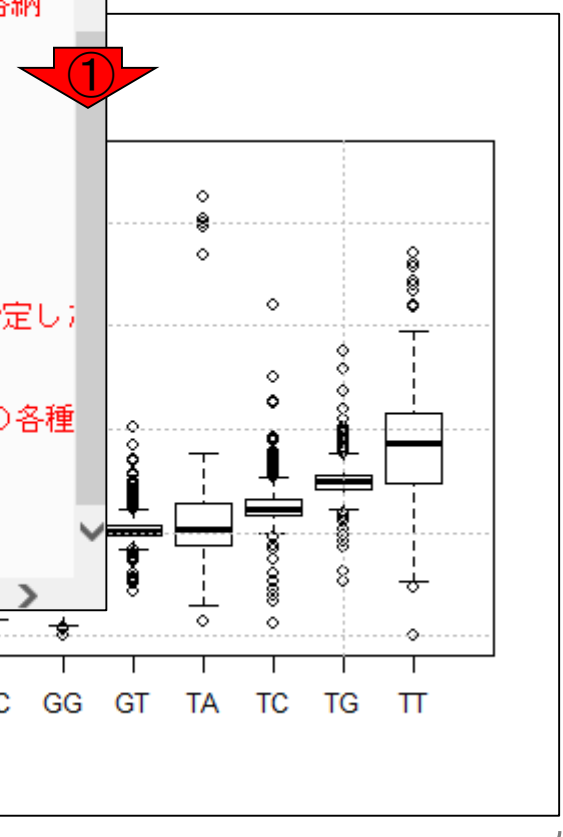
```
library(param_bsgenome, character.only=T)#指定したパッケージの読み込み

#前処理(指定したパッケージ中のオブジェクト名をgenomeに統一)
tmp <- ls(paste("package", param_bsgenome, sep=":"))#指定したパッケージで利用可能なオブ
genome <- eval(parse(text=tmp)) #文字列tmpをRオブジェクトとしてgenomeに格納(パッ
fasta <- getSeq(genome) #ゲノム塩基配列情報を抽出した結果をfastaに格納
names(fasta) <- seqnames(genome) #description情報を追加している
fasta #確認してるだけです

#本番
out <- dinucleotideFrequency(fasta, as.prob=T)#連続塩基の出現確率情報をoutに格納

#ファイルに保存(テキストファイル)
tmp <- cbind(names(fasta), out) #保存したい情報をtmpに格納
write.table(tmp, out_f1, sep="\t", append=F, quote=F, row.names=F)#tmpの中身を指定し;

#ファイルに保存(pngファイル)
png(out_f2, pointsize=13, width=param_fig[1], height=param_fig[2])#出力ファイルの各種
boxplot(out, ylab="Probability") #描画
grid(col="gray", lty="dotted") #指定したパラメータでグリッドを表示
dev.off() #おまじない
```



400 pixels

700 pixels

作図(box plot): 基本形

- ① dinucleotideFrequency関数実行結果のoutオブジェクトを入力として、
- ② boxplot関数を実行しているだけ

10. ヒトゲノム配列パッケージ(BSgenome.Hsapiens.NCBI.GRCh38)の場合:

7.と基本的に同じですが、box plotのPNGファイルも出力しています。

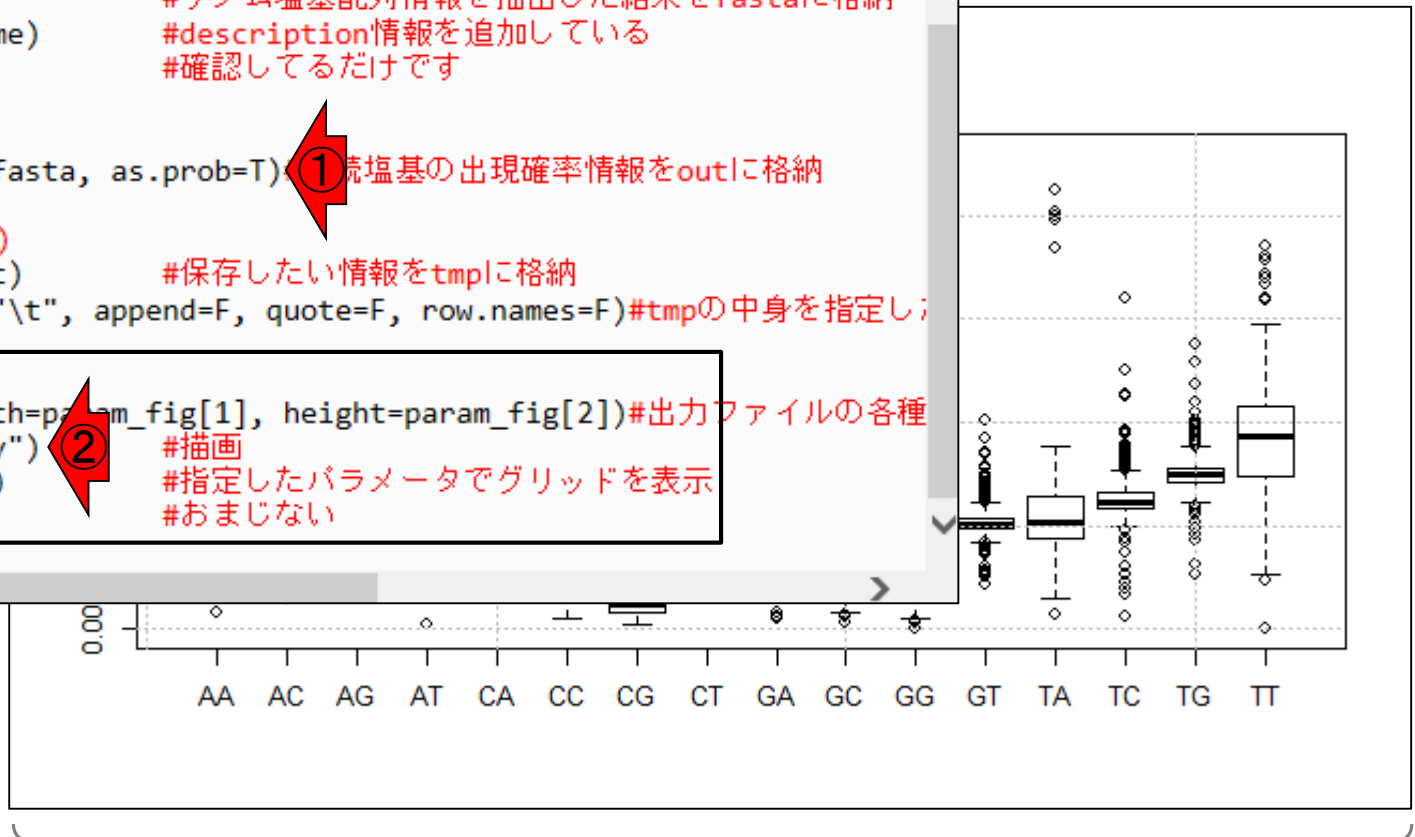
```
library(param_bsgenome, character.only=T)#指定したパッケージの読み込み

#前処理(指定したパッケージ中のオブジェクト名をgenomeに統一)
tmp <- ls(paste("package", param_bsgenome, sep=":"))#指定したパッケージで利用可能なオブ
genome <- eval(parse(text=tmp)) #文字列tmpをRオブジェクトとしてgenomeに格納(パッ
fasta <- getSeq(genome) #ゲノム塩基配列情報を抽出した結果をfastaに格納
names(fasta) <- seqnames(genome) #description情報を追加している
fasta #確認してるだけです

#本番
out <- dinucleotideFrequency(fasta, as.prob=T) #① 塩基の出現確率情報をoutに格納

#ファイルに保存(テキストファイル)
tmp <- cbind(names(fasta), out) #保存したい情報をtmpに格納
write.table(tmp, out_f1, sep="\t", append=F, quote=F, row.names=F)#tmpの中身を指定し

#ファイルに保存(pngファイル)
png(out_f2, pointsize=13, width=param_fig[1], height=param_fig[2])#出力ファイルの各種
boxplot(out, ylab="Probability") #② 描画
grid(col="gray", lty="dotted") #指定したパラメータでグリッドを表示
dev.off() #おまじない
```




400 pixels

700 pixels

②ヒトゲノムのGC含量は50%未満なので、塩基ごとの出現確率で見ると、A or TのほうがC or Gに比べて高い。そのため、2連続塩基の出現確率の期待値は③A or Tのみ、④C or Gのみ、⑤混合、の種類ごとに異なるという点に注意

注意点

- ① 解析したパッケージ名 : `BSgenome.Hsapiens.NCBI.GRCh38`
- ② ヒトゲノムの全体のGC含量 : 約41% 
- 各塩基(A, C, G, T)の出現確率 : (0.295, 0.205, 0.205, 0.295)
- ③ AA, AT, TA, TTの出現確率の期待値 = $0.295 \times 0.295 = 8.7\%$
- ④ CC, CG, GC, GGの出現確率の期待値 = $0.205 \times 0.205 = 4.2\%$
- ⑤ AC, AG, CA, CT, GA, GT, TC, TGの出現確率の期待値 = $0.205 \times 0.295 = 6.0\%$

期待値	8.7	6.0	6.0	8.7	6.0	4.2	4.2	6.0	6.0	4.2	4.2	6.0	8.7	6.0	6.0	8.7
連続塩基	AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	9.5%	5.0%	7.1%	7.4%	7.3%	5.4%	1.0%	7.1%	6.0%	4.4%	5.4%	5.0%	6.3%	6.0%	7.3%	9.6%
2	10.0%	5.0%	7.0%	7.9%	7.2%	5.0%	0.9%	7.0%	5.9%	4.1%	5.0%	5.0%	6.7%	5.9%	7.2%	10.0%
3	10.1%	5.0%	6.9%	8.0%	7.2%	4.9%	0.8%	6.9%	5.9%	4.0%	4.9%	5.0%	6.9%	5.9%	7.2%	10.2%
4	10.6%	5.0%	6.7%	8.5%	7.1%	4.5%	0.8%	6.7%	5.9%	3.8%	4.5%	5.0%	7.3%	5.8%	7.1%	10.6%
5	10.2%	5.0%	6.9%	8.1%	7.2%	4.8%	0.9%	6.9%	5.9%	4.0%	4.8%	5.1%	6.9%	5.9%	7.2%	10.3%
6	10.2%	5.0%	6.9%	8.1%	7.2%	4.8%	0.9%	6.9%	5.9%	4.0%	4.9%	5.0%	6.9%	5.9%	7.2%	10.2%
7	9.8%	5.0%	7.0%	7.7%	7.3%	5.1%	1.0%	7.0%	6.0%	4.2%	5.1%	5.1%	6.5%	5.9%	7.3%	10.0%
8	10.0%	5.1%	6.9%	7.9%	7.2%	5.0%	0.9%	6.9%	6.0%	4.1%	5.0%	5.0%	6.7%	5.9%	7.2%	10.0%
9	9.7%	5.1%	7.0%	7.6%	7.3%	5.3%	1.0%	7.0%	6.0%	4.3%	5.3%	5.0%	6.4%	6.0%	7.3%	9.7%
10	9.6%	5.0%	7.1%	7.5%	7.3%	5.3%	1.0%	7.1%	6.0%	4.4%	5.3%	5.1%	6.3%	6.0%	7.4%	9.7%
11	9.5%	5.1%	7.1%	7.5%	7.3%	5.3%	1.0%	7.1%	6.1%	4.3%	5.4%	5.0%	6.3%	6.0%	7.3%	9.6%

① 作図(box plot): 色づけ

①例題11。②colorという列名のところに2連続塩基の種類ごとに色を指定した、③タブ区切りファイル(human_2mer.txt)を与えて利用。④このファイルの情報を利用しているのは、コードの下のほう

11. ヒトゲノム配列パッケージ(BSgenome.Hsapiens.NCBI.GRCh38)の場合:

10.と基本的に同じですが、連続塩基の種類ごとの期待値とボックスプロット(box plot)上のファイル(human_2mer.txt)を入力として利用し、色情報のみを取り出して利用しています。

```

in_f <- "human_2mer.txt"
out_f1 <- "hoge11.txt"
out_f2 <- "hoge11.png"
param_bsgenome <- "BSgenome.Hsapiens.NCBI.GRCh38"
param_fig <- c(700, 400)

#必要なパッケージをロード
library(Biostrings)
library(param_bsgenome, character.only=T)

#入力ファイルの読み込み
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")

#前処理(指定したパッケージ中のオブジェクト名をgenomeに統一)
tmp <- ls(paste("package", param_bsgenome, sep=":"))
genome <- eval(parse(text=tmp))
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)
fasta

#本番
out <- dinucleotideFrequency(fasta, as.prob=T)
    
```

③

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_f1に格納
#出力ファイル名を指定してout_f2に格納
#パッケージ名を指定(BSgenome系の)
#ファイル出力時の横幅と縦幅を指定(単位はピク)

#パッケージの読み込み
#指定したパッケージの読み込み

#in_fで指定

#文字列tmpをRオブジェクトとしてgenomeに格納
#ゲノム塩基配列情報を抽出した結果をfastaに格納
#description情報を追加している
#確認してるだけです

#連続塩基の出現確率情報をoutに格納

④

②

type	expected	color
AA	0.087025	red
AC	0.060475	skyblue
AG	0.060475	skyblue
AT	0.087025	red
CA	0.060475	skyblue
CC	0.042025	black
CG	0.042025	black
CT	0.060475	skyblue
GA	0.060475	skyblue
GC	0.042025	black
GG	0.042025	black
GT	0.060475	skyblue
TA	0.087025	red
TC	0.060475	skyblue
TG	0.060475	skyblue
TT	0.087025	red

作図(box plot): 色づい

①boxplot関数実行時の②colオプション部分で③color列の情報を利用していることがわかる。④expected列情報は、例題11では利用していない

11. ヒトゲノム配列パッケージ(BSgenome.Hsapiens.NCBI.GRCh38)の場合:

10.と基本的に同じですが、連続塩基の種類ごとの期待値とボックスプロット(box plot)上での色情報を含むファイル(human_2mer.txt)を入力として利用し、色情報のみを取り出して利用しています。

```
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")#in_fで指
#前処理(指定したパッケージ中のオブジェクト名をgenomeに統一)
tmp <- ls(paste("package", param_bsgenome, sep=":"))#指定したパッケージで利用可能な
genome <- eval(parse(text=tmp)) #文字列tmpをRオブジェクトとしてgenomeに格納
fasta <- getSeq(genome) #ゲノム塩基配列情報を抽出した結果をfastaに格
names(fasta) <- seqnames(genome) #description情報を追加している
fasta #確認してるだけです

#本番
out <- dinucleotideFrequency(fasta, as.prob=T)#連続塩基の出現確率情報をoutに格納

#ファイルに保存(テキストファイル)
tmp <- cbind(names(fasta), out) #保存したい情報をtmpに格納
write.table(tmp, out_f1, sep="\t", append=F, quote=F, row.names=F)#tmpの中身を指

#ファイルに保存(pngファイル)
png(out_1, pointsize=13, width=param_fig[1], height=param_fig[2])#出力ファイルの
boxplot(out, ylab="Probability", col=as.character(data$color))#描画
grid(col="gray", lty="dotted") #指定したパラメータでグリッドを表示
dev.off() #おまじない
```

type	expected	color
AA	0.087025	red
AC	0.060475	skyblue
AG	0.060475	skyblue
AT	0.087025	red
CA	0.060475	skyblue
CC	0.042025	black
CG	0.042025	black
CT	0.060475	skyblue
GA	0.060475	skyblue
GC	0.042025	black
GG	0.042025	black
GT	0.060475	skyblue
TA	0.087025	red
TC	0.060475	skyblue
TG	0.060475	skyblue
TT	0.087025	red



作図(box plot): 色づけ

色づけてきている。同じ期待値(4.2%)の4種類間で比較しても、①CGは、他の②CC, ③GC, ④GGと比べて低い

11. ヒトゲノム配列パッケージ(BSgenome.Hsapiens.NCBI.GRCh38)の場合:

10.と基本的に同じですが、連続塩基の種類ごとの期待値とボックスプロット(box plot)上での色情報を含むファイル ([human_2mer.txt](#))を入力として利用し、色情報のみを取り出して利用しています。

```

in_f <- "human_2mer.txt"
out_f1 <- "hoge11.txt"
out_f2 <- "hoge11.png"
param_bsgenome <- "BSgenome.Hsapiens.NCBI.GRCh38"
param_fig <- c(700, 400)

```

#入力ファイル名を指定してin_fに格納
 #出力ファイル名を指定してout_f1に格納
 #出力ファイル名を指定してout_f2に格納
 #パッケージ名を指定(BSgenome系の)
 #ファイル出力時の横幅と縦幅を指定(単位はピク)

```

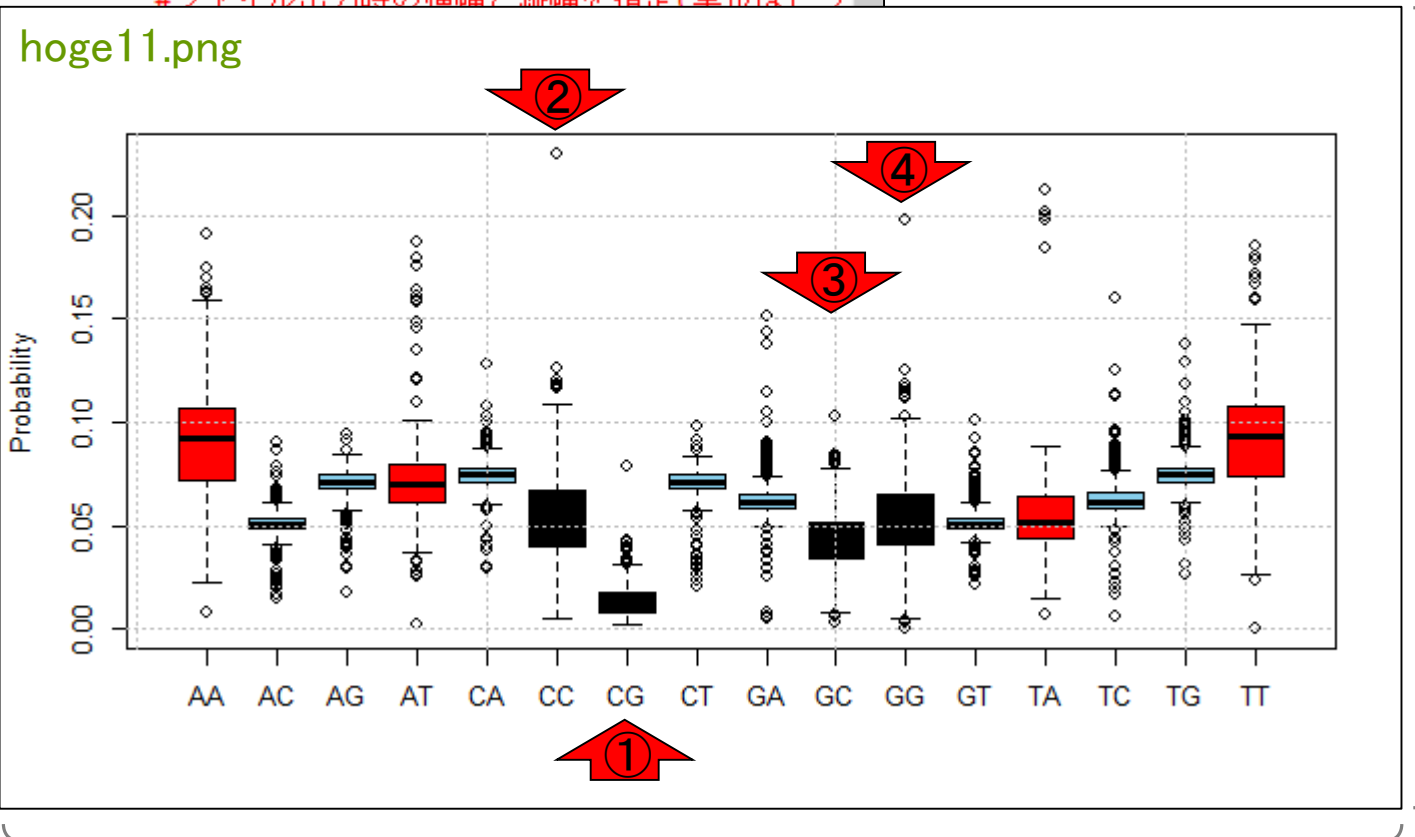
#必要なパッケージをロード
library(Biostrings)
library(param_bsgenome, charac

#入力ファイルの読み込み
data <- read.table(in_f, heade

#前処理(指定したパッケージ中のオ
tmp <- ls(paste("package", par
genome <- eval(parse(text=tmp)
fasta <- getSeq(genome)
names(fasta) <- seqnames(genom
fasta

#本番
out <- dinucleotideFrequency(f

```



700 pixels

400 pixels

期待値との差分

2連続塩基の出現確率の期待値は③A or Tのみ、④C or Gのみ、⑤混合、の種類ごとに異なるため、期待値を揃えてbox plotを描きたい

① 解析したパッケージ名 : `BSgenome.Hsapiens.NCBI.GRCh38`

② ヒトゲノムの全体のGC含量 : 約41%

各塩基(A, C, G, T)の出現確率 : (0.295, 0.205, 0.205, 0.295)

③ AA, AT, TA, TTの出現確率の期待値 = $0.295 \times 0.295 = 8.7\%$

④ CC, CG, GC, GGの出現確率の期待値 = $0.205 \times 0.205 = 4.2\%$

⑤ AC, AG, CA, CT, GA, GT, TC, TGの出現確率の期待値 = $0.205 \times 0.295 = 6.0\%$

期待値	8.7	6.0	6.0	8.7	6.0	4.2	4.2	6.0	6.0	4.2	4.2	6.0	8.7	6.0	6.0	8.7
連続塩基	AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	9.5%	5.0%	7.1%	7.4%	7.3%	5.4%	1.0%	7.1%	6.0%	4.4%	5.4%	5.0%	6.3%	6.0%	7.3%	9.6%
2	10.0%	5.0%	7.0%	7.9%	7.2%	5.0%	0.9%	7.0%	5.9%	4.1%	5.0%	5.0%	6.7%	5.9%	7.2%	10.0%
3	10.1%	5.0%	6.9%	8.0%	7.2%	4.9%	0.8%	6.9%	5.9%	4.0%	4.9%	5.0%	6.9%	5.9%	7.2%	10.2%
4	10.6%	5.0%	6.7%	8.5%	7.1%	4.5%	0.8%	6.7%	5.9%	3.8%	4.5%	5.0%	7.3%	5.8%	7.1%	10.6%
5	10.2%	5.0%	6.9%	8.1%	7.2%	4.8%	0.9%	6.9%	5.9%	4.0%	4.8%	5.1%	6.9%	5.9%	7.2%	10.3%
6	10.2%	5.0%	6.9%	8.1%	7.2%	4.8%	0.9%	6.9%	5.9%	4.0%	4.9%	5.0%	6.9%	5.9%	7.2%	10.2%
7	9.8%	5.0%	7.0%	7.7%	7.3%	5.1%	1.0%	7.0%	6.0%	4.2%	5.1%	5.1%	6.5%	5.9%	7.3%	10.0%
8	10.0%	5.1%	6.9%	7.9%	7.2%	5.0%	0.9%	6.9%	6.0%	4.1%	5.0%	5.0%	6.7%	5.9%	7.2%	10.0%
9	9.7%	5.1%	7.0%	7.6%	7.3%	5.3%	1.0%	7.0%	6.0%	4.3%	5.3%	5.0%	6.4%	6.0%	7.3%	9.7%
10	9.6%	5.0%	7.1%	7.5%	7.3%	5.3%	1.0%	7.1%	6.0%	4.4%	5.3%	5.1%	6.3%	6.0%	7.4%	9.7%
11	9.5%	5.1%	7.1%	7.5%	7.3%	5.3%	1.0%	7.1%	6.1%	4.3%	5.4%	5.0%	6.3%	6.0%	7.3%	9.6%

作図(box plot): 発展形

期待値との差分を評価すべく、①縦軸を $\log(\text{観測値}/\text{期待値})$ としてプロット。②0付近にある2連続塩基は、観測値(実測された出現確率)が期待値とほぼ同じことを意味する。この縦軸のような表現方法は一般的です

12. ヒトゲノム配列パッケージ(BSgenome.Hsapiens.NCBI.GRCh38)の場合:

11. と基本的に同じですが、human_2mer.txt というファイルを入力として与えて、連続塩基の箱プロット(box plot)上での色情報を利用しています。また、重要なのは期待値からの偏差(観測値(expected)と同程度の観測値(observed)であればゼロ、観測値のほうが大きければプラス、観測値のほうが小さければマイナス)といった具合で表現したほうがスマートです。それゆえ、box plotの縦軸を $\log(\text{observed}/\text{expected})$ として表現しています。CG以外の連続塩基は縦軸上で0付近に位置していることが分かります。

```
in_f <- "human_2mer.txt"
out_f1 <- "hoge12.txt"
out_f2 <- "hoge12.png"
param_bsgenome <- "BSgenome.Hsapiens.NCBI.GRCh38"
param_fig <- c(700, 400)
```

```
#必要なパッケージをロード
library(Biostrings)
library(param_bsgenome, character.only = TRUE)
```

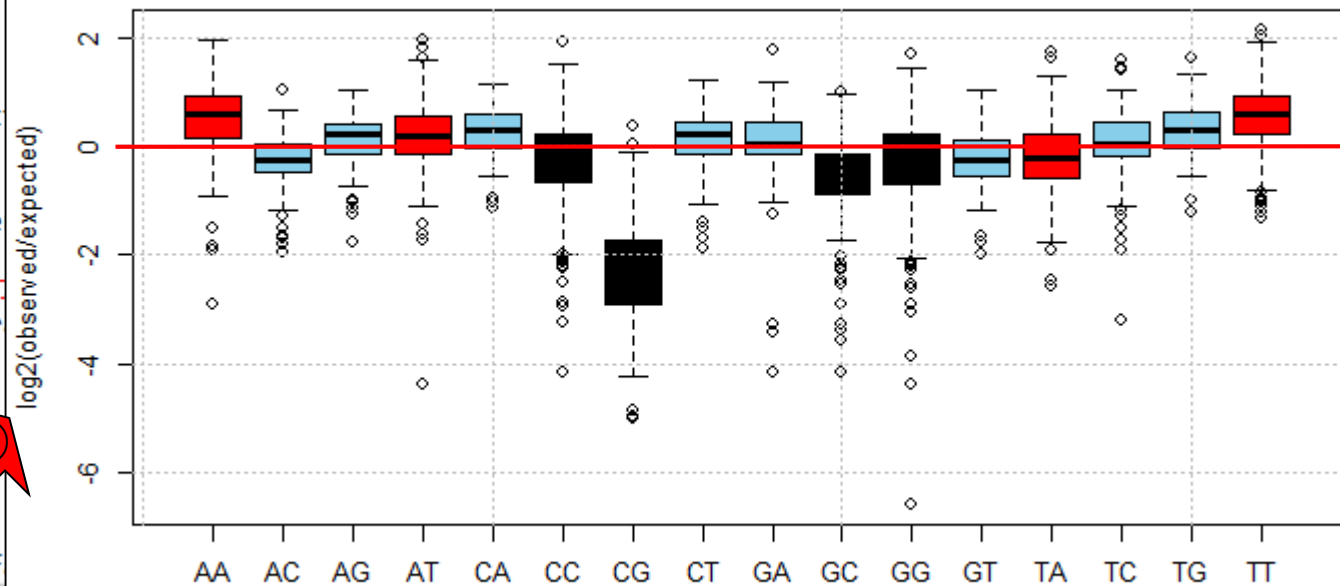
```
#入力ファイルの読み込み
data <- read.table(in_f, header = TRUE, as.is = TRUE)
```

```
#前処理(指定したパッケージ中のオブジェクトを生成)
tmp <- ls(paste("package", param_bsgenome))
genome <- eval(parse(text=tmp))
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)
fasta
```

```
#本番
out <- dinucleotideFrequency(fasta)
```

#入力ファイル名を指定してin_f1に格納
#出力ファイル名を指定してout_f1に格納

hoge12.png



作図(box plot): 発展形

- ①コピー実行時に警告メッセージが出た。これは、どこかに $-\text{Inf}$ (マイナス無限大のこと)があり、それはプロットされないということ。
- ②プロット対象はlogratioオブジェクトであり、
- ③でoutオブジェクトを期待値(data\$expected)で割ることで作成している

12. ヒトゲノム配列パッケージ(BSgenome.Hsapiens.NCBI.GRCh38)の場合:

11と基本的に同じですが、human_2mer.txtというファイルを入力として与えて、連続塩基プロット(box plot)上での色情報を利用しています。また、重要なのは期待値から観測値(expected)と同程度の観測値(observed)であればゼロ、観測値のほうが大きければ「プラス」、小さければ「マイナス」といった具合で表現したほうがスマートです。それゆえ、box plotの縦軸が $\log_2(\text{observed}/\text{expected})$ として表現しています。CG以外の連続塩基は縦軸上で0近辺に位置していることが分かります。

```
in_f <- "human_2mer.txt"
out_f1 <- "hoge12.txt"
out_f2 <- "hoge12.png"
param_bsgenome <- "BSgenome.Hsapiens.NCBI.GRCh38"
param_fig <- c(700, 400)

#必要なパッケージをロード
library(Biostrings)
library(param_bsgenome, character.only=TRUE)

#入力ファイルの読み込み
data <- read.table(in_f, header=TRUE, as.is=TRUE)

#前処理(指定したパッケージ中のオブジェクトを読み込み)
tmp <- ls(paste("package", param_bsgenome))
genome <- eval(parse(text=tmp))
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)
fasta

#本番
out <- dinucleotideFrequency(fasta,
```

#入力ファイル名を指定してin_fに格納

```
R Console
> logratio <- log2(out/data$expected) #log(observed/expected)
>
> #ファイルに保存(テキストファイル)
> tmp <- cbind(names(fasta), logratio) #保存したい情報をtmpに格納
> write.table(tmp, out_f1, sep="\t", append=F, quote=F, row.names=FALSE)
>
> #ファイルに保存(pngファイル)
> png(out_f2, pointsize=13, width=param_fig[1], height=param_fig[2])
> boxplot(logratio, ylab="log2(observed/expected)", #描画
+         col=as.character(data$color)) #描画
警告メッセージ:
bplot(at[i], wid = width[i], stats = z$stats[, i], out = z$stats[, i])
Outlier (-Inf) in boxplot 16 is not drawn
> grid(col="gray", lty="dotted") #指定したパラメータで描画
> dev.off() #おまじない
null device
      1
> |
```

作図(box plot): 発展形

-Inf(マイナス無限大のこと)なので、①apply関数でlogratioオブジェクトの②各列(MARGIN=2)に対してmin関数を実行し、どの列に-Infがあったかを確認しているだけ。③logratioの行数と列数を表示。ヒトゲノムとはいえ、断片を合わせると455配列あるため、TTの実測値が0ということもあるのだろう。スライドを見るだけ

12. ヒトゲノム配列パッケージ(BSgenome.Hsapiens.NCBI.GRCh38)の場合:

11.と基本的に同じですが、human_2mer.txtというファイルを入力として与えて、連鎖プロット(box plot)上での色情報を利用しています。また、重要なのは期待値(expected)と同程度の観測値(observed)であればゼロ、観測値のほうが大きければ「マイナス」といった具合で表現したほうがスマートです。それゆえ、box plotとして表現しています。CG以外の連続塩基は縦軸上では0近辺に位置しているこ

```
in_f <- "human_2mer.txt"
out_f1 <- "hoge12.txt"
out_f2 <- "hoge12.png"
param_bsgenome <- "BSgenome.Hsapiens.NCBI.GRCh38"
param_fig <- c(700, 400)

#必要なパッケージをロード
library(Biostrings)
library(param_bsgenome, character.only=TRUE)

#入力ファイルの読み込み
data <- read.table(in_f, header=TRUE, as.is=TRUE)

#前処理(指定したパッケージ中のオブジェクトを利用)
tmp <- ls(paste("package", param_bsgenome))
genome <- eval(parse(text=tmp))
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)
fasta

#本番
out <- dinucleotideFrequency(fasta)
```

#入力ファイル名を指定してin_fに格納

```
R Console
> apply(logratio, MARGIN=2, max)
      AA      AC      AG      AT      CA      CC
1.9553181 1.0365555 1.0380107 1.9660966 1.1534243 1.9262110
      CG      CT      GA      GC      GG      GT
0.3787232 1.2207788 1.7721007 1.0111684 1.7080309 1.0141437
      TA      TC      TG      TT
1.7346511 1.5812906 2.247335 2.1372668
> apply(logratio, MARGIN=2, min)
      AA      AC      AG      AT      CA      CC
-2.906679 -1.947401 -1.766627 -4.381588 -1.123647 -4.159398
      CG      CT      GA      GC      GG      GT
-5.021894 -1.888509 -4.157054 -4.172619 -6.606857 -1.981918
      TA      TC      TG      TT
-2.574233 -3.219269 -1.207795      -Inf
> dim(logratio)
[1] 455 16
> |
```



Contents

- NGS解析手段、ウェブツール(DDBJ Pipeline)との連携
- DDBJ PipelineでPlatanusを実行
- k-mer解析(k個の連続塩基に基づく各種解析)の基礎
 - 塩基ごとの出現頻度解析(k=1)、一気に計算
 - 2連続塩基の出現頻度解析(k=2)、基本スキルの復習や作図
- de novoアセンブリ時のエラー補正やゲノムサイズ推定の基本的な考え方
 - ランダムな塩基配列(仮想ゲノムおよび仮想NGSデータ)の作成
 - k-mer解析の応用、ゲノムサイズ推定の基礎
 - ゲノムサイズ推定(1,000 bpの仮想ゲノムの場合)
 - ゲノムサイズ推定(1,000 bpの仮想ゲノム; 4X → 10X coverageの場合)
 - k-mer出現頻度分布
 - シークエンスエラーを含む場合
 - 最終確認



この後に用いるランダムな塩基配列、および仮想NGSデータの作成について解説します。①サンプルデータの、②例題32。コピペはスライド75で行う

ランダムデータ生成

- インストール | Rパッケージ | [必要最小限\(数GB?!\)](#) (last modified 2015/05/25)
- インストール | Rパッケージ | [個別](#) (last modified 2015/06/10)
- (削除予定) [Rのインストールと起動](#) (last modified 2016/02/21)
- (削除予定) [個別パッケージのインストール](#) (last modified 2015/02/20)
- [基本的な利用法](#) (last modified 2015/04/03)
- [サンプルデータ](#) (last modified 2015/06/15)
- バイオインフォマティクス人材育成カリキュラム(次世代シーケンサ) | [NGSハンズオン講習会2016](#)
- バイオインフォマティクス人材育成カリキュラム(次世代シーケンサ) | [NGSハンズオン講習会2015](#)

サンプルデータ

1. Illumina/36bp/single-end/human (SRA000299) data (Marioni et al., Genome Res., 2008)

「Kidney 7 samples vs Liver 7 samples」のRNA-seqの遺伝子発現行列データ(SupplementaryTable2.txt)です。サンプルは二つの濃度(1.5 pM and 3 pM)でシーケンスされており、「3 pMのものが5 samples vs. 5 samples」、「1.5 pMのものが5 samples vs. 5 samples」に発現データが提供されています。

```
7列目 : R1L1K:
8列目 : R1L2L:
9列目 : R1L3K:
10列目 : R1L4L:
11列目 : R1L6L:
12列目 : R1L7L:
13列目 : R1L8L:
```

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル(sample32_ref.fastaとsample32_ngs.fasta)です。50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を10リード分だけランダム抽出したものです。塩基の存在比はAが22%, Cが28%, Gが28%, Tが22%にしています。リファレンス配列(仮想ゲノム配列)がsample32_ref.fastaで、10リードからなる仮想NGSデータがsample32_ngs.fastaです。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。イントロ | [NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成から](#)と基本的に同じです。

```
out_f1 <- "sample32_ref.fasta" #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample32_ngs.fasta" #出力ファイル名を指定してout_f2に格納
param_len_ref <- 50 #リファレンス配列の長さを指定
narabi <- c("A", "C", "G", "T") #以下の数値指定時にACGTの並びを間違えないようにする
param_composition <- c(22, 28, 28, 22) #(A,C,G,Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20 #リード長を指定
param_num_ngs <- 10 #リード数を指定
param_desc <- "kkk" #FASTA形式ファイルのdescription行に記述する内容

#必要なパッケージをロード
library(Biostrings) #パッケージの読み込み

#本番(リファレンス配列生成)
set.seed(1010) #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam_compositionで指定した数
```


ランダムデータ生成

①例題32。ここではまず②A, C, G, Tの並びで、③塩基の存在比率が22, 28, 28, 22のランダム塩基配列を生成させようとしている。④最終的に得られるのは2つのFASTA形式ファイル

32. Rで解析用のランダム配列から生成したFASTA形式ファイル(sample32_ref.fasta, sample32_ngs.fasta)を生成する。50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を10リードからランダム抽出したものである。塩基の存在比はAが22%, Cが28%, Gが28%, Tが22%にしています。リファレンス配列(仮想ゲノム配列)がsample32_ref.fastaで、10リードからなる仮想NGSデータがsample32_ngs.fastaです。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。イントロ | NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成からと基本的に同じです。

```

out_f1 <- "sample32_ref.fasta" #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample32_ngs.fasta" #出力ファイル名を指定してout_f2に格納
param_len_ref <- 50 #リファレンス配列の長さを指定
narabi <- c("A", "C", "G", "T") #以下の数値指定時にACGTの並びを間違えないようにする
param_composition <- c(22, 28, 28, 22) # (A, C, G, Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20 #リード長を指定
param_num_ngs <- 10 #リード数を指定
param_desc <- "kkk" #FASTA形式ファイルのdescription行に記述する内容

#必要なパッケージをロード
library(Biostrings) #パッケージの読み込み

#本番(リファレンス配列生成)
set.seed(1010) #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam_compositionで指定した数
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="") #ACGTsetから
reference <- DNAStringSet(reference) #DNA塩基配列だと認識させるDNAStringSet関数を適用し
names(reference) <- param_desc #description行に相当する記述を追加している
reference #確認してるだけです

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T) #部分塩基配
    
```

仮想リファレンス

①はsingle-FASTA形式の仮想リファレンス(ゲノム)ファイル。②配列長は50 bp、③description部分はkkk。④この部分で作成しており、referenceという名前のオブジェクトを最後に①のファイル名で出力している

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル(sample32_ref.fasta)は50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を10リードからランダム抽出したものです。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)がsample32_ref.fastaで、10リードからなる仮想NGSデータがsample32_ngs.fastaです。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。イントロ | NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成からと基本的に同じです。

```

out_f1 <- "sample32_ref.fasta"
out_f2 <- "sample32_ngs.fasta"
param_len_ref <- 50
narabi <- c("A", "C", "G", "T")
param_composition <- c(22, 28, 28, 22)
param_len_ngs <- 20
param_num_ngs <- 10
param_desc <- "kkk"
  
```

① #出力ファイル名を指定してout_f1に格納
 ② #出力ファイル名を指定してout_f2に格納
 #リファレンス配列の長さを指定
 #以下の数値指定時にACGTの並びを間違えないようにする
 #(A,C,G,Tの並びで)各塩基の存在比率を指定
 #リード長を指定
 #リード数を指定
 #FASTA形式ファイルのdescription行に記述する内容

```

#必要なパッケージをロード
library(Biostrings)
  
```

#パッケージの読み込み

```

#本番(リファレンス配列生成)
set.seed(1010)
ACGTset <- rep(narabi, param_composition)
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="")
reference <- DNASTringSet(reference)
names(reference) <- param_desc
reference
  
```

#おまじない(同じ乱数になるようにするため)
 #narabi中の塩基がparam_compositionで指定した数
 #ACGTsetから
 #DNA塩基配列だと認識させるDNASTringSet関数を適用し
 #description行に相当する記述を追加している
 #確認してるだけです

```

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T)
  
```

#部分塩基配

仮想リファレンス

赤枠部分をコピー。①ランダムデータ生成と云いつつ、得られるreferenceオブジェクトは、②全員同じGAA...GCAという配列になっているはず。この理由は③set.seed関数で1010という同じタネ番号を指定しているから。この数字はなんでもよい

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル(sample32_ref.fasta)と50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列です。塩基の存在比はAが22%, Cが28%, Gが28%, Tが22%にしています。sample32_ref.fastaで、10リードからなる仮想NGSデータがsample32_ngs.fastaです。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることになります(4X coverageに相当)。イントロ | NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成からと基本的に同じです。

```
out_f1 <- "sample32_ref.fasta"
out_f2 <- "sample32_ngs.fasta"
param_len_ref <- 50
narabi <- c("A","C","G","T")
param_composition <- c(0.22, 0.28, 0.28, 0.22)
param_len_ngs <- 20
param_num_ngs <- 10
param_desc <- "kkk"

#必要なパッケージをロード
library(Biostrings)

#本番(リファレンス配列生成)
set.seed(1010)
ACGTset <- rep(narabi, param_composition)
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="")
reference <- DNASTringSet(reference)
names(reference) <- param_desc

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref * param_num_ngs), param_num_ngs)
```

```
R Console

setdiff, sort, table, tapply, union, unique, unlist,
unsplit

要求されたパッケージ S4Vectors をロード中です
要求されたパッケージ stats4 をロード中です
要求されたパッケージ IRanges をロード中です
要求されたパッケージ XVector をロード中です

>
> #本番(リファレンス配列生成)
> set.seed(1010) #おまじない(同じ乱数になす)
> ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam$
> reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="")
> reference <- DNASTringSet(reference) #DNA塩基配列だと認識させす
> names(reference) <- param_desc #description行に相当するす
> reference #確認してるだけです

A DNASTringSet instance of length 1
width seq names
[1] 50 GAAGTTGTCTTTTGAA...GAAGACGCGGACGGCA kkk
> |
```

仮想NGSリード

①はmulti-FASTA形式の仮想NGSリードファイル。
 ②長さは20 bp、③リード数は10。これらのリードは、仮想リファレンス(ゲノム)ファイルとなるreferenceオブジェクトから乱数を発生させて、部分配列として切り出している。
 ④下のほうに移動

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル(sample32_ref.fasta)と50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を生成する。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。sample32_ref.fastaで、10リードからなる仮想NGSデータがsample32_ngs.fastaです。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。[イントロ](#) | [NGS](#) | [配列取得](#) | [シミュレーションデータ](#) | [ランダムな塩基配列の生成から](#)と基本的に同じです。

```

out_f1 <- "sample32_ref.fasta"
out_f2 <- "sample32_ngs.fasta"
param_len_ref <- 50
narabi <- c("A","C","G","T")
param_composition <- c(22, 28, 28, 22)
param_len_ngs <- 20
param_num_ngs <- 10
param_desc <- "kkk"

#必要なパッケージをロード
library(Biostrings)

#本番(リファレンス配列生成)
set.seed(1010)
ACGTset <- rep(narabi, param_composition)
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="")
reference <- DNAStringSet(reference)
names(reference) <- param_desc
reference

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T)
    
```

① #出力ファイル名を指定してout_f1に格納
 #出力ファイル名を指定してout_f2に格納
 #リファレンス配列の長さを指定
 #以下の数値指定時にACGTの並びを間違えないようにする
 #(A,C,G,Tの並びで)各塩基の存在比率を指定

② #リード長を指定

③ #リード数を指定
 #FASTA形式ファイルのdescription行に記述する内容

④ #パッケージの読み込み
 #おまじない(同じ乱数になるようにするため)
 #narabi中の塩基がparam_compositionで指定した数
 #ACGTsetから
 #DNA塩基配列だと認識させるDNAStringSet関数を適用し
 #description行に相当する記述を追加している
 #確認してるだけです

仮想NGSリード

①sample関数で乱数を発生させて、得られた乱数s_posiを②subseq関数でreferenceから切り出すときの④start位置として利用している。⑤20 bpというparam_len_ngsはwidthオプションとして与えている

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル(sample32_ref.fasta)

50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を10リードからランダム抽出したものです。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)がsample32_ref.fastaで、10リードからなる仮想NGSデータがsample32_ngs.fastaです。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。イントロ | NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成からと基本的に同じです。

```
reference <- DNASTringSet(reference) #DNA塩基配列だと認識させるDNASTringSet関数を適用し
names(reference) <- param_desc #description行に相当する記述を追加している
reference #確認してるだけです

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T)#部分塩基配
s_posi #確認してるだけです
hoge <- NULL #最終的に得る結果を格納するためのベクトルホルダhoge
for(i in 1:length(s_posi)){ #length(s_posi)回だけループを回す
  hoge <- append(hoge, subseq(reference, start=s_posi[i], width=param_len_ngs))#sub
}
fasta <- hoge #hogeの中身をfastaに格納

#後処理(description部分の作成)
description <- paste(param_desc, s_posi, (s_posi+param_len_ngs-1), sep="_")#param_desc
names(fasta) <- description #description行に相当する記述を追加している
fasta #確認してるだけです

#ファイルに保存(仮想リファレンス配列と仮想NGS配列)
writeXStringSet(reference, file=out_f1, format="fasta", width=50)#referenceの中身を指定
writeXStringSet(fasta, file=out_f2, format="fasta", width=50)#fastaの中身を指定したファ
```



仮想NGSリード

①出力ファイルのdescription部分にリファレンス配列のどの領域由来リードかを書き込んでいる

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル([sample32_ref.fasta](#)と[sample32_ngs.fasta](#))です。50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を10リード分だけランダム抽出したものです。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)が[sample32_ref.fasta](#)で、10リードからなる仮想NGSデータが[sample32_ngs.fasta](#)です。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。[イントロ](#) | [NGS](#) | [配列取得](#) | [シミュレーションデータ](#) | [ランダムな塩基配列の生成から](#)と基本的に同じです。

```
reference <- DNASTringSet(reference) #DNA塩基配列だと認識させるDNASTringSet関数を適用し
names(reference) <- param_desc #description行に相当する記述を追加している
reference #確認してるだけです

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T)#部分塩基配
s_posi #確認してるだけです
hoge <- NULL #最終的に得る結果を格納するためのブレースホルダhoge
for(i in 1:length(s_posi)){ #length(s_posi)回だけループを回す
  hoge <- append(hoge, subseq(reference, start=s_posi[i], width=param_len_ngs))#sub
}
fasta <- hoge #hogeの中身をfastaに格納

#後処理(description部分の作成)
description <- paste(param_desc, s_posi, (s_posi+param_len_ngs-1), sep="_")#param_desc
names(fasta) <- description #description行に相当する記述を追加している
fasta #確認してるだけです

#ファイルに保存(仮想リファレンス配列と仮想NGS配列)
writeXStringSet(reference, file=out_f1, format="fasta", width=50)#referenceの中身を指定
writeXStringSet(fasta, file=out_f2, format="fasta", width=50)#fastaの中身を指定したファ
```



最後までコピー

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル([sample32_ref.fasta](#)と[sample32_ngs.fasta](#))です。50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を10リード分だけランダム抽出したものです。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)が[sample32_ref.fasta](#)で、10リードからなる仮想NGSデータが[sample32_ngs.fasta](#)です。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。[イントロ](#) | [NGS](#) | [配列取得](#) | [シミュレーションデータ](#) | [ランダムな塩基配列の生成から](#)と基本的に同じです。

```
reference <- DNASTringSet(
names(reference) <- param
reference

#本番(シミュレーションデータ)
s_posi <- sample(1:(param
s_posi
hoge <- NULL
for(i in 1:length(s_posi)
  hoge <- append(hoge,
)
fasta <- hoge

#後処理(description部分の付
description <- paste(para
names(fasta) <- descripti
fasta

#ファイルに保存(仮想リファレン
writeXStringSet(reference, fi
writeXStringSet(fasta, fi
```

```
R Console
> names(fasta) <- description #description行に相当する$
> fasta #確認してるだけです
A DNASTringSet instance of length 10
width seq names
[1] 20 CACCAGGACATGAAGACGCG kkk_24_43
[2] 20 AGGACATGAAGACGCGGACG kkk_28_47
[3] 20 TTGAACTCACTACACCAGGA kkk_12_31
[4] 20 GTTGTCTTTTGAACCTACTA kkk_4_23
[5] 20 TTGTCTTTTGAACCTACTAC kkk_5_24
[6] 20 GTCTTTTGAACCTACTACAC kkk_7_26
[7] 20 TTGAACTCACTACACCAGGA kkk_12_31
[8] 20 ACACCAGGACATGAAGACGC kkk_23_42
[9] 20 GAAGTTGTCTTTTGAACCTCA kkk_1_20
[10] 20 AGTTGTCTTTTGAACCTACT kkk_3_22
>
> #ファイルに保存(仮想リファレンス配列と仮想NGS配列)
> writeXStringSet(reference, file=out_f1, format="fasta", width=5$
> writeXStringSet(fasta, file=out_f2, format="fasta", width=50)#f$
> |
```

出力ファイル

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル([sample32_ref.fasta](#)と[sample32_ngs.fasta](#))です。50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を10リード分だけランダム抽出したものです。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)が[sample32_ref.fasta](#)で、10リードからなる仮想NGSデータが[sample32_ngs.fasta](#)です。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。[イントロ | NGS | ① | シミュレーションデータ | ランダムな塩基配列の生成から](#)と基本的に同じです。

```

out_f1 <- "sample32_ref.fasta" #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample32_ngs.fasta" #出力ファイル名を指定してout_f2に格納
param_len_ref <- 50 #リファレンス配列の長さを指定
narabi <- c("A","C","G","T") #以下の数値指定時にACGTの並びを間違えないようにする
param_composition <- c(22, 28, 28, 22) #(A,C,G,Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20 #リード長を指定
param_num_ngs <- 10 #リード数を指定
param_desc <- "kkk" #FASTA形式ファイルのdescription行に記述する内容

#必要なパッケージ
library(Bios) > kkk
library(Bios) GAAGTTGTCTTTTGAACCTACTACACCAGGACATGAAGACGCGGACGGCA

#本番(リファレンス配列生成)
set.seed(1010) #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam_compositionで指定した数
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="") #ACGTsetから
reference <- DNASTringSet(reference) #DNA塩基配列だと認識させるDNASTringSet関数を適用し
names(reference) <- param_desc #description行に相当する記述を追加している
reference #確認してるだけです

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T) #部分塩基配
    
```


出力ファイル

①multi-FASTA形式の仮想NGSリードファイル。②長さは20 bp、③リード数は10

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル([sample32_ref.fasta](#)と[sample32_ngs.fasta](#))です。50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を10リード分だけランダム抽出したものです。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)が[sample32_ref.fasta](#)で、10リードからなる仮想NGSデータが[sample32_ngs.fasta](#)です。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。イントロ | NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成からと基本的に同じです。

```

out_f1 <- "sample32_ref.fasta"      #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample32_ngs.fasta"     #出力ファイル名を指定してout_f2に格納
param_len_ref <- 50                #リファレンス配列の長さを指定
narabi <- c("A", "C", "G", "T")   #以下の数値指定時にACGTの並びを間違えないようにする
param_composition <- c(22, 28, 28, 22) # (A,C,G,Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20                #リード長を指定
param_num_ngs <- 10               #リード数を指定
param_desc <- "kkk"               #FASTA形式ファイルのdescription行に記述する内容

#必要なパッケージをロード
library(Biostrings)               #パッケージの読み込み

#本番(リファレンス配列生成)
set.seed(1010)                    #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam_compositionで指定した数
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="") #ACGTsetから
reference <- DNAStringSet(reference) #DNA塩基配列だと認識させるDNAStringSet関数を適用し
names(reference) <- param_desc     #description行に相当する記述を追加している
reference                           #確認してるだけです

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T) #部分塩基配

```

```

>kkk_24_43
CACCAGGACATGAAGACGCG
>kkk_28_47
AGGACATGAAGACGCGGACG
>kkk_12_31
TTGAACTCACTACACCAGGA
>kkk_4_23
GTTGTCTTTTGAACCTACTA
>kkk_5_24
TTGTCTTTTGAACCTACTAC
>kkk_7_26
GTCTTTTGAACCTACTACAG
>kkk_12_31
TTGAACTCACTACACCAGGA
>kkk_23_42
ACACCAGGACATGAAGACGC
>kkk_1_20
GAAGTTGTCTTTTGAACCTCA
>kkk_3_22
AGTTGTCTTTTGAACCTACT

```



出力ファイル

仮想NGSリードファイルの①description部分を見れば、②そのリードが③仮想リファレンス配列中のどの領域由来かがわかる

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル(`sample32_ref.fasta`と`sample32_ngs.fasta`)です。50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を10リード分だけランダム抽出したものです。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)が`sample32_ref.fasta`で、10リードからなる仮想NGSデータが`sample32_ngs.fasta`です。リード長20塩基で10リードでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverage相当)。イントロ | NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成からと基本的に同じです。

```

out_f1 <- "sample32_ref.fasta" #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample32_ngs.fasta" #出力ファイル名を指定してout_f2に格納
param_len_ref <- 50 #リファレンス配列の長さを指定
narabi <- c("A","C","G","T") #以下の数値指定時にACGTの並びを間違えないようにする
param_composition <- c(22, 28, 28, 22) #(A,C,G,Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20 #リード長を指定
param_num_ngs <- 10 #リード数を指定
param_desc <- "kkk" #FASTA形式ファイルのdescription行に記述する内容
    
```

```

#必要なパッケージ
library(Bios)
>kkk
GAAGTTGTCTTTTGAAGTCACTACACCAGGACATGAAGACGCGGACGGCA
    
```

```

#本番(リファレンス配列生成)
set.seed(1010) #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam_compositionで指定した数
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="") #ACGTsetから
reference <- DNASTringSet(reference) #DNA塩基配列だと認識させるDNASTringSet関数を適用し
names(reference) <- param_desc #description行に相当する記述を追加している
reference #確認してるだけです
    
```

```

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T) #部分塩基配
    
```

```

>kkk_24_43
CACCAGGACAAGACGCGG
>kkk_28_47
AGGACATGAAGACGCGGACG
>kkk_12_31
TTGAACTCACTACACCAGGA
>kkk_4_23
GTTGTCTTTTGAAGTCACTA
>kkk_5_24
TTGTCTTTTGAAGTCACTAC
>kkk_7_26
GTCTTTTGAAGTCACTACAG
>kkk_12_31
TTGAACTCACTACACCAGGA
>kkk_23_42
ACACCAGGACATGAAGACGCG
>kkk_1_20
GAAGTTGTCTTTTGAAGTCA
>kkk_3_22
AGTTGTCTTTTGAAGTCACT
    
```



出力ファイル

総塩基数は仮想リファレンス配列が50 bp、仮想NGSリードが20 bp × 10リード = 200 bp。NGSデータはリファレンスの4倍の塩基数、つまり4X coverage。そしてsequence errorはゼロという解釈が可能

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル([sample32_ref.fasta](#)と[sample32_ngs.fasta](#))です。50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を10リード分だけランダム抽出したものです。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)が[sample32_ref.fasta](#)で、10リードからなる仮想NGSデータが[sample32_ngs.fasta](#)です。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。[イントロ](#) | [NGS](#) | [配列取得](#) | [シミュレーションデータ](#) | [ランダムな塩基配列の生成から](#)と基本的に同じです。

```
out_f1 <- "sample32_ref.fasta" #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample32_ngs.fasta" #出力ファイル名を指定してout_f2に格納
param_len_ref <- 50 #リファレンス配列の長さを指定
narabi <- c("A","C","G","T") #以下の数値指定時にACGTの並びを間違えないようにする
param_composition <- c(22, 28, 28, 22) #(A,C,G,Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20 #リード長を指定
param_num_ngs <- 10 #リード数を指定
param_desc <- "kkk" #FASTA形式ファイルのdescription行に記述する内容
```

```
#必要なパッケージ
library(Bios) > kkk
GAAGTTGTCTTTTGAAGTCACTACACCAGGACATGAAGACGCGGACGGCA
```

```
#本番(リファレンス配列生成)
set.seed(1010) #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam_compositionで指定した数
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="") #ACGTsetから
reference <- DNASTringSet(reference) #DNA塩基配列だと認識させるDNASTringSet関数を適用し
names(reference) <- param_desc #description行に相当する記述を追加している
reference #確認してるだけです
```

```
#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T) #部分塩基配
```

```
> kkk_24_43
CACCAGGACATGAAGACGCG
> kkk_28_47
AGGACATGAAGACGCGGACG
> kkk_12_31
TTGAACTCACTACACCAGGA
> kkk_4_23
GTTGTCTTTTGAAGTCACTA
> kkk_5_24
TTGTCTTTTGAAGTCACTAC
> kkk_7_26
GTCTTTTGAAGTCACTACAG
> kkk_12_31
TTGAACTCACTACACCAGGA
> kkk_23_42
ACACCAGGACATGAAGACGC
> kkk_1_20
GAAGTTGTCTTTTGAAGTCA
> kkk_3_22
AGTTGTCTTTTGAAGTCACT
```

Contents

- NGS解析手段、ウェブツール(DDBJ Pipeline)との連携
- DDBJ PipelineでPlatanusを実行
- k-mer解析(k個の連続塩基に基づく各種解析)の基礎
 - 塩基ごとの出現頻度解析(k=1)、一気に計算
 - 2連続塩基の出現頻度解析(k=2)、基本スキルの復習や作図
- de novoアセンブリ時のエラー補正やゲノムサイズ推定の基本的な考え方
 - ランダムな塩基配列(仮想ゲノムおよび仮想NGSデータ)の作成
 - k-mer解析の応用、ゲノムサイズ推定の基礎
 - ゲノムサイズ推定(1,000 bpの仮想ゲノムの場合)
 - ゲノムサイズ推定(1,000 bpの仮想ゲノム; 4X → 10X coverageの場合)
 - k-mer出現頻度分布
 - シークエンスエラーを含む場合
 - 最終確認



k-mer解析

■ 比較ゲノム解析

- k=3 or 4付近の値を用いてゲノムごとの頻度情報を取得し、類似性尺度として利用

■ アセンブリ(ゲノムやトランスクリプトーム)

- k=25~200付近の値を用いてde Bruijnグラフを作成
- k-mer頻度グラフを作成して眺め、Heterozygosityの有無などを調査



■ モチーフ解析

- 転写開始点の上流配列解析。古細菌の上流50塩基に絞ってk=4で出現頻度解析すると、おそらくTATAが上位にランクイン

■ 発現量推定

- RNA-seq解析で、リファレンスにリードをマップしてリード数をカウントするのが主流だが、マッピング作業をすっ飛ばしてk-merに基づく方法で定量。Sailfish (Patro et al., *Nat Biotechnol.*, 2014)やRNA-Skim (Zhang and Wang, *Bioinformatics*, 2014)。

ゲノムアセンブリ概略

大きく分けて4つの手順からなる。k-merは、主にIlluminaに代表される～数百塩基程度のリード長のNGSデータに適用される。2014.06.25の講義資料をベースに作成

1. 前処理(pre-processing filtering)

- クオリティの低いリードやコンタミを除去するステップ。塩基置換(substitution)やインデル(indels; insertion/deletion)を含むリードの除去や補正(error correction)。
- 4つのアプローチ: **k-mer**, suffix tree/array, multiple sequence alignment, hybrid

2. グラフ構築(graph construction)

- 前処理後のリードを用いて、リード間のオーバーラップ(overlap)を頼りにつなげていくステップ。シーケンスエラー(sequencing error)と多型(polymorphism)の違いを見るべく、グラフ構築時にエラー補正を行うものもある。
- 4つのアプローチ: OLC, de Bruijn graph (**k-mer**), greedy, hybrid

3. グラフ簡易化(graph simplification)

- グラフ構築後に、複雑化したグラフをシンプルにしていくステップ。連続したノード(nodes; 頂点) やバブルのマージ作業に相当。

4. 後処理(post-processing)

- contigsやscaffoldsを得るステップ。ミスアセンブリの同定も含む。

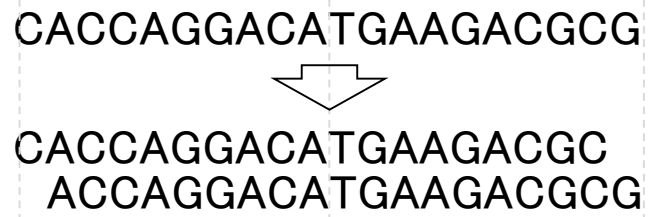
アセンブリ時の前処理部分でよく用いられる。kの値は概ねリード長の半分~2/3程度が採用されるが、解析戦略は多様化してきているので参考程度

k-mer出現頻度解析

■ 入力はNGSデータ

- リード長より短いk連続塩基からなる部分文字列を発生させるのが最初のステップ。発生させたk-merの出現頻度情報をもとに、カバレッジやゲノムサイズ推定、コンタミリードの除去などを行う

- 例1: 20塩基長のリードをk=19で分割すると、2個のk-merを発生可能



- 例2: 20塩基長のリードをk=17で分割すると4個のk-merを発生可能



- 例3: L塩基長のリードをk-merに分割すると(L - k + 1)個のk-merを発生可能

ゲノムサイズ推定

ここでのゲノムサイズ推定とは、①NGSリードファイルを入力としてk-mer解析を行い、②元のゲノム配列のサイズ(総塩基数)を予想すること。通常、*de novo*アセンブリの前処理段階で行う。この場合50 bpが正解

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル(sample32_ref.fasta)

50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を10リードからランダム抽出したものです。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)がsample32_ref.fastaで、10リードからなる仮想NGSデータがsample32_ngs.fastaです。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。イントロ | NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成からと基本的に同じです。

```

out_f1 <- "sample32_ref.fasta" #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample32_ngs.fasta" #出力ファイル名を指定してout_f2に格納
param_len_ref <- 50 #リファレンス配列の長さを指定
narabi <- c("A","C","G","T") #以下の数値指定時にACGTの並びを間違えないようにする
param_composition <- c(22, 28, 28, 22) #(A,C,G,Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20 #リード長を指定
param_num_ngs <- 10 #リード数を指定
param_desc <- "kkk" #FASTA形式ファイルのdescription行に記述する内容

```

```

#必要なパッケージ
library(Biostrings)
>kkk
GAAGTTGTCTTTTGAAGTCACTACACCAGGACATGAAGACGCGGACGGCA

```

```

#本番(リファレンス配列生成)
set.seed(1010) #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam_compositionで指定した数
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="") #ACGTsetからランダムに抽出
reference <- DNASTringSet(reference) #DNA塩基配列だと認識させるDNASTringSet関数を適用
names(reference) <- param_desc #description行に相当する記述を追加している
reference #確認してるだけです

```

```

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T) #部分塩基配列の位置をランダムに抽出

```

```

>kkk_28_47
AGGACATGAAGACGCGGACG
>kkk_12_31
TTGAACTCACTACACCAGGA
>kkk_4_23
GTTGTCTTTTGAAGTCACTA
>kkk_5_24
TTGTCTTTTGAAGTCACTAC
>kkk_7_26
GTCTTTTGAAGTCACTACAG
>kkk_12_31
TTGAACTCACTACACCAGGA
>kkk_23_42
ACACCAGGACATGAAGACGC
>kkk_1_20
GAAGTTGTCTTTTGAAGTCA
>kkk_3_22
AGTTGTCTTTTGAAGTCACT

```



ゲノムサイズ推定

意義: *de novo*アセンブリの前処理段階で行うため、アセンブリ結果が推定ゲノムサイズ(や近縁種のゲノムサイズ)と似た結果になっているかどうかや、PacBioデータのアセンブリ実行時に事前情報としても利用

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル(sample32_ref.fasta)と、50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列と10リードからなるNGSデータ生成のためのFASTA形式ファイル(sample32_ngs.fasta)を生成します。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)がsample32_ref.fastaで、10リードからなる仮想NGSデータがsample32_ngs.fastaです。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。イントロ | NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成からと基本的に同じです。

```

out_f1 <- "sample32_ref.fasta" #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample32_ngs.fasta" #出力ファイル名を指定してout_f2に格納
param_len_ref <- 50 #リファレンス配列の長さを指定
narabi <- c("A","C","G","T") #以下の数値指定時にACGTの並びを間違えないようにする
param_composition <- c(22, 28, 28, 22) # (A,C,G,Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20 #リード長を指定
param_num_ngs <- 10 #リード数を指定
param_desc <- "kkk" #FASTA形式ファイルのdescription行に記述する内容
    
```

```

#必要なパッケージ
library(Bios)
>kkk
GAAGTTGTCTTTTGAAGTCACTACACCAGGACATGAAGACGCGGACGGCA
    
```

```

#本番(リファレンス配列生成)
set.seed(1010) #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam_compositionで指定した数
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="") #ACGTsetから
reference <- DNASTringSet(reference) #DNA塩基配列だと認識させるDNASTringSet関数を適用し
names(reference) <- param_desc #description行に相当する記述を追加している
reference #確認してるだけです
    
```

```

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T) #部分塩基配
    
```

```

>kkk_28_47
AGGACATGAAGACGCGGACG
>kkk_12_31
TTGAACTCACTACACCAGGA
>kkk_4_23
GTTGTCTTTTGAAGTCACTA
>kkk_5_24
TTGTCTTTTGAAGTCACTAC
>kkk_7_26
GTCTTTTGAAGTCACTACAG
>kkk_12_31
TTGAACTCACTACACCAGGA
>kkk_23_42
ACACCAGGACATGAAGACGC
>kkk_1_20
GAAGTTGTCTTTTGAAGTCA
>kkk_3_22
AGTTGTCTTTTGAAGTCACT
    
```



NGSデータからのゲノムサイズ推定の考え方を説明すべく、k=2の場合のk-mer出現頻度解析から順に述べる。①項目名に注意！②例題7。③入力ファイル

k-mer解析(k=2)

- ・イントロ | 一般 | [逆鎖\(reverse\)を取得](#) (last modified 2013/06/14)
- ・イントロ | 一般 | k-mer解析 | k=1(塩基ごとの出現頻度解析) | [Biostrings](#) (last modified 2016/04/27) **NEW**
- ・イントロ | 一般 | k-mer解析 | k=2(2連続塩基の出現頻度解析) | [Biostrings](#) (last modified 2016/01/28)
- ・イントロ | 一般 | k-mer解析 | k=3(3連続塩基の出現頻度解析) | [Biostrings](#) (last modified 2016/01/28)
- ・イントロ | 一般 | k-mer解析 | k=n(n連続塩基の出現頻度解析) | [Biostrings](#) (last modified 2016/04/30) **NEW**
- ・(削除予定)イントロ | 一般 | [2連続塩基の出現頻度情報取得](#) (last modified 2015/04/20)
- ・(削除予定)イントロ | 一般 | [2連続塩基の出現頻度情報取得](#) (last modified 2015/02/19)

イントロ | 一般 | k-mer解析 | k=n(n連続塩基の出現頻度解析) | Biostrings **NEW**

イントロ | 一般 | [Biostrings](#)パッケージを用いて、4ⁿ通りの任意のn連続塩基の出現頻度を調べるやり方を示します。k-mer解析のk=nに相当します。

② 7. サンプルデータの例題32を実行して得られたmulti-FASTAファイル(sample32_ngs.fasta)「ファイル」

1. イントロ | 一般 | [4連続塩基\(k=4\)の出現頻度情報取得](#)

4連続塩基(k=4)の出現頻度情報取得

```
in_f <- "sample32_ngs.fasta"
out_f <- "hoge7.txt"
param_kmer <- 4

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNAStringSet(in_f, format="fasta")
```

```
in_f <- "sample32_ngs.fasta"
out_f <- "hoge7.txt"
param_kmer <- 2
```

```
#必要なパッケージをロード
library(Biostrings)
```

```
#入力ファイルの読み込み
fasta <- readDNAStringSet(in_f, format="fasta")
```

```
#本番
out <- oligonucleotideFrequency(fasta, width=param_kmer)
```

```
#ファイルに保存
tmp <- cbind(names(fasta), out)
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F)
```

#入力ファイル名を指定してin_f
#出力ファイル名を指定してout_f
#k-merのkの値を指定

#パッケージの読み込み

#最初の列にID情報、そのあとに出現頻度情報のoutを結合したtmpの中身を指定したファイル名に格納

```
>KKK_24_43
CACCAGGACATGAAGACGCG
>kkk_28_47
AGGACATGAAGACGCGGACG
>kkk_12_31
TTGAACTCACTACACCAGGA
>kkk_4_23
GTTGTCTTTTGAAGCTACTA
>kkk_5_24
TTGTCTTTTGAAGCTACTAC
>kkk_7_26
GTCTTTTGAAGCTACTACAG
>kkk_12_31
TTGAACTCACTACACCAGGA
>kkk_23_42
ACACCAGGACATGAAGACGC
>kkk_1_20
GAAGTTGTCTTTTGAAGCTCA
>kkk_3_22
AGTTGTCTTTTGAAGCTCACT
```

①例題7コピペ実行後に、②出力ファイルの中身に相当するoutオブジェクトを表示。これは、 $4^2 = 16$ 種類の2連続塩基の出現頻度をリードごとに算出したもの。例えば③は3番目のリードの結果

k-mer解析(k=2)

7. サンプルデータの例題32を実行して得られたmulti-FASTAファイル(sample32.fasta)

2連続塩基(k=2)の出現頻度情報を得るやり方です。 $4^2 = 16$ 通りのk-merの出現頻度を計算することになります。リードごとに出現頻度を算出しています。

```
in_f <- "sample32_ngs.fasta"
out_f <- "hoge7.txt"
param_kmer <- 2

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")

#本番
out <- oligonucleotideFrequency(fasta, k=2)
```

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

#パッケージの読み込み

```
R Console
> #ファイルに保存
> tmp <- cbind(names(fasta), out) #最初の列にID情報$
> write.table(tmp, out_f, sep="\t", append=F, quote=F, r$
> out
```

	AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
[1,]	1	3	2	1	3	1	2	0	3	1	1	0	0	0	1	0
[2,]	1	3	2	1	1	0	3	0	4	1	2	0	0	0	1	0
[3,]	1	4	1	0	3	1	0	2	2	0	1	0	1	1	1	1
[4,]	1	2	0	0	1	0	0	3	1	0	0	2	1	2	2	4
[5,]	1	3	0	0	1	0	0	3	1	0	0	1	1	2	2	4
[6,]	1	4	0	0	2	0	0	3	1	0	0	1	1	2	1	3
[7,]	1	4	1	0	3	1	0	2	2	0	1	0	1	1	1	1
[8,]	1	4	2	1	3	1	1	0	3	1	1	0	0	0	1	0
[9,]	2	1	1	0	1	0	0	2	2	0	0	2	0	2	2	4
[10,]	1	2	1	0	1	0	0	3	1	0	0	2	0	2	2	4

```
>kkk_28_47
AGGACATGAAGACGCGGACG
>kkk_12_31
TTGAACTCACTACACCAGGA
>kkk_4_23
GTTGTCTTTTGAACCTACTA
```

ゲノムサイズ推定の場合は、①リードごとではなく全体をまとめて考える。②例えばTTという2連続塩基は、入力ファイル中に全部で21回出現したと解釈する

k-mer解析(k=2)

7. サンプルデータの例題32を実行して得られたmulti-FASTAファイル(sample32_ngs.fasta)の場合:

2連続塩基(k=2)の出現頻度情報を得るやり方です。4² = 16通りのk-merの出現頻度を計算することになります。リード毎に出現頻度を算出しています。

```
in_f <- "sample32_ngs.fasta"
out_f <- "hoge7.txt"
param_kmer <- 2
```

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

```
#必要なパッケージをロード
library(Biostrings)
```

```
#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")
```

```
#本番
out <- oligonucleotideFrequency(fasta,
```

```
#ファイルに保存
tmp <- cbind(names(fasta), out)
write.table(tmp, out_f, sep="\t", asper
```

#パッケージの読み込み

```
R Console
> out
      AA AC AG AT CA CC CG CT GA GC GG GT TA TC TG TT
[1,]  1  3  2  1  3  1  2  0  3  1  1  0  0  0  1  0
[2,]  1  3  2  1  1  0  3  0  4  1  2  0  0  0  1  0
[3,]  1  4  1  0  3  1  0  2  2  0  1  0  1  1  1  1
[4,]  1  2  0  0  1  0  0  3  1  0  0  2  1  2  2  4
[5,]  1  3  0  0  1  0  0  3  1  0  0  1  1  2  2  4
[6,]  1  4  0  0  2  0  0  3  1  0  0  1  1  2  1  3
[7,]  1  4  1  0  3  1  0  2  2  0  1  0  1  1  1  1
[8,]  1  4  2  1  3  1  1  0  3  1  1  0  0  0  1  0
[9,]  2  1  1  0  1  0  0  2  2  0  0  2  0  2  2  4
[10,] 1  2  1  0  1  0  0  3  1  0  0  2  0  2  2  4
> colSums(out)
      AA AC AG AT CA CC CG CT GA GC GG GT TA TC TG TT
      11 30 10  3 19  4  6 18 20  3  6  8  5 12 14 21
> |
```



①例題8は、全リードの合算結果を出力するコード。②outオブジェクトが合算後の数値ベクトル。③ $4^2 = 16$ 通り全ての2連続塩基が④1回以上出現していることがわかる

k-mer解析(k=2)

8. サンプルデータの例題32を実行して得られたmulti-FASTAファイル(sample32_ngs.fasta)の場合:

2連続塩基(k=2)の出現頻度情報を得るやり方です。 $4^2 = 16$ 通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を出力するやり方です。

```
in_f <- "sample32_ngs.fasta"
out_f <- "hoge8.txt"
param_kmer <- 2
```

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

```
#必要なパッケージをロード
library(Biostrings)
```

#パッケージの読み込み

```
#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")
```

```
#本番
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge) #列ごとの総和
```

```
#ファイルに保存
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=FALSE)
```

```
R Console
> #必要なパッケージをロード
> library(Biostrings) #パッケージの読み込み
> #入力ファイルの読み込み
> fasta <- readDNASTringSet(in_f, format="fasta") #in_fで読み込み
> #本番
> hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
> out <- colSums(hoge) #列ごとの総和をoutに格納
> #ファイルに保存
> write.table(out, out_f, sep="\t", append=F, quote=F, row.names=FALSE)
> out
```

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
11	30	10	3	19	4	6	18	20	3	6	8	5	12	14	21

ゲノムサイズ推定

ゲノムサイズ推定は、①kの値を大きくして、1回以上出現しているk連続塩基(k-mer)の種類数をカウントするのが基本

8. サンプルデータの例題32を実行して得られたmulti-FASTAファイル(sample32_ngs.fasta)の場合:

2連続塩基(k=2)の出現頻度情報を得るやり方です。4² = 16通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を出力するやり方です。

```

in_f <- "sample32_ngs.fasta"
out_f <- "hoge8.txt"
param_kmer <- 2

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")

#本番
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)

#ファイルに保存
write.table(out, out_f, sep="\t", append=F, quote=F, r$

```

#入力ファイル名を指定してin_fに格納
 #出力ファイル名を指定してout_fに格納
 #k-merのkの値を指定

#パッケージの読み込み

```

R Console
> #必要なパッケージをロード
> library(Biostrings) #パッケージの読み込み
>
> #入力ファイルの読み込み
> fasta <- readDNASTringSet(in_f, format="fasta") #in_fで$
>
> #本番
> hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
> out <- colSums(hoge) #列ごとの総和をoutに$
>
> #ファイルに保存
> write.table(out, out_f, sep="\t", append=F, quote=F, r$
> out
AA AC AG AT CA CC CG CT GA GC GG GT TA TC TG TT
11 30 10 3 19 4 6 18 20 3 6 8 5 12 14 21
> |

```

k-mer解析(k=3)

①例題9は、②k=3のときの、 $4^3 = 64$ 通りの3連続塩基ごとの出現頻度。いくつかの3連続塩基(例えばAAA)は1回も出現していないことがわかる

9. サンプルデータの例題32を実行して得られたmulti-FASTAファイル(sample32_ngs.fasta)の場合:

3連続塩基(k=3)の出現頻度情報を得るやり方です。 $4^3 = 64$ 通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を入力するやり方です。

```
in_f <- "sample32_ngs.fasta"
out_f <- "hoge9.txt"
param_kmer <- 3
```

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

```
#必要なパッケージをロード
library(Biostrings)
```

#パッケージの読み込み

```
#入力ファイルの読み込み
```

```
fasta <- readDNAStringSet(in_f, format="fasta")
```

```
#本番
```

```
hoge <- oligonucleotideFrequency(fasta, k=3)
```

```
out <- colSums(hoge)
```

```
> out <- colSums(hoge)
```

#列ごとの総和を\$

```
>
```

```
> #ファイルに保存
```

```
> write.table(out, out_f, sep="\t", append=F, quote=F, r$
```

```
> out
```

```
#ファイルに保存
```

```
write.table(out, out_f, sep="\t", apper
```

```
AAA AAC AAG AAT ACA ACC ACG ACT AGA AGC AGG AGT ATA ATC
  0  7  4  0  7  4  4 13  3  0  5  2  0  0
ATG ATT CAA CAC CAG CAT CCA CCC CCG CCT CGA CGC CGG CGT
  3  0  0 11  4  3  4  0  0  0  0  3  1  0
CTA CTC CTG CTT GAA GAC GAG GAT GCA GCC GCG GCT GGA GGC
  5  7  0  5 11  7  0  0  0  0  2  0  6  0
GGG GGT GTA GTC GTG GTT TAA TAC TAG TAT TCA TCC TCG TCT
  0  0  0  5  0  3  0  4  0  0  7  0  0  5
TGA TGC TGG TGT TTA TTC TTG TTT
 10  0  0  4  0  0 11 10
```

```
> |
```

k-mer解析(k=3)

当たり前だが、①outオブジェクトの要素数は64。
 ②1回以上出現した3連続塩基(k-mer)の種類数は32。この種類数が、ゲノムサイズの推定値に相当する。この場合、推定値は32 bpで正解は50 bp

9. サンプルデータの例題32を実行して得られたmulti-FASTAファイル(sam

3連続塩基(k=3)の出現頻度情報を得るやり方です。4³ = 64通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を出力するやり方です。

```
in_f <- "sample32_ngs.fasta"
out_f <- "hoge9.txt"
param_kmer <- 3

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fastq")

#本番
hoge <- oligonucleotideFrequency(fasta, k=3)
out <- colSums(hoge)

#ファイルに保存
write.table(out, out_f, sep="\t", as.is=TRUE)
```

#入力ファイル名を指定してin_fに格納
 #出力ファイル名を指定してout_fに格納
 #k-merのkの値を指定

#パッケージの読み込み

```
R Console
> out
AAA AAC AAG AAT ACA ACC ACG ACT AGA AGC AGG AGT ATA ATC
  0  7  4  0  7  4  4 13  3  0  5  2  0  0
ATG ATT CAA CAC CAG CAT CCA CCC CCG CCT CGA CGC CGG CGT
  3  0  0 11  4  3  4  0  0  0  0  3  1  0
CTA CTC CTG CTT GAA GAC GAG GAT GCA GCC GCG GCT GGA GGC
  5  7  0  5 11  7  0  0  0  0  2  0  6  0
GGG GGT GTA GTC GTG GTT TAA TAC TAG TAT TCA TCC TCG TCT
  0  0  0  5  0  3  0  4  0  0  7  0  0  5
TGA TGC TGG TGT TTA TTC TTG TTT
10  0  0  4  0  0 11 10
> length(out)
[1] 64
> sum(out > 0)
[1] 32
> |
```



この仮想データは

リマインド。①入力のNGSリードファイルは、②50 bpの③仮想リファレンスゲノム配列からランダム抽出して得られたもの。①の総塩基数は200 bpなので4X coverage

32. k-mer解析用のランダム配列から生成したFASTA形式ファイル([sample32_ref.fasta](#)と[sample32_ngs.fasta](#))です。50塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を10リード分だけランダム抽出したものです。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)が[sample32_ref.fasta](#)で、10リードからなる仮想NGSデータが[sample32_ngs.fasta](#)です。リード長20塩基で10リードなのでトータル200塩基となり、50塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。[イントロ](#) | [NGS](#) | [配列取得](#) | [シミュレーションデータ](#) | [ランダムな塩基配列の生成から](#)と基本的に同じです。

```
out_f1 <- "sample32_ref.fasta" #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample32_ngs.fasta" #出力ファイル名を指定してout_f2に格納
param_len_ref <- 50 #リファレンス配列の長さを指定
narabi <- c("A","C","G","T") #以下の数値指定時にACGTの並びを間違えないようにする
param_composition <- c(22, 28, 28, 22) #(A,C,G,Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20 #リード長を指定
param_num_ngs <- 10 #リード数を指定
param_desc <- "kkk" #FASTA形式ファイルのdescription行に記述する内容
```

```
#必要なパッケージ
library(Bios) > kkk
library(Bios) GAAGTTGTCTTTTGAAGTCACTACACCAGGACATGAAGACGCGGACGGCA
```

```
#本番(リファレンス配列生成)
set.seed(1010) #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam_compositionで指定した数
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="") #ACGTsetから
reference <- DNASTringSet(reference) #DNA塩基配列だと認識させるDNASTringSet関数を適用し
names(reference) <- param_desc #description行に相当する記述を追加している
reference #確認してるだけです
```

```
#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T) #部分塩基配
```

```
> kkk_24_43
CACCAGGACATGAAGACGCGG
> kkk_28_47
AGGACATGAAGACGCGGACG
> kkk_12_31
TTGAACTCACTACACCAGGA
> kkk_4_23
GTTGTCTTTTGAAGTCACTA
> kkk_5_24
TTGTCTTTTGAAGTCACTAC
> kkk_7_26
GTCTTTTGAAGTCACTACAG
> kkk_12_31
TTGAACTCACTACACCAGGA
> kkk_23_42
ACACCAGGACATGAAGACGCG
> kkk_1_20
GAAGTTGTCTTTTGAAGTCA
> kkk_3_22
AGTTGTCTTTTGAAGTCACT
```



Contents

- NGS解析手段、ウェブツール(DDBJ Pipeline)との連携
- DDBJ PipelineでPlatanusを実行
- k-mer解析(k個の連続塩基に基づく各種解析)の基礎
 - 塩基ごとの出現頻度解析(k=1)、一気に計算
 - 2連続塩基の出現頻度解析(k=2)、基本スキルの復習や作図
- de novoアセンブリ時のエラー補正やゲノムサイズ推定の基本的な考え方
 - ランダムな塩基配列(仮想ゲノムおよび仮想NGSデータ)の作成
 - k-mer解析の応用、ゲノムサイズ推定の基礎
 - ゲノムサイズ推定(1,000 bpの仮想ゲノムの場合)
 - ゲノムサイズ推定(1,000 bpの仮想ゲノム; 4X → 10X coverageの場合)
 - k-mer出現頻度分布
 - シークエンスエラーを含む場合
 - 最終確認



ランダムデータ生成

NGSリードファイル中に存在するk-merの種類数でゲノムサイズ推定ができることの実感を持ってもらうべく、1,000 bpの仮想ゲノム配列、および仮想NGSデータを作成する。①サンプルデータの、②例題33

- インストール | Rパッケージ | [必要最小限\(数GB?!\)](#) (last modified 2015/05/25)
- インストール | Rパッケージ | [個別](#) (last modified 2015/06/10)
- (削除予定) [Rのインストールと起動](#) (last modified 2016/02/21)
- (削除予定) [個別パッケージのインストール](#) (last modified 2015/02/20)
- [基本的な利用法](#) (last modified 2015/04/03)
- [サンプルデータ](#) (last modified 2015/06/15)
- バイオインフォマティクス人材育成カリキュラム(次世代シーケンサ) | [NGSハンズオン講習会2016](#)
- バイオインフォマティクス人材育成カリキュラム(次世代シーケンサ) | [NGSハンズオン講習会2015](#)
- バイオインフォマティクス人材育成カリキュラム(次世代シーケンサ) | [NGSハンズオン講習会2014](#)
- [書籍 | トランスクリプトミクス](#)
- [書籍 | トランスクリプトミクス](#)
- [書籍 | トランスクリプトミクス](#)
- [書籍 | トランスクリプトミクス](#)
- [書籍 | トランスクリプトミクス](#)

サンプルデータ

1. Illumina/36bp/single-end/human (SRA000299) data (Marioni et al., Genome Res., 2008)

「Kidney 7 samples vs Liver 7 samples」のRNA-seqの遺伝子発現行列データ(SupplementaryTable2.txt)です。サンプルは二つの濃度(1.5 pM and 3 pM)でシーケンスされており、「3 pMのものが5 samples vs. 5 samples」、「1.5 pMのものが5 samples vs. 5 samples」の2つのデータセットに発現データが提供されています。

7列目 : R1L1
8列目 : R1L2
9列目 : R1L3
10列目 : R1L4
11列目 : R1L5
12列目 : R1L6
13列目 : R1L7

33. k-mer解析用のランダム配列から生成したFASTA形式ファイル(sample33_ref.fastaとsample33_ngs.fasta)です。1000塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を200リード分だけランダム抽出したものです。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)がsample33_ref.fastaで、200リードからなる仮想NGSデータがsample33_ngs.fastaです。リード長20塩基で200リードなのでトータル4,000塩基となり、1,000塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。イントロ | [NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成から](#)と基本的に同じです。

```

out_f1 <- "sample33_ref.fasta" #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample33_ngs.fasta" #出力ファイル名を指定してout_f2に格納
param_len_ref <- 1000 #リファレンス配列の長さを指定
narabi <- c("A", "C", "G", "T") #以下の数値指定時にACGTの並びを間違えないようにするた
param_composition <- c(22, 28, 28, 22) #(A,C,G,Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20 #リード長を指定
param_num_ngs <- 200 #リード数を指定
param_desc <- "kkk" #FASTA形式ファイルのdescription行に記述する内容

#必要なパッケージをロード
library(Biostrings) #パッケージの読み込み

#本番(リファレンス配列生成)
set.seed(1010) #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam_compositionで指定した数だけ

```

ランダムデータ生成

①例題33。②1000 bpの仮想リファレンスゲノム配列、および③仮想NGSデータを作成する。20 bpのリードが200個なので、計4,000 bp、つまり4X coverageのデータを作成していることになる

33. k-mer解析用のランダム配列から生成したFASTA形式ファイル(sample33_ref.fasta)と、200リードからなる仮想NGSデータ(sample33_ngs.fasta)を作成する。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%としています。リファレンス配列(仮想ゲノム配列)がsample33_ref.fastaで、200リードからなる仮想NGSデータがsample33_ngs.fastaです。リード長20塩基で200リードなのでトータル4,000塩基となり、1,000塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。イントロ | NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成からと基本的に同じです。

```

out_f1 <- "sample33_ref.fasta" #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample33_ngs.fasta" #出力ファイル名を指定してout_f2に格納
param_len_ref <- 1000 #リファレンス配列の長さを指定
narabi <- c("A","C","G","T") #以下の数値指定時にACGTの並びを間違えないようにするた
param_composition <- c(22, 28, 28, 22) # (A,C,G,Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20 #リード長を指定
param_num_ngs <- 200 #リード数を指定
param_desc <- "kkk" #FASTA形式ファイルのdescription行に記述する内容

#必要なパッケージをロード
library(Biostrings) #パッケージの読み込み

#本番(リファレンス配列生成)
set.seed(1010) #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam_compositionで指定した数だけ
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="") #ACGTsetからpa
reference <- DNASTringSet(reference) #DNA塩基配列だと認識させるDNASTringSet関数を適用した
names(reference) <- param_desc #description行に相当する記述を追加している
reference #確認してるだけです

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T) #部分塩基配列
    
```

ランダムデータ生成

33. k-mer解析用のランダム配列から生成したFASTA形式ファイル([sample33_ref.fasta](#)と[sample33_nginx.fasta](#))です。1000塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を200リード分だけランダム抽出したものです。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)が[sample33_ref.fasta](#)で、200リードからなる仮想NGSデータが[sample33_nginx.fasta](#)です。リード長20塩基で200リードなのでトータル4,000塩基となり、1,000塩基からなる元のゲノム配列の4倍シーケンスしていることとなります(4X coverageに相当)。[イントロ](#) | [NGS](#) | [配列取得](#) | [シミュレーションデータ](#) | [ランダムな塩基配列の生成から](#)と基本的に同じです。

```
out_f1 <- "sample33_ref.fasta"
out_f2 <- "sample33_nginx.fasta"
param_len_ref <- 1000
narabi <- c("A","C","G","T")
param_composition <- c(22, 28, 28, 22)
param_len_nginx <- 20
param_num_nginx <- 200
param_desc <- "kkk"

#必要なパッケージをロード
library(Biostrings)

#本番(リファレンス配列生成)
set.seed(1010)
ACGTset <- rep(narabi, param_composition)
reference <- paste(sample(ACGTset, param_len_ref), collapse="")
reference <- DNASTringSet(reference)
names(reference) <- param_desc
reference

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_nginx), param_num_nginx)
```

```
R Console
> fasta #確認してるだけ$
A DNASTringSet instance of length 200
  width seq names
[1] 20 GGTACGGTTCGGTTGCCGA kkk_658_677
[2] 20 TTATCCGGCAGTCCTATATG kkk_197_216
[3] 20 CAATAGACACCACGCGCCAC kkk_559_578
[4] 20 ACCGAAAGTTCTGGAGGTCA kkk_55_74
[5] 20 TTGGGCCCGATATCTATAAA kkk_681_700
...
[196] 20 GGACACGGAACGGGGCTGAG kkk_283_302
[197] 20 TACGGTTCGGTTGCCGATC kkk_660_679
[198] 20 TCACAACGGCCGTCTGCATT kkk_101_120
[199] 20 TTAGGAAGGTACGGTTCCGG kkk_651_670
[200] 20 CTTTATTGGGACGCTTTTCA kkk_925_944
>
> #ファイルに保存(仮想リファレンス配列と仮想NGS配列)
> writeXStringSet(reference, file=out_f1, format="fasta")
> writeXStringSet(fasta, file=out_f2, format="fasta", width=20)
> |
```

k=3の場合のk-mer出現頻度解析を行う。①項目名に注意！②例題10。③入力ファイル

k-mer解析(k=3)

- ・ [イントロ](#) | [一般](#) | [逆鎖\(reverse\)を取得](#) (last modified 2013/06/14)
- ・ [イントロ](#) | [一般](#) | [k-mer解析](#) | [k=1\(塩基ごとの出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/04/27) **NEW**
- ・ [イントロ](#) | [一般](#) | [k-mer解析](#) | [k=2\(2連続塩基の出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/01/28)
- ・ [イントロ](#) | [一般](#) | [k-mer解析](#) | [k=3\(3連続塩基の出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/01/28)
- ・ [イントロ](#) | [一般](#) | [k-mer解析](#) | [k=n\(n連続塩基の出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/04/30) **NEW**
- ・ (削除予定) [イントロ](#) | [一般](#) | [2連続塩基の出現頻度情報取得](#) (last modified 2015/04/20)
- ・ (削除予定) [イントロ](#) | [一般](#) | [2連続塩基の出現頻度情報取得](#) (last modified 2015/02/19)

イントロ | 一般 | k-mer解析 | k=n(n連続塩基の出現頻度解析) | Biostrings **NEW**

[Biostrings](#)パッケージを用いて、 4^n 通りの任意のn連続塩基の出現頻度を調べるやり方を示します。k-mer解析のk=nの場合に

相当します。例えば、[例題10](#)の[サンプルデータの例題33](#)を実行して得られたmulti-FASTAファイル([sample33_ngs.fasta](#))の場合:

1. イントロ | 一般 | [ラン](#)

4連続塩基(k=4)の出

```
in_f <- "hoge4.
out_f <- "hoge1
param_kmer <- 4
```

```
#必要なパッケージ
library(Biostrin
```

```
#入力ファイルの読
fasta <- readDN
```

3連続塩基(k=3)の出現頻度情報を得るやり方です。 $4^3 = 64$ 通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を出力するやり方です。

```
in_f <- "sample33_ngs.fasta"
out_f <- "hoge10.txt"
param_kmer <- 3
```

```
#必要なパッケージをロード
library(Biostrings)
```

```
#入力ファイルの読み込み
```

```
fasta <- readDNAStringSet(in_f, format="fasta")#in_fで指定したファイルの読み込み
```

```
#本番
```

```
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)#k連続塩基の出現頻度情報をhoge1に格納
out <- colSums(hoge) #列ごとの総和をoutに格納
```

```
#ファイルに保存
```

```
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=T, col.names=F)#outの中身
```

①例題10は、②k=3のときの、③ $4^3 = 64$ 通りの3連続塩基ごとの出現頻度。④全ての3連続塩基(例えばAAA)は1回以上出現していることがわかる

k-mer解析(k=3)

10. サンプルデータの例題33を実行して得られたmulti-FASTAファイル(sample33_ngs.fasta)の場合:

3連続塩基(k=3)の出現頻度情報を得るやり方です。 $4^3 = 64$ 通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を入力するやり方です。

```
in_f <- "sample33_ngs.fasta"
out_f <- "hoge10.txt"
param_kmer <- 3
```

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

```
#必要なパッケージをロード
library(Biostrings)
```

#パッケージの読み込み

#入力ファイルの読み込み

```
fasta <- readDNAStringSet(in_f, format="fasta")
```

#本番

```
hoge <- oligonucleotideFrequency(fasta, width=3)
out <- colSums(hoge)
```

> out
AAA AAC AAG AAT ACA ACC ACG ACT AGA AGC AGG AGT ATA
36 52 42 44 58 89 59 45 36 47 55 39 39
ATC ATG ATT CAA CAC CAG CAT CCA CCC CCG CCT CGA CGC
49 62 32 34 90 62 18 56 124 121 54 66 70
CGG CGT CTA CTC CTG CTT GAA GAC GAG GAT GCA GCC GCG
85 47 50 35 93 49 55 75 53 44 30 68 64
GCT GGA GGC GGG GGT GTA GTC GTG GTT TAA TAC TAG TAT
94 66 85 90 53 54 52 39 43 51 42 21 76
TCA TCC TCG TCT TGA TGC TGG TGT TTA TTC TTG TTT
63 61 19 32 61 49 67 56 46 45 40 68

#ファイルに保存

```
write.table(out, out_f, sep="\t", append=FALSE)
```

```
> out
AAA AAC AAG AAT ACA ACC ACG ACT AGA AGC AGG AGT ATA
36 52 42 44 58 89 59 45 36 47 55 39 39
ATC ATG ATT CAA CAC CAG CAT CCA CCC CCG CCT CGA CGC
49 62 32 34 90 62 18 56 124 121 54 66 70
CGG CGT CTA CTC CTG CTT GAA GAC GAG GAT GCA GCC GCG
85 47 50 35 93 49 55 75 53 44 30 68 64
GCT GGA GGC GGG GGT GTA GTC GTG GTT TAA TAC TAG TAT
94 66 85 90 53 54 52 39 43 51 42 21 76
TCA TCC TCG TCT TGA TGC TGG TGT TTA TTC TTG TTT
63 61 19 32 61 49 67 56 46 45 40 68

> length(out)
[1] 64
> sum(out > 0)
[1] 64
> |
```

ゲノムサイズ推定

ゲノムサイズ推定は、①kの値を大きくして、②1回以上出現しているk連続塩基(k-mer)の種類数をカウントするのが基本。しかし、③k=3だと種類数の最大は $4^3 = 64$ しかないため、このデータの正解である1000 bpにはなりようがない。つまり①kの値を大きくする必要がある

10. サンプルデータの例題33を実行して得られたmulti-FASTAファイルの3連続塩基(k=3)の出現頻度情報を得るやり方です。4³ = 64通りの全リードを合算した出現頻度を入力するやり方です。

```

in_f <- "sample33_ngs.fasta"
out_f <- "hoge10.txt"
param_kmer <- 3

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")

#本番
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)

#ファイルに保存
write.table(out, out_f, sep="\t", append=FALSE)

```

#入力ファイル名を指定してin_fに格納
 #出力ファイル名を指定してout_fに格納
 #k-merのkの値を指定

#パッケージの読み込み

```

R Console
> out
AAA AAC AAG AAT ACA ACC ACG ACT AGA AGC AGG AGT ATA
36 52 42 44 58 89 59 45 36 47 55 39 39
ATC ATG ATT CAA CAC CAG CAT CCA CCC CCG CCT CGA CGC
49 62 32 34 90 62 18 56 124 121 54 66 70
CGG CGT CTA CTC CTG CTT GAA GAC GAG GAT GCA GCC GCG
85 47 50 35 93 49 55 75 53 44 30 68 64
GCT GGA GGC GGG GGT GTA GTC GTG GTT TAA TAC TAG TAT
94 66 85 90 53 54 52 39 43 51 42 21 76
TCA TCC TCG TCT TGA TGC TGG TGT TTA TTC TTG TTT
63 61 19 36 61 49 67 56 46 45 40 68

> length(out)
[1] 64

> sum(out > 0)
[1] 64

> |

```



ゲノムサイズ推定

kをどの程度まで大きくする必要があるかの感覚をつかむ。まず、1KB (1,000 bp)のゲノムサイズを推定するためには、①最低でもk=5以上にする必要がある。②k=10にすると1MB (1,000,000 bp)だが、これだと原理的に3GBのヒトゲノム(3,000,000,000 bp)には適用できない

10. サンプルデータの例題33を実行して得られたmulti-FASTAファイルの3連続塩基(k=3)の出現頻度情報を得るやり方です。4³ = 64通り全リードを合算した出現頻度を入力するやり方です。

```

in_f <- "sample33_ngs.fasta"      #入力ファイル名を指定してin_fに格納
out_f <- "hoge10.txt"            #出力ファイル名を指定してout_fに格納
param_kmer <- 3                  #k-merのkの値を指定

#必要なパッケージをロード
library(Biostrings)              #パッケージの読み込み

#入力ファイルの読み込み
fasta <- readDNAStringSet(in_f, format="fasta")

#本番
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)              #列ごとの合計

#ファイルに保存
write.table(out, out_f, sep="\t", append=FALSE)

```

```

R Console
> length(out)
[1] 64
> sum(out > 0)
[1] 64
> 4^5
[1] 1024
> 4^10
[1] 1048576
> 4^15
[1] 1073741824
> 4^16
[1] 4294967296
> 4^20
[1] 1.099512e+12
> |

```



ゲノムサイズ推定

3GBのヒトゲノム(3,000,000,000 bp)に適用可能な下限は、①k=16であることがわかる。②このデータの大元はゲノムサイズ1,000 bpなので、とりあえず③k=10にして1回以上出現するk-merの種類数を調べる

10. サンプルデータの例題33を実行して得られたmulti-FASTAファイル

3連続塩基(k=3)の出現頻度情報を得るやり方です。4³ = 64通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を入力するやり方です。

```
in_f <- "sample33_ngs.fasta"
out_f <- "hoge10.txt"
param_kmer <- 3
```

②

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

```
#必要なパッケージをロード
library(Biostrings)
```

#パッケージの読み込み

```
#入力ファイルの読み込み
fasta <- readDNAStringSet(in_f, format="fasta")
```

```
#本番
hoge <- oligonucleotideFrequency(fasta, width=3)
out <- colSums(hoge)
```

#列

```
#ファイルに保存
write.table(out, out_f, sep="\t", append=FALSE)
```

```
R Console
> length(out)
[1] 64
> sum(out > 0)
[1] 64
> 4^5
[1] 1024
> 4^10
[1] 1048576
> 4^15
[1] 1073741824
> 4^16
[1] 4294967296
> 4^20
[1] 1.099512e+12
> |
```

③

①

①例題11は、②k=10のときの、 $4^{10} = 1,048,576$ 通りの10連続塩基ごとの出現頻度を返すコード。③1回以上出現している10連続塩基は904個。つまり推定ゲノムサイズは904 bp。正解は1,000 bpなので近い値

k-mer解析(k=10)

11. サンプルデータの例題33を実行して得られたmulti-FASTAファイル10連続塩基(k=10)の出現頻度情報を得るやり方です。 $4^{10} = 1,048,576$ 通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を出力するやり方です。

```

in_f <- "sample33_ngs.fasta"      #入力ファイル名を指定してin_fに格納
out_f <- "hoge11.txt"            #出力ファイル名を指定してout_fに格納
param_kmer <- 10                 #k-merのkの値を指定

#必要なパッケージをロード
library(Biostrings)             #パッケージの読み込み

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")

#本番
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)             #列ごとの総計

#ファイルに保存
write.table(out, out_f, sep="\t", append=F, quote=F,
            length(out)         #4^k
            sum(out > 0)        #1回以上出現$)

```

```

R Console
>
> #入力ファイルの読み込み
> fasta <- readDNASTringSet(in_f, format="fasta")#in_ $
>
> #本番
> hoge <- oligonucleotideFrequency(fasta, width=param_ $
> out <- colSums(hoge)             #列ごとの総計
>
> #ファイルに保存
> write.table(out, out_f, sep="\t", append=F, quote=F$
> length(out)                     #4^param_kme$
[1] 1048576
> sum(out > 0)                     #1回以上出現$
[1] 904
> |

```

Contents

- NGS解析手段、ウェブツール(DDBJ Pipeline)との連携
- DDBJ PipelineでPlatanusを実行
- k-mer解析(k個の連続塩基に基づく各種解析)の基礎
 - 塩基ごとの出現頻度解析(k=1)、一気に計算
 - 2連続塩基の出現頻度解析(k=2)、基本スキルの復習や作図
- de novoアセンブリ時のエラー補正やゲノムサイズ推定の基本的な考え方
 - ランダムな塩基配列(仮想ゲノムおよび仮想NGSデータ)の作成
 - k-mer解析の応用、ゲノムサイズ推定の基礎
 - ゲノムサイズ推定(1,000 bpの仮想ゲノムの場合)
 - ゲノムサイズ推定(1,000 bpの仮想ゲノム; 4X → 10X coverageの場合)
 - k-mer出現頻度分布
 - シークエンスエラーを含む場合
 - 最終確認



ランダムデータ生成

例題33は4X coverageデータで904 bpだった(正解は1,000 bp)。Coverageの違いによるゲノムサイズ推定精度の違いを把握するべく、10X coverageデータを作成する。①サンプルデータの、②例題34

- インストール | Rパッケージ | [必要最小限\(数GB?!\)](#) (last modified 2015/05/25)
- インストール | Rパッケージ | [個別](#) (last modified 2015/06/10)
- (削除予定)Rのインストールと起動 (last modified 2016/02/21)
- (削除予定)個別パッケージのインストール (last modified 2015/02/20)
- [基本的な利用法](#) (last modified 2015/04/03)
- [サンプルデータ](#) (last modified 2015/06/15)
- バイオインフォマティクス人材育成カリキュラム(次世代シーケンサ) | [NGSハンズオン講習会2016](#)
- バイオインフォマティクス人材育成カリキュラム(次世代シーケンサ) | [NGSハンズオン講習会2015](#)
- バイオインフォマティクス人材育成カリキュラム(次世代シーケンサ) | [NGSハンズオン講習会2015](#)
- [書籍 | トランスクリプトミクス](#)
- [書籍 | トランスクリプトミクス](#)
- [書籍 | トランスクリプトミクス](#)
- [書籍 | トランスクリプトミクス](#)
- [書籍 | トランスクリプトミクス](#)

サンプルデータ

1. Illumina/36bp/single-end/human (SRA000299) data ([Marioni et al., Genome Res., 2008](#))

「Kidney 7 samples vs Liver 7 samples」のRNA-seqの遺伝子発現行列データ([Supplementary Table 2.txt](#))です。サンプルは2つの濃度(1.5 pM and 3 pM)でシーケンスされており、「3 pMのものが5 samples vs. 5 samples」、「1.5 pMのものが5 samples vs. 5 samples」の2つのデータセットに発現データがあります。

7列目 : R1L
8列目 : R1L
9列目 : R1L
10列目 : R1
11列目 : R1
12列目 : R1
13列目 : R1



34. k-mer解析用のランダム配列から生成したFASTA形式ファイル([sample34_ref.fasta](#)と[sample34_ngs.fasta](#))です。1000塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を500リード分だけランダム抽出したものです。塩基の存在比はAが22%, Cが28%, Gが28%, Tが22%にしています。リファレンス配列(仮想ゲノム配列)が[sample34_ref.fasta](#)で、500リードからなる仮想NGSデータが[sample34_ngs.fasta](#)です。リード長20塩基で500リードなのでトータル10,000塩基となり、1,000塩基からなる元のゲノム配列の10倍シーケンスしていることとなります(10X coverageに相当)。[イントロ | NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成から](#)と基本的に同じです。

```

out_f1 <- "sample34_ref.fasta" #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample34_ngs.fasta" #出力ファイル名を指定してout_f2に格納
param_len_ref <- 1000 #リファレンス配列の長さを指定
narabi <- c("A", "C", "G", "T") #以下の数値指定時にACGTの並びを間違えないようにするため
param_composition <- c(22, 28, 28, 22) #(A,C,G,Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20 #リード長を指定
param_num_ngs <- 500 #リード数を指定
param_desc <- "kkk" #FASTA形式ファイルのdescription行に記述する内容

#必要なパッケージをロード
library(Biostrings) #パッケージの読み込み

#本番(リファレンス配列生成)
set.seed(1010) #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition) #narabi中の塩基がparam_compositionで指定した数だけ
    
```

ランダムデータ生成

①例題34。②1000 bpの仮想リファレンスゲノム配列、および③仮想NGSデータを作成する。20 bpのリードが500個なので、計10,000 bp、つまり10X coverageのデータを作成していることになる

34. k-mer解析用のランダム配列から生成したFASTA形式ファイル(sample34_ref.fasta)と、500リードからなる仮想NGSデータ(sample34_ngs.fasta)を作成する。塩基の存在比はAが22%、Cが28%、Gが28%、Tが22%にしています。リファレンス配列(仮想ゲノム配列)がsample34_ref.fastaで、500リードからなる仮想NGSデータがsample34_ngs.fastaです。リード長20塩基で500リードなのでトータル10,000塩基となり、1,000塩基からなる元のゲノム配列の10倍シーケンスしていることとなります(10X coverageに相当)。イントロ | NGS | 配列取得 | シミュレーションデータ | ランダムな塩基配列の生成からと基本的に同じです。

```

out_f1 <- "sample34_ref.fasta" #出力ファイル名を指定してout_f1に格納
out_f2 <- "sample34_ngs.fasta" #出力ファイル名を指定してout_f2に格納
param_len_ref <- 1000 #リファレンス配列の長さを指定
narabi <- c("A","C","G","T") #以下の数値指定時にACGTの並びを間違えないようにするための並び
param_composition <- c(22, 28, 28, 22) #(A,C,G,Tの並びで)各塩基の存在比率を指定
param_len_ngs <- 20 #リード長を指定
param_num_ngs <- 500 #リード数を指定
param_desc <- "kkk" #FASTA形式ファイルのdescription行に記述する内容

#必要なパッケージをロード
library(Biostrings) #パッケージの読み込み

#本番(リファレンス配列生成)
set.seed(1010) #おまじない(同じ乱数になるようにするため)
ACGTset <- rep(narabi, param_composition)#narabi中の塩基がparam_compositionで指定した数だけ
reference <- paste(sample(ACGTset, param_len_ref, replace=T), collapse="")#ACGTsetからparam_len_ref個の塩基をランダムに抽出して連結
reference <- DNASTringSet(reference) #DNA塩基配列だと認識させるDNASTringSet関数を適用した
names(reference) <- param_desc #description行に相当する記述を追加している
reference #確認してるだけです

#本番(シミュレーションデータ生成)
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs, replace=T)#部分塩基配列抽出

```

ランダムデータ生成

34. k-mer解析用のランダム配列から生成したFASTA形式ファイル([sample34_ref.fasta](#)と[sample34_ngs.fasta](#))です。
 1000塩基の長さのリファレンス配列を生成したのち、20塩基長の部分配列を500リード分だけランダム抽出したものです。塩基の存在比はAが22%, Cが28%, Gが28%, Tが22%にしています。リファレンス配列(仮想ゲノム配列)が[sample34_ref.fasta](#)で、500リードからなる仮想NGSデータが[sample34_ngs.fasta](#)です。リード長20塩基で500リードなのでトータル10,000塩基となり、1,000塩基からなる元のゲノム配列の10倍シーケンスしていることとなります(10X coverageに相当)。[イントロ](#) | [NGS](#) | [配列取得](#) | [シミュレーションデータ](#) | [ランダムな塩基配列の生成から](#)と基本的に同じです。

```
out_f1 <- "sample34_ref.fasta"
out_f2 <- "sample34_ngs.fasta"
param_len_ref <- 1000
narabi <- c("A", "C", "G", "T")
param_composition <- c(22, 28, 28, 22)
param_len_ngs <- 20
param_num_ngs <- 500
param_desc <- "kkk"
```

#必要なパッケージをロード
 library(Biostrings)

#本番(リファレンス配列生成)

```
set.seed(1010)
ACGTset <- rep(narabi, param_composition)
reference <- paste(sample(ACGTset, param_len_ref), collapse="")
reference <- DNASTringSet(reference)
names(reference) <- param_desc
reference
```

#本番(シミュレーションデータ生成)

```
s_posi <- sample(1:(param_len_ref-param_len_ngs), param_num_ngs)
```

```
R Console
> fasta #確認してるだ$
A DNASTringSet instance of length 500
width seq names $
[1] 20 GGTACGGTTCCGGTTGCCGA kkk_658_677
[2] 20 TTATCCGGCAGTCTATATG kkk_197_216
[3] 20 CAATAGACACCACGCGCCAC kkk_559_578
[4] 20 ACCGAAAGTTCTGGAGGTCA kkk_55_74
[5] 20 TTGGGCCCGATATCTATAAA kkk_681_700
...
[496] 20 GAGCTATCAATGTACCCCAG kkk_433_452
[497] 20 AACCTAATCGCCCCGTGGCG kkk_805_824
[498] 20 CGCCGCCTTCGAGCTCGCTT kkk_140_159
[499] 20 GACAGCGGGAGTGATCCGGC kkk_401_420
[500] 20 TGTAAGCCGTACCCTACCCC kkk_904_923
>
> #ファイルに保存(仮想リファレンス配列と仮想NGS配列)
> writeXStringSet(reference, file=out_f1, format="fasta$
> writeXStringSet(fasta, file=out_f2, format="fasta", w$
> |
```

k-mer解析(k=10)

- ・ [イントロ](#) | [一般](#) | [逆鎖\(reverse\)を取得](#) (last modified 2013/06/14)
- ・ [イントロ](#) | [一般](#) | [k-mer解析](#) | [k=1\(塩基ごとの出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/04/27) **NEW**
- ・ [イントロ](#) | [一般](#) | [k-mer解析](#) | [k=2\(2連続塩基の出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/01/28)
- ・ [イントロ](#) | [一般](#) | [k-mer解析](#) | [k=3\(3連続塩基の出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/01/28)
- ・ [イントロ](#) | [一般](#) | [k-mer解析](#) | [k=n\(n連続塩基の出現頻度解析\)](#) | [Biostrings](#) (last modified 2016/04/30) **NEW**
- ・ (削除予定) [イントロ](#) | [一般](#) | [2連続塩基の出現頻度情報取得](#) (last modified 2015/04/20)
- ・ (削除予定) [イントロ](#) | [一般](#) | [2連続塩基の出現頻度情報取得](#) (last modified 2015/02/19)

イントロ | 一般 | k-mer解析 | k=n(n連続塩基の出現頻度解析) | Biostrings **NEW**

[Biostrings](#)パッケージを用いて、 4^n 通りの任意のn連続塩基の出現頻度を調べるやり方を示します。k-mer解析のk=nの場合に相当します。例題12. サンプルデータの例題34を実行して得られたmulti-FASTAファイル(sample34_ngs.fasta)の場合:

1. イントロ | 一般 |

4連続塩基(k=4)

```
in_f <- "hog
out_f <- "ho
param_kmer <

#必要なパッケージ
library(Bios
```

```
#入力ファイル
fasta <- rea
```

12. サンプルデータの例題34を実行して得られたmulti-FASTAファイル(sample34_ngs.fasta)の場合:

10連続塩基(k=10)の出現頻度情報を得るやり方です。 $4^{10} = 1,048,576$ 通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を出力するやり方です。

```
in_f <- "sample34_ngs.fasta"
out_f <- "hoge12.txt"
param_kmer <- 10
```

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

```
#必要なパッケージをロード
library(Biostrings)
```

#パッケージの読み込み

```
#入力ファイルの読み込み
```

```
fasta <- readDNAStringSet(in_f, format="fasta")#in_fで指定したファイルの読み込み
```

```
#本番
```

```
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)#k連続塩基の出現頻度情報をhogeに格納
out <- colSums(hoge) #列ごとの総和をoutに格納
```

```
#ファイルに保存
```

```
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=T, col.names=F)#outの中身を出力
length(out) #4^param_kmerの値を表示
sum(out > 0) #1回以上出現したk-merの種類数を表示
```


k-mer解析 (k=10)

①例題12は、②k=10のときの、 $4^{10} = 1,048,576$ 通りの10連続塩基ごとの出現頻度を返すコード。③1回以上出現している10連続塩基は985個。つまり10X coverageのときの推定ゲノムサイズは985 bp。4X coverageのときの推定値(904 bp)に比べて、より正解(1,000 bp)に近い値になっていることがわかる

12. **サンプルデータの例題34**を実行して得られたmulti-FASTAファイルから10連続塩基(k=10)の出現頻度情報を得るやり方です。 $4^{10} = 1,048,576$ 通りです。全リードを合算した出現頻度を出力するやり方です。

```

in_f <- "sample34_ngs.fasta"
out_f <- "hoge12.txt"
param_kmer <- 10

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")

#本番
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)

#ファイルに保存
write.table(out, out_f, sep="\t", append=F, quote=F, length(out))
sum(out > 0)

```

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

```

R Console
> param_kmer <- 10 #k-merのkの値$
> 
> #必要なパッケージをロード
> library(Biostrings) #パッケージの$
> 
> #入力ファイルの読み込み
> fasta <- readDNASTringSet(in_f, format="fasta") #in_f$
> 
> #本番
> hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
> out <- colSums(hoge) #列ごとの総和$
> 
> #ファイルに保存
> write.table(out, out_f, sep="\t", append=F, quote=F, $
> length(out) #4^param_kmer$
[1] 1048576
> sum(out > 0) #1回以上出現し$
[1] 985
> |

```

Contents

- NGS解析手段、ウェブツール(DDBJ Pipeline)との連携
- DDBJ PipelineでPlatanusを実行
- k-mer解析(k個の連続塩基に基づく各種解析)の基礎
 - 塩基ごとの出現頻度解析(k=1)、一気に計算
 - 2連続塩基の出現頻度解析(k=2)、基本スキルの復習や作図
- de novoアセンブリ時のエラー補正やゲノムサイズ推定の基本的な考え方
 - ランダムな塩基配列(仮想ゲノムおよび仮想NGSデータ)の作成
 - k-mer解析の応用、ゲノムサイズ推定の基礎
 - ゲノムサイズ推定(1,000 bpの仮想ゲノムの場合)
 - ゲノムサイズ推定(1,000 bpの仮想ゲノム; 4X → 10X coverageの場合)
 - k-mer出現頻度分布
 - シークエンスエラーを含む場合
 - 最終確認



k-mer出現頻度分布

①outベクトルの要素数は1,048,576個で、ほとんどの要素(1,048,576 - 985)が0。②1回以上出現している10連続塩基(10-mer)のみで取り扱いたい場合は、out[out > 0]とすればよい

12. サンプルデータの例題34を実行して得られたmulti-FASTAファイル(sample34

10連続塩基(k=10)の出現頻度情報を得るやり方です。4¹⁰ = 1,048,576通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を出力するやり方です。

```

in_f <- "sample34_ngs.fasta"
out_f <- "hoge12.txt"
param_kmer <- 10

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNAStringSet(in_f, format="fastq")

#本番
hoge <- oligonucleotideFrequency(fasta,
out <- colSums(hoge)

#ファイルに保存
write.table(out, out_f, sep="\t", append=F,
length(out)
sum(out > 0)

```

```

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

```

```

R Console
> out <- colSums(hoge) #列ごとの総和$
>
> #ファイルに保存
> write.table(out, out_f, sep="\t", append=F, quote=F, $
> length(out) #4^param_kmer$
[1] 1048576
> sum(out > 0) #1回以上出現し$
[1] 985
① > head(out)
AAAAA AAAAAAAC AAAAAAAG AAAAAAAT AAAAAACA
      0          0          0          0          0
AAAAAAC
      0
② > head(out[out > 0])
AAACTTGGC AAAAGTTTT AAACAGCTGC AAACGGGAGC AAACTTGGCG
      8          6          7          2          7
AAAGAGCTAT
      8
> |

```

①1回以上出現している10連続塩基(10-mer)のみをkmerとして取り扱う。②当然kmerの要素数は985

k-mer出現頻度分布

12. サンプルデータの例題34を実行して得られたmulti-FASTAファイル(sample34_ngs.fasta)の場合:

10連続塩基(k=10)の出現頻度情報を得るやり方です。4¹⁰ = 1,048,576通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を出力するやり方です。

```

in_f <- "sample34_ngs.fasta"
out_f <- "hoge12.txt"
param_kmer <- 10

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNAStringSet(in_f, format="fastq")

#本番
hoge <- oligonucleotideFrequency(fasta, param_kmer)
out <- colSums(hoge)

#ファイルに保存
write.table(out, out_f, sep="\t", append=TRUE,
            length(out),
            sum(out > 0))

```

```

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

```

```

R Console
> head(out)
AAAAA      AAAAAAAC  AAAAAAAG  AAAAAAAT  AAAAAACA
0          0          0          0          0
AAAAAAACC
0
> head(out[out > 0])
AAACTTGGC  AAAAGTTTT  AAACAGCTGC  AAACGGGAGC  AAACTTGGCG
8          6          7          2          7
AAAGAGCTAT
8
① > kmer <- out[out > 0]
> head(kmer)
AAACTTGGC  AAAAGTTTT  AAACAGCTGC  AAACGGGAGC  AAACTTGGCG
8          6          7          2          7
AAAGAGCTAT
8
② > length(kmer)
[1] 985
> |

```

k-mer出現頻度分布

k-mer出現頻度分布は、①hist(kmer)を実行して得られる、②横軸が出現回数、縦軸が出現頻度からなるヒストグラムのこと

12. サンプルデータの例題34を実行して得られたmulti-FASTAファイル(sample34_ngs.fasta)の場合:

10連続塩基(k=10)の出現頻度情報を得るやり方です。4¹⁰ = 1,048,576通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を出力するやり方です。

```

in_f <- "sample34_ngs.fasta"
out_f <- "hoge12.txt"
param_kmer <- 10

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNAStringSet(in_f, format="fasta")

#本番
hoge <- oligonucleotideFrequency(fasta, k=param_kmer)
out <- colSums(hoge)

#ファイルに保存
write.table(out, out_f, sep="\t", append=TRUE,
            length(out),
            sum(out > 0))

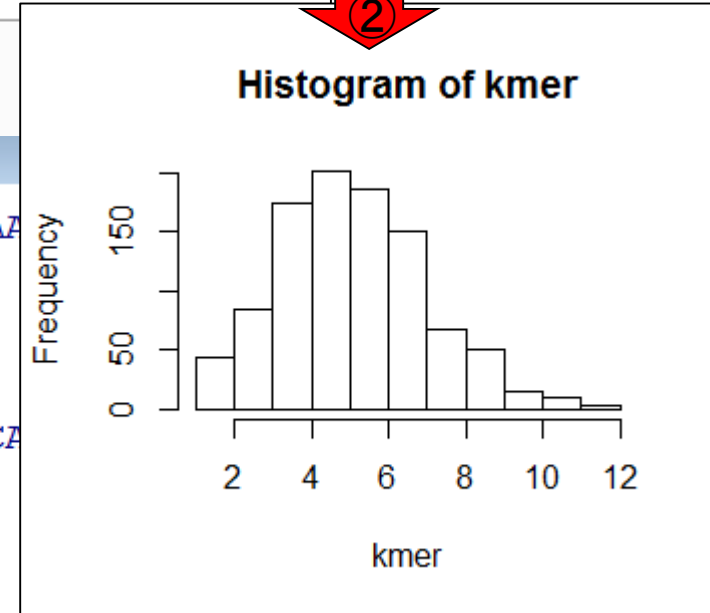
```

```

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

R Console
# AAAAAAAAAA AAAAAAAAAAC AAAAA
#           0           0
# AAAAAAAAAACC
#           0
# > head(out[out > 0])
# AAAACTTGGC AAAAGTTTTT AAACA
#           8           6
# AAAGAGCTAT
#           8
# > kmer <- out[out > 0]
# > head(kmer)
# AAAACTTGGC AAAAGTTTTT AAACAGCTGC AAACGGGAGC AAACTTGGCG
#           8           6           7           2           7
# AAAGAGCTAT
#           8
# > length(kmer)
# [1] 985
# > hist(kmer)
# > |

```



k-mer出現頻度分布

①ヒストグラムでピンとこないヒトは、① table(kmer)と同じものという理解でよい。例えば ②は、12回出現したk-merが3個(3種類)あったということ。気になるヒトは③のようにして確認

12. サンプルデータの例題34を実行して得られたmulti-FASTAファイル(sample34.fasta)の出現頻度情報を得るやり方です。4¹⁰ = 1,048,576通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を入力するやり方です。

10連続塩基(k=10)の出現頻度情報を得るやり方です。4¹⁰ = 1,048,576通りのk-merの出現頻度を計算することになります。全リードを合算した出現頻度を入力するやり方です。

```
in_f <- "sample34_ngs.fasta"
out_f <- "hoge12.txt"
param_kmer <- 10

#必要なパッケージをロード
library(Biostrings)

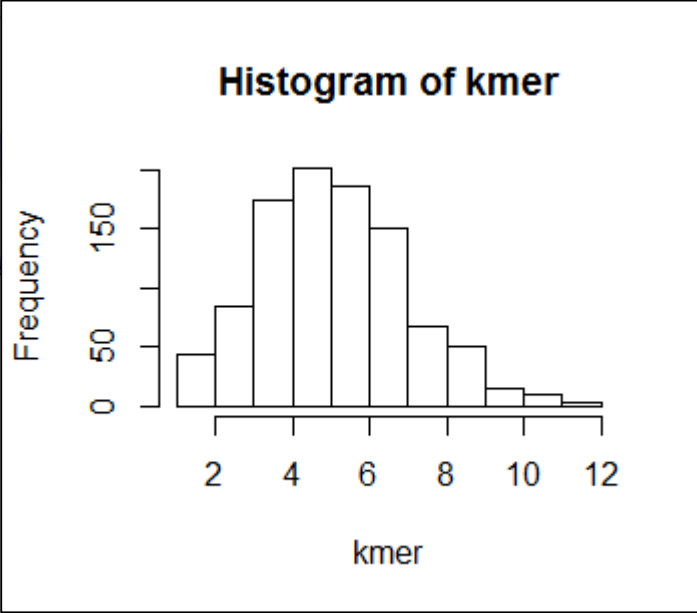
#入力ファイルの読み込み
fasta <- readDNAStringSet(in_f, format="fasta")

#本番
hoge <- oligonucleotideFrequency(fasta, k=param_kmer)
out <- colSums(hoge)

#ファイルに保存
write.table(out, out_f, sep="\t", append=TRUE)
length(out)
sum(out > 0)
```

```
#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

R Console
> head(kmer)
AAACTTGGC AAAAGTTTT AAACA
      8          6
AAAGAGCTAT
      8
> length(kmer)
[1] 985
> hist(kmer)
> table(kmer)
kmer
 1   2   3   4   5   6   7   8   9  10  11  12
3  40  84 174 201 186 151  67  51  15  10   3
> kmer[kmer == 1]
GACGTGTAAG TAAACGGGAG TGACGTGTAA
      1          1          1
> kmer[kmer == 12]
CACACGACCG CACCCTCAGA TCACACGACC
      12          12          12
> |
```



k-mer出現頻度分布

ちなみに②k-mer出現回数の中央値(median)は5。③入力ファイルは、ゲノムサイズ1,000 bpのリファレンス配列から、20 bpのリードを500個ランダム抽出したもの。計10,000 bpなので、10X coverageのデータ。おかしくない?

12. サンプルデータの例題34を実行して得られたmulti-FASTAファイル(sample34

10連続塩基(k=10)の出現頻度情報を得るやり方です。4¹⁰ = 1,048,576通りのk-merがあります。全リードを合算した出現頻度を出力するやり方です。

```
in_f <- "sample34_ngs.fasta"
out_f <- "hoge12.txt"
param_kmer <- 10

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNAStringSet(in_f, format="fasta")

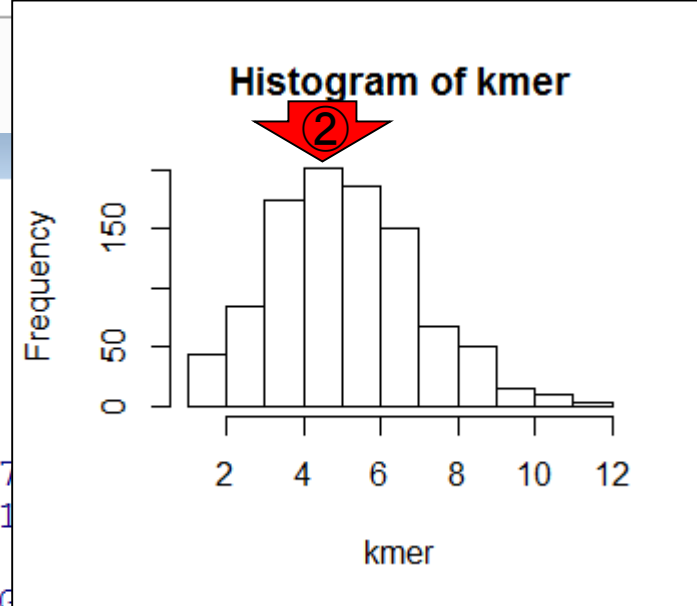
#本番
hoge <- oligonucleotideFrequency(fasta, k=param_kmer)
out <- colSums(hoge)

#ファイルに保存
write.table(out, out_f, sep="\t", append=TRUE)
length(out)
sum(out > 0)
```



#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#k-merのkの値を指定

```
R Console
> length(kmer)
[1] 985
> hist(kmer)
> table(kmer)
kmer
 1  2  3  4  5  6  7
3 40 84 174 201 186 151
> kmer[kmer == 1]
GACGTGTAAA TAAACGGGAG TGACCGGAGG
 1 1 1
> kmer[kmer == 12]
CACACGACCG CACCCTCAGA TCACACGACC
 12 12 12
> mean(kmer)
[1] 5.583756
> median(kmer)
[1] 5
> |
```



k=6, 8, 10, 12

k値の違いによるk-mer出現頻度分布の傾向を把握。①k=6のときの②k-mer出現回数の中央値(median)は8。③最多出現回数は39回。作業としては、赤枠部分のみコピペしていくのが基本だが、全部一度にコピペ実行してもよい。尚、k=6, 8, 10まではエラーは出ないが、12でエラーが出る(後述)

• k=6, 8, 10, 12

必要最小限のコードにして、k値の違い(k=6, 8, 10, and 12)による影響の全体像を

```
in_f <- "sample34_ngs.fasta" #入力ファイル名を指定し
library(Biostrings) #パッケージの読み込み
fasta <- readDNASTringSet(in_f, format="fasta")#in_fで指定し
```

```
param_kmer <- 6
hoge <- oligonucleotideFrequency(fasta,
out <- colSums(hoge)
kmer <- out[out > 0]
length(kmer)
table(kmer)
median(kmer)
```

```
param_kmer <- 8
hoge <- oligonucleotideFrequency(fasta,
out <- colSums(hoge)
kmer <- out[out > 0]
length(kmer)
table(kmer)
median(kmer)
```

```
param_kmer <- 10
hoge <- oligonucleotideFrequency(fasta,
```

```
R Console
> in_f <- "sample34_ngs.fasta" #入力ファイル$
> library(Biostrings) #パッケージの$
> fasta <- readDNASTringSet(in_f, format="fasta")#in_f$
>
> param_kmer <- 6 #k-merのkの値$
> hoge <- oligonucleotideFrequency(fasta, width=param_k$
> out <- colSums(hoge) #列ごとの総和$
> kmer <- out[out > 0] #1回以上出現し$
> length(kmer) #1回以上出現し$
[1] 862
> table(kmer) # (k-merの種類)$
kmer
 2  3  4  5  6  7  8  9 10 11 12 13 14 15
 5 10 36 74 109 134 146 112 68 43 29 11 17 11
16 17 18 19 20 21 22 24 26 27 28 29 30 39
13 7 10 4 6 4 4 1 2 2 1 1 1 1
> median(kmer) #出現回数の中$
[1] 8
>
```


k=6, 8, 10, 12

k値の違いによるk-mer出現頻度分布の傾向を把握。①k=8のときの②k-mer出現回数の中央値(median)は7。③最多出現回数は16回

- k=6, 8, 10, 12

必要最小限のコードで、k値の違い(k=6, 8, 10, and 12)による影響の全体像を把握。

```
param_kmer <- 8
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)
kmer <- out[out > 0]
length(kmer)
table(kmer)
median(kmer)
```

```
param_kmer <- 10
hoge <- oligonucleotideFrequency(fasta,
out <- colSums(hoge)
kmer <- out[out > 0]
length(kmer)
table(kmer)
median(kmer)
```

```
param_kmer <- 12
hoge <- oligonucleotideFrequency(fasta,
out <- colSums(hoge)
kmer <- out[out > 0]
length(kmer)
table(kmer)
```

```
#k-merのkの値を指定
#k連続塩基の出現頻度
#列ごとの総和をoutに格納
#1回以上出現したk-merのみをkmerに格納
#1回以上出現したk-merの種類数を表示
```

```
R Console
1  2  3  4  5  6  7  8  9 10 11 12
3 40 84 174 201 186 151 67 51 15 10 3
> median(kmer)
[1] 5
> param_kmer <- 8
> hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
> out <- colSums(hoge)
> kmer <- out[out > 0]
> length(kmer)
[1] 975
> table(kmer)
# (k-merの種類$
kmer
 2  3  4  5  6  7  8  9 10 11 12 13 14 15
18 25 99 142 185 178 148 78 47 22 12 4 4 1
16
 2
> median(kmer)
[1] 7
>
```

k=6, 8, 10, 12

k値の違いによるk-mer出現頻度分布の傾向を把握。①k=10のときの②k-mer出現回数の中央値(median)は5。③最多出現回数は12回

- k=6, 8, 10, 12

必要最小限のコードにして、k値の違い(k=6, 8, 10, and 12)による影響の全体像を把握。

```
out <- colSums(hoge)
kmer <- out[out > 0]
length(kmer)
table(kmer)
median(kmer)
```

①

```
param_kmer <- 10
hoge <- oligonucleotideFrequency(fasta,
out <- colSums(hoge)
kmer <- out[out > 0]
length(kmer)
table(kmer)
median(kmer)
```

```
param_kmer <- 12
hoge <- oligonucleotideFrequency(fasta,
out <- colSums(hoge)
kmer <- out[out > 0]
length(kmer)
table(kmer)
median(kmer)
```

#列ごとの総和をoutに格納
#1回以上出現したk-merのみをkmerに格納
#1回以上出現したk-merの種類数を表示
#(k-merの種類は問わずに)k-merの出現回数

R Console

```
2 3 4 5 6 7 8 9 10 11 12 13 14 15
18 35 99 142 185 178 148 78 47 22 12 4 4 1
16
2
```

> median(kmer) #出現回数の中\$

[1] 7

```
> param_kmer <- 10 #k-merのkの値$
> hoge <- oligonucleotideFrequency(fasta, width=param_k$
> out <- colSums(hoge) #列ごとの総和$
> kmer <- out[out > 0] #1回以上出現し$
> length(kmer) #1回以上出現し$
```

[1] 985

> table(kmer) #(k-merの種類\$

```
kmer
1 2 3 4 5 6 7 8 9 10 11 12
3 40 84 174 201 186 151 67 51 15 10 3
```

> median(kmer) #出現回数の中\$

[1] 5

>

②

③

k=6, 8, 10, 12

k値の違いによるk-mer出現頻度分布の傾向を把握。①k=12のときはメモリ不足の問題でエラーが出ていることがわかる。それゆえ、後半の結果はk=10のときのものがそのまま表示されている点に注意!

• k=6, 8, 10, 12

必要最小限のコードにして、k値の違い(k=6, 8, 10, and 12)による影響の全体像を

```
out <- colSums(hoge)
kmer <- out[out > 0]
length(kmer)
table(kmer)
median(kmer)

param_kmer <- 10
hoge <- oligonucleotideFrequency(fasta,
out <- colSums(hoge)
kmer <- out[out > 0]
length(kmer)
table(kmer)
median(kmer)
```

①

```
param_kmer <- 12
hoge <- oligonucleotideFrequency(fasta,
out <- colSums(hoge)
kmer <- out[out > 0]
length(kmer)
table(kmer)
median(kmer)
```

②

```
R Console
> param_kmer <- 12 #k-merのkの値$
> hoge <- oligonucleotideFrequency(fasta, width=param_k$
エラー: サイズ 31.2 Gb のベクトルを割り当てることができ$
追加情報: 警告メッセージ:
1: .Call2("XStringSet_oligo_frequency", x, width, step$
Reached total allocation of 8106Mb: see help(memory.s$
2: .Call2("XStringSet_oligo_frequency", x, width, step$
Reached total allocation of 8106Mb: see help(memory.s$
3: .Call2("XStringSet_oligo_frequency", x, width, step$
Reached total allocation of 8106Mb: see help(memory.s$
4: .Call2("XStringSet_oligo_frequency", x, width, step$
Reached total allocation of 8106Mb: see help(memory.s$
> out <- colSums(hoge) #列ごとの総和$
> kmer <- out[out > 0] #1回以上出現し$
> length(kmer) #1回以上出現し$
[1] 985
> table(kmer) #(k-merの種類$
kmer
 1  2  3  4  5  6  7  8  9 10 11 12
 3 40 84 174 201 186 151 67 51 15 10 3
> median(kmer) #出現回数の中$
[1] 5
> |
```

作図(k=6, 8, 10)

k値が小さくなる(図の下から上)にしたがって、分布が右にシフトしていることがわかる。つまり、出現回数の中央値は、k-merのk値次第だということ。横軸の範囲を揃えないとわからないので、①1-20に揃えるべく、②のxlimオプション部分で指定

作図(k=6, 8, 10)

上記のk=6, 8, 10の結果より、横軸の最大値は39であることがわかっていて、遷が把握しづらいので1から20の範囲を指定している。このあたりは結果合わせて微調整する。

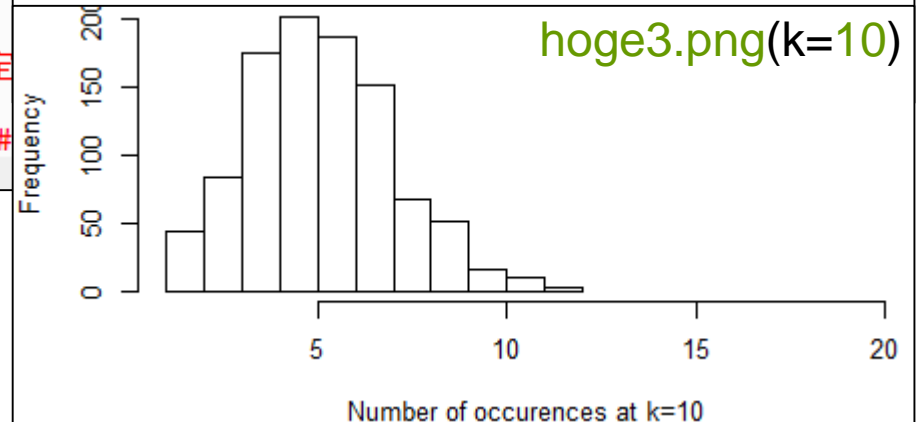
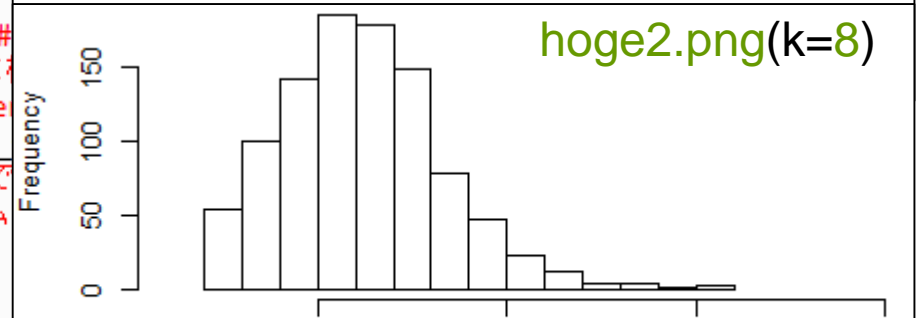
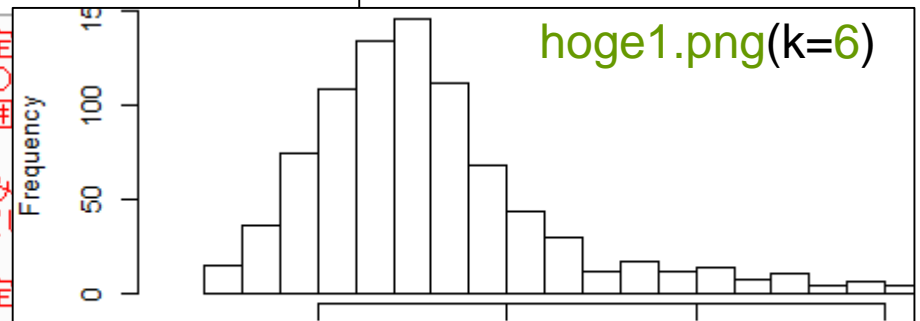
```
in_f <- "sample34_ngs.fasta" #入力ファイル名を指定
param_fig <- c(450, 210) #ヒストグラム描画時の
param_range <- c(1, 20) #表示したい横軸の範囲 ①

library(Biostrings) #パッケージの読み込み
fasta <- readDNAStringSet(in_f, format="fasta") #in_fで指定

out_f <- "hoge1.png" #出力ファイル名を指定
param_kmer <- 6 #k-merのkの値を指定
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #
out <- colSums(hoge) #列ごとの総和をoutに
kmer <- out[out > 0] #1回以上出現したk-mer

png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],
par(mar=c(4, 4, 0, 0)) #下、左、上、右の順で
hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム
xlab=paste("Number of occurrences at k=", param_kmer,
xlim=param_range, main="")) #ヒストグラムを描画
dev.off() #おまじない

out_f <- "hoge2.png" #出力ファイル名を指定
param_kmer <- 8 #k-merのkの値を指定
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #
<
```



おさらい

k-mer出現頻度解析は、このように1塩基づつずらして部分配列を発生させて出現回数を解析しています。L塩基長のリードをk-merで分割すると、 $(L - k + 1)$ 個のk-merを発生可能

k=6の場合: 15個のk-mer

GGTACGGTTCGGTTGCCGA

```
GGTACG
GTACGG
TACGGT
ACGGTT
CGGTTC
GGTTCC
GTTCCG
TTCCGG
TCCGGT
CCGGTT
CGGTTG
GGTTGC
GTTGCC
TTGCCG
TGCCGA
```

k=8の場合: 13個のk-mer

GGTACGGTTCGGTTGCCGA

```
GGTACGGT
GTACGGTT
TACGGTTC
ACGGTTCC
CGGTTCGG
GGTTCGGG
GTTCCGGT
TTCCGGTT
TCCGGTTG
CCGGTTGC
CGGTTGCC
GGTTGCCG
GTTGCCGA
```

k=10の場合: 11個のk-mer

GGTACGGTTCGGTTGCCGA

```
GGTACGGTTC
GTACGGTTCC
TACGGTTCCG
ACGGTTCCGG
CGGTTCGGGT
GGTTCGGGTT
GTTCCGGTTG
TTCCGGTTGC
TCCGGTTGCC
CCGGTTGCCG
CGGTTGCCGA
```

Contents

- NGS解析手段、ウェブツール(DDBJ Pipeline)との連携
- DDBJ PipelineでPlatanusを実行
- k-mer解析(k個の連続塩基に基づく各種解析)の基礎
 - 塩基ごとの出現頻度解析(k=1)、一気に計算
 - 2連続塩基の出現頻度解析(k=2)、基本スキルの復習や作図
- de novoアセンブリ時のエラー補正やゲノムサイズ推定の基本的な考え方
 - ランダムな塩基配列(仮想ゲノムおよび仮想NGSデータ)の作成
 - k-mer解析の応用、ゲノムサイズ推定の基礎
 - ゲノムサイズ推定(1,000 bpの仮想ゲノムの場合)
 - ゲノムサイズ推定(1,000 bpの仮想ゲノム; 4X → 10X coverageの場合)
 - k-mer出現頻度分布
 - シークエンスエラーを含む場合
 - 最終確認



作図 (k=6, 8, 10)

①は1000 bpの仮想リファレンス配列からランダム抽出した、長さ20 bpのリードが500個からなるシーケンスエラーのない仮想データ

作図(k=6, 8, 10)

上記のk=6, 8, 10の結果より、横軸の最大値は39であることがわかっているが、ヒストグラム全体の変遷が把握しづらいので1から20の範囲を指定している。このあたりは結果を眺めながら自分の好みに合わせて微調整する。

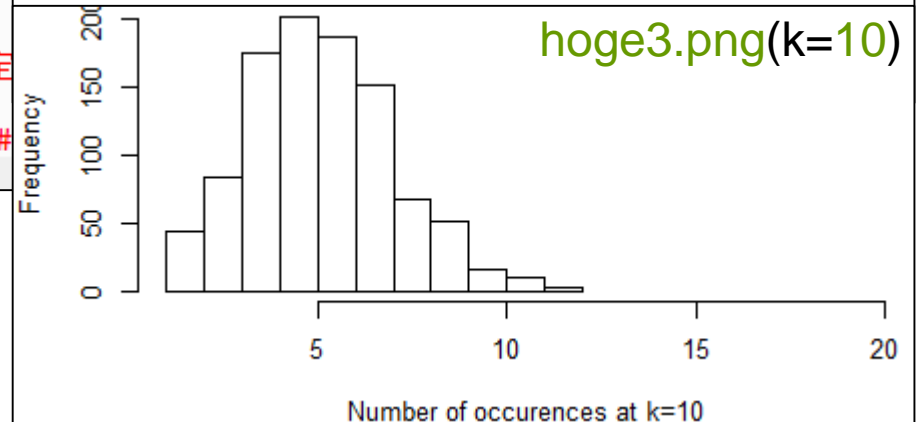
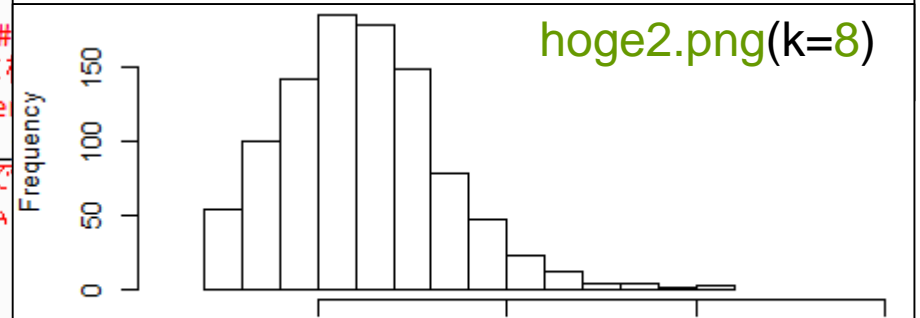
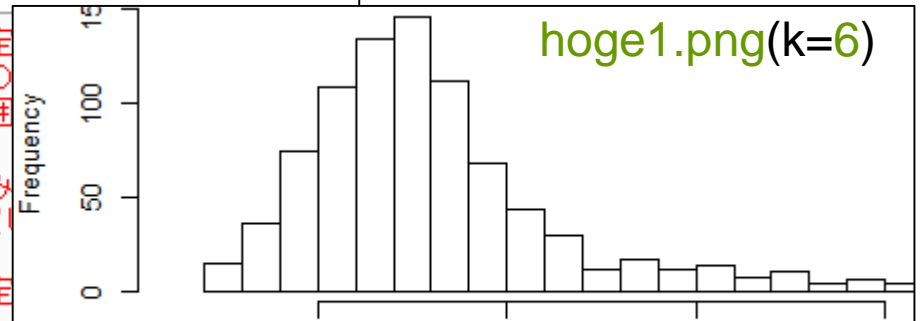
```
in_f <- "sample34_ngs.fasta" ① #入力ファイル名を指定
param_fig <- c(450, 210) #ヒストグラム描画時の
param_range <- c(1, 20) #表示したい横軸の範囲

library(Biostrings) #パッケージの読み込み
fasta <- readDNAStringSet(in_f, format="fasta") #in_fで指定

out_f <- "hoge1.png" #出力ファイル名を指定
param_kmer <- 6 #k-merのkの値を指定
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #
out <- colSums(hoge) #列ごとの総和をoutに
kmer <- out[out > 0] #1回以上出現したk-mer

png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],
par(mar=c(4, 4, 0, 0)) #下、左、上、右の順で
hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム
xlab=paste("Number of occurrences at k=", param_kmer,
xlim=param_range, main="")) #ヒストグラムを描画
dev.off() #おまじない

out_f <- "hoge2.png" #出力ファイル名を指定
param_kmer <- 8 #k-merのkの値を指定
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #
```



塩基置換

赤枠内をコピーした結果のR Console画面。①は長さ20 bpのリードが500個からなるシークエンスエラーのない仮想データ。500個の全リードについて、②18番目のポジションに1塩基だけ「TならA、CならG」となるような塩基置換(仮想シークエンスエラー)を入れて、k-mer出現頻度分布がどうなるかをみる

塩基置換

GならC、AならTのような塩基置換を行うenkichikan関数の塩基置換を行う。「書籍 | トランスクリプトーム解析 | 2.」を参照してください。

```
in_f <- "sample34_ngs.fasta"
param_fig <- c(450, 210)
param_range <- c(1, 20)

library(Biostrings)
fasta <- readDNAStringSet(in_f, format="fasta")
fasta

#塩基置換関数の作成
enkichikan <- function(fa, p) {
  t <- substring(fa, p, p)
  t_c <- chartr("CGAT", "GCTA", t)
  substring(fa, p, p) <- t_c
  return(fa)
}

#enkichikan関数の動作確認
as.character(fasta[500])
enkichikan(as.character(fasta[500]), 18)

#塩基置換本番
fasta_org <- fasta
```

#入力ファイル名を指定してin_fに格納
#ヒストグラム描画時の横幅と縦幅を指定
#表示したk-merの範囲を指定

#パッケージ
#確認して

#関数名や
#置換した
#置換後の
#置換
#置換後の

#500番目
#18番目

#塩基置換

```
R Console
[499] 20 GACAGCGGGAGTGATCCGGC kkk_401_420
[500] 20 TGTAAGCCGTACCCTACCCC kkk_904_923
>
> #塩基置換関数の作成
> enkichikan <- function(fa, p) {
+   t <- substring(fa, p, p)
+   t_c <- chartr("CGAT", "GCTA", t)
+   substring(fa, p, p) <- t_c
+   return(fa)
+ }
>
> #enkichikan関数の動作確認
> as.character(fasta[500])
      kkk_904_923
"TGTAAGCCGTACCCTACCCC"
> enkichikan(as.character(fasta[500]), 18)
[1] "TGTAAGCCGTACCCTACGCC"
> |
```


①はfastaオブジェクトの500番目の要素を文字列として (as.character)表示させた結果。つまり、②と同じもの

塩基置換

塩基置換
GならC、AならTのような塩基置換を行うenkichikan関数(孫 建強氏作成)を利用して任意の位置の塩基置換を行う。「書籍 | トランスクリプトーム解析 | 2.3.4 マッピング(準備)」のp81と同じ関数です。

```
in_f <- "sample34_ngs.fasta" #入力ファイル名を指定してin_fに格納
param_fig <- c(450, 210) #ヒストグラム描画時の横幅と縦幅を指定
param_range <- c(1, 20) #表示したヒストグラムの範囲を指定

library(Biostrings) #パッケージ
fasta <- readDNAStringSet(in_f, format="fasta") #確認して
fasta

#塩基置換関数の作成
enkichikan <- function(fa, p) { #関数名や
  t <- substring(fa, p, p) #置換した
  t_c <- chartr("CGAT", "GCTA", t) #置換後の
  substring(fa, p, p) <- t_c #置換
  return(fa) #置換後の
}

#enkichikan関数の動作確認
as.character(fasta[500]) #500番目
enkichikan(as.character(fasta[500]), 18) #18番目

#塩基置換本番
fasta_org <- fasta #塩基置換
```

```
R Console
[499] 20 GACAGCGGGAGTGTATCCGGC kkk_401_420
[500] 20 TGTAAGCCGTACCCTACCCC kkk_904_923
>
> #塩基置換関数の作成
> enkichikan <- function(fa, p) { #関数名$
+ t <- substring(fa, p, p) #置換し$
+ t_c <- chartr("CGAT", "GCTA", t) #置換後$
+ substring(fa, p, p) <- t_c #置換
+ return(fa) #置換後$
+ }
>
> #enkichikan関数の動作確認
> as.character(fasta[500]) #500$
kkk_904_923
"TGTAAGCCGTACCCTACCCC"
> enkichikan(as.character(fasta[500]), 18) #18番$
[1] "TGTAAGCCGTACCCTACGCC"
> |
```

①の部分で作成しておいたenkichikan関数を用いて、②18番目の塩基を置換。③確かにCがGに置換されている

塩基置換

- 塩基置換
GならC、AならTのような塩基置換を行うenkichikan関数(孫 建強氏作成)を利用して任意の位置の塩基置換を行う。「書籍 | トランスクリプトーム解析 | 2.3.4 マッピング(準備)」のp81と同じ関数です。

```
in_f <- "sample34_ngs.fasta" #入力ファイル名を指定してin_fに格納
param_fig <- c(450, 210) #ヒストグラム描画時の横幅と縦幅を指定
param_range <- c(1, 20) #表示した塩基の範囲を指定
```

```
library(Biostrings) #パッケージ
fasta <- readDNASTringSet(in_f, format="fasta"); #確認して
fasta
```

#塩基置換関数の作成

```
enkichikan <- function(fa, p) { #関数名や #置換した #置換後の #置換 #置換後の
  t <- substring(fa, p, p)
  t_c <- chartr("CGAT", "GCTA", t)
  substring(fa, p, p) <- t_c
  return(fa)
}
```

#enkichikan関数の動作確認

```
as.character(fasta[500]) #500番目
enkichikan(as.character(fasta[500]), 18) #18番目
```

#塩基置換本番

```
fasta_org <- fasta #塩基置換
```

```
R Console
[499] 20 GACAGCGGGAGTGTATCCGGC kkk_401_420
[500] 20 TGTAAGCCGTACCCTACCC kkk_904_923
>
> #塩基置換関数の作成
> enkichikan <- function(fa, p) { #関数名$
+ t <- substring(fa, p, p) #置換し$
+ t_c <- chartr("CGAT", "GCTA", t) #置換後$
+ substring(fa, p, p) <- t_c #置換
+ return(fa) #置換後$
+ }
> #enkichikan関数の動作確認
> as.character(fasta[500]) #500$
kkk_904_923
"TGTAAGCCGTACCCTACCC"
> enkichikan(as.character(fasta[500]), 18) #18番$
[1] "TGTAAGCCGTACCCTACGC"
> |
```

塩基置換

①コード下部に移動。②全リードについて18番目の塩基を置換したい場合には、入力部分をas.character(fasta)とすればよい

塩基置換

GならC、AならTのような塩基置換を行うenkichikan関数(孫建強氏作成)を利用して任意の位置の塩基置換を行う。「書籍 | トランスクリプトーム解析 | 2.3.4 マッピング(準備)」のp81と同じ関数です。

```
enkichikan(as.character(fasta[500]), 18) #18番目の塩基を置換

#塩基置換本番
fasta_org <- fasta #塩基置換前のfastaをfasta_orgに格納
hoge <- enkichikan(as.character(fasta), 18) #2番目の塩基を置換
fasta <- DNASTringSet(hoge) #DNA塩基配列だと認識させるDNASTringSet
names(fasta) <- names(fasta_org) #fasta_orgのdescription情報をコピー
fasta #確認してるだけです

#k-mer出現頻度分布
out_f <- "hoge2_18.png" #出力ファイル名を指定してout_fに格納
param_kmer <- 8 #k-merのkの値を指定
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #k連続塩基の出現歩
out <- colSums(hoge) #列ごとの総和をoutに格納
kmer <- out[out > 0] #1回以上出現したk-merのみをkmerに格納
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2]) #出力ファイル名を指定してout_fに格納
par(mar=c(4, 4, 0, 0)) #下、左、上、右の順で余白(行)を指定してpar(mar=c(4, 4, 0, 0))
hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラムを描画
      xlab=paste("Number of occurrences at k=", param_kmer, sep=""), #ヒストグラムを描画
      xlim=param_range, main="") #ヒストグラムを描画
dev.off() #おまじない
```



塩基置換

②実行後のhogeオブジェクトは、文字列ベクトル。③
それゆえDNAStrngSet関数を実行して、見慣れた
DNAStrngSet形式のfastaオブジェクトに変換している

塩基置換

GならC、AならTのような塩基置換を行うenkichikan関数(孫 建強氏作成)を利用して任意の位置の塩基置換を行う。「書籍 | トランスクリプトーム解析 | 2.3.4 マッピング(準備)」のp81と同じ関数です。

```
enkichikan(as.character(fasta[500]), 18) #18番目の塩基を置換
```

```
#塩基置換本番
```

```
fasta_org <- fasta  
hoge <- enkichikan(as.character(fasta), 18)  
fasta <- DNAStrngSet(hoge)  
names(fasta) <- names(fasta_org)  
fasta
```

```
#k-mer出現頻度分布
```

```
out_f <- "hoge2_18.png"  
param_kmer <- 8  
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)  
out <- colSums(hoge)  
kmer <- out[out > 0]  
png(out_f, pointsize=13, width=param_fig[1],  
par(mar=c(4, 4, 0, 0))  
hist(kmer, breaks=max(kmer), ylab="Frequency",  
xlab=paste("Number of occurrences at k=", param_kmer),  
xlim=param_range, main="")  
dev.off()
```

```
#塩基置換
```

```
#確認し
```

```
#出力ファイル名
```

```
#k-merの幅
```

```
#列ごとの出現回数
```

```
#1回以上出現するk-mer
```

```
#下、左
```

```
#ヒストグラムの幅
```

```
#おまじな
```

```
R Console  
> #塩基置換本番  
> fasta_org <- fasta #塩基置$  
> hoge <- enkichikan(as.character(fasta), 18) #18$  
> head(hoge)  
[1] "GGTACGGTTCCGGTTGCGGA" "TTATCCGGCAGTCCTATTTG"  
[3] "CAATAGACACCACGCGCGAC" "ACCGAAAGTTCTGGAGGACA"  
[5] "TTGGGCCCGATATCTATTAA" "TATGATATGGGCCGCTGGGG"  
> length(hoge)  
[1] 500  
> fasta <- DNAStrngSet(hoge) #DNA塩$  
> fasta  
A DNAStrngSet instance of length 500  
width seq  
[1] 20 GGTACGGTTCCGGTTGCGGA  
[2] 20 TTATCCGGCAGTCCTATTTG  
[3] 20 CAATAGACACCACGCGCGAC  
[4] 20 ACCGAAAGTTCTGGAGGACA  
[5] 20 TTGGGCCCGATATCTATTAA
```

塩基置換

ただし、変換直後の①fastaオブジェクトはdescription情報が
ついていない。(k-mer出現頻度解析自体には無関係だが、
気持ち悪いと思うヒトのために)②で塩基置換前の情報を保
持するfasta_orgのdescription情報部分を③でコピーしている

塩基置換

GならC、AならTのような塩基置換を行うenkichikan関数(孫建雄氏作成)を利用して任意の位置
の塩基置換を行う。「書籍 | トランスクリプトーム解析 | 2.3.4 マッピング(準備)」のp81と同じ関数で
す。

```
enkichikan(as.character(fasta[500]), 18) #18番目の塩基を置換

#塩基置換本番
fasta_org <- fasta #塩基置換
hoge <- enkichikan(as.character(fasta), 18) #18番 #塩基置換
fasta <- DNASTringSet(hoge) #DNA塩基置
names(fasta) <- names(fasta_org) #fasta_or #確認して
fasta #確認して

#k-mer出現頻度分布
out_f <- "hoge2_18.png" #出力ファイル名
param_kmer <- 8 #k-merの幅
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #列ごとの出現頻度
out <- colSums(hoge) #1回以上出現するk-mer
kmer <- out[out > 0] #下、左
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2], par=mar=c(4, 4, 0, 0)) #ヒストグラム
hist(kmer, breaks=max(kmer), ylab="Frequency", xlab=paste("Number of occurrences at k=", param_kmer), xlim=param_range, main="") #おまじな
dev.off()
```

```
R Console
> #塩基置換本番
> fasta_org <- fasta #塩基置$
> hoge <- enkichikan(as.character(fasta), 18) #18$
> head(hoge)
[1] "GGTACGGTTCGGTTGCGGA" "TTATCCGGCAGTCCTATTTG"
[3] "CAATAGACACCACGCGCGAC" "ACCGAAAGTTCTGGAGGACA"
[5] "TTGGGCCCGATATCTATTAA" "TATGATATGGGCCGCTGGGG"
> length(hoge)
[1] 500
> fasta <- DNASTringSet(hoge) #DNA塩$
> fasta
A DNASTringSet instance of length 500
width seq
[1] 20 GGTACGGTTCGGTTGCGGA
[2] 20 TTATCCGGCAGTCCTATTTG
[3] 20 CAATAGACACCACGCGCGAC
[4] 20 ACCGAAAGTTCTGGAGGACA
[5] 20 TTGGGCCCGATATCTATTAA
```

塩基置換後

①description情報をコピーした塩基置換後のfastaオブジェクトの中身は②

塩基置換
GならC、AならTのような塩基置換を行うenkichikan関数(孫 建強氏作成)を利用して任意の位置の塩基置換を行う。「書籍 | トランスクリプトーム解析 | 2.3.4 マッピング(準備)」のp81と同じ関数です。

```
enkichikan(as.character(fasta[500]), 18) #18番目の塩基を置換

#塩基置換本番
fasta_org <- fasta
hoge <- enkichikan(as.character(fasta), 18) #18番
fasta <- DNAStringSet(hoge)
names(fasta) <- names(fasta_org)
fasta

#k-mer出現頻度分布
out_f <- "hoge2_18.png"
param_kmer <- 8
hoge <- oligonucleotideFrequency(fasta, width=pa
out <- colSums(hoge)
kmer <- out[out > 0]
png(out_f, pointsize=13, width=param_fig[1], he
par(mar=c(4, 4, 0, 0))
hist(kmer, breaks=max(kmer), ylab="Frequency", #
      xlab=paste("Number of occurrences at k=", pa
      xlim=param_range, main="")
dev.off()
```

```
R Console
[499] 20 GACAGCGGGAGTGATCCCCG
[500] 20 TGTAAGCCGTACCCTACGCC
> names(fasta) <- names(fasta_org)
> fasta
A DNAStringSet instance of length 500
      width seq          names
[1] 20 GGTACGG...TGCGGA kkk_658_677
[2] 20 TTATCCG...TATTTG kkk_197_216
[3] 20 CAATAGA...CGCGAC kkk_559_578
[4] 20 ACCGAAA...AGGACA kkk_55_74
[5] 20 TTGGGCC...TATTAA kkk_681_700
... ..
[496] 20 GAGCTAT...CCCGAG kkk_433_452
[497] 20 AACCTAA...GTGCCG kkk_805_824
[498] 20 CGCCGCC...TCGGTT kkk_140_159
[499] 20 GACAGCG...TCCCGC kkk_401_420
[500] 20 TGTAAGC...TACGCC kkk_904_923
> |
```

塩基置換前

①塩基置換前のfasta_objオブジェクトの中身は②。意図通りに置換できていることがわかる

塩基置換

GならC、AならTのような塩基置換を行うenkichikan関数(孫 建強氏作成)を利用して任意の位置の塩基置換を行う。「書籍 | トランスクリプトーム解析 | 2.3.4 マッピング(準備)」のp81と同じ関数です。

```
enkichikan(as.character(fasta[500]), 18) #18番目の塩基を置換
```

#塩基置換本番

```
fasta_obj <- fasta #塩基置換  
hoge <- enkichikan(as.character(fasta), 18) #18番  
fasta <- DNAStringSet(hoge) #DNA塩基置換  
names(fasta) <- names(fasta_obj) #fasta_obj  
fasta #確認して
```

#k-mer出現頻度分布

```
out_f <- "hoge2_18.png" #出力ファイル名  
param_kmer <- 8 #k-merの幅  
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #k-merの出現頻度  
out <- colSums(hoge) #列ごとの出現頻度  
kmer <- out[out > 0] #1回以上出現するk-mer  
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2]) #下、左、右に出力  
par(mar=c(4, 4, 0, 0)) #ヒストグラムのマージン  
hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム  
      xlab=paste("Number of occurrences at k=", param_kmer), #おまじな  
      xlim=param_range, main="")  
dev.off()
```

```
R Console  
[498] 20 CGCCGCC...TCGGTT kkk_140_159  
[499] 20 GACAGCG...TCCCGC kkk_401_420  
[500] 20 TGTAAGC...TACGCC kkk_904_923  
> fasta_obj  
A DNAStringSet instance of length 500  
width seq names  
[1] 20 GGTACGG...TGCCGA kkk_658_677  
[2] 20 TTATCCG...TATATG kkk_197_216  
[3] 20 CAATAGA...CGCCAC kkk_559_578  
[4] 20 ACCGAAA...AGGTCA kkk_55_74  
[5] 20 TTGGGCC...TATAAA kkk_681_700  
... ..  
[496] 20 GAGCTAT...CCCCAG kkk_433_452  
[497] 20 AACCTAA...GTGGCG kkk_805_824  
[498] 20 CGCCGCC...TCGCTT kkk_140_159  
[499] 20 GACAGCG...TCCCGC kkk_401_420  
[500] 20 TGTAAGC...TACCCC kkk_904_923  
> |
```

塩基置換後の...

- 塩基置換
GならC、AならTのような塩基置換を行うenkichikan関数(孫 建強氏作成)を利用して任意の位置の塩基置換を行う。「書籍 | トランスクリプトーム解析 | 2.3.4 マッピング(準備)」のp81と同じ関数です。

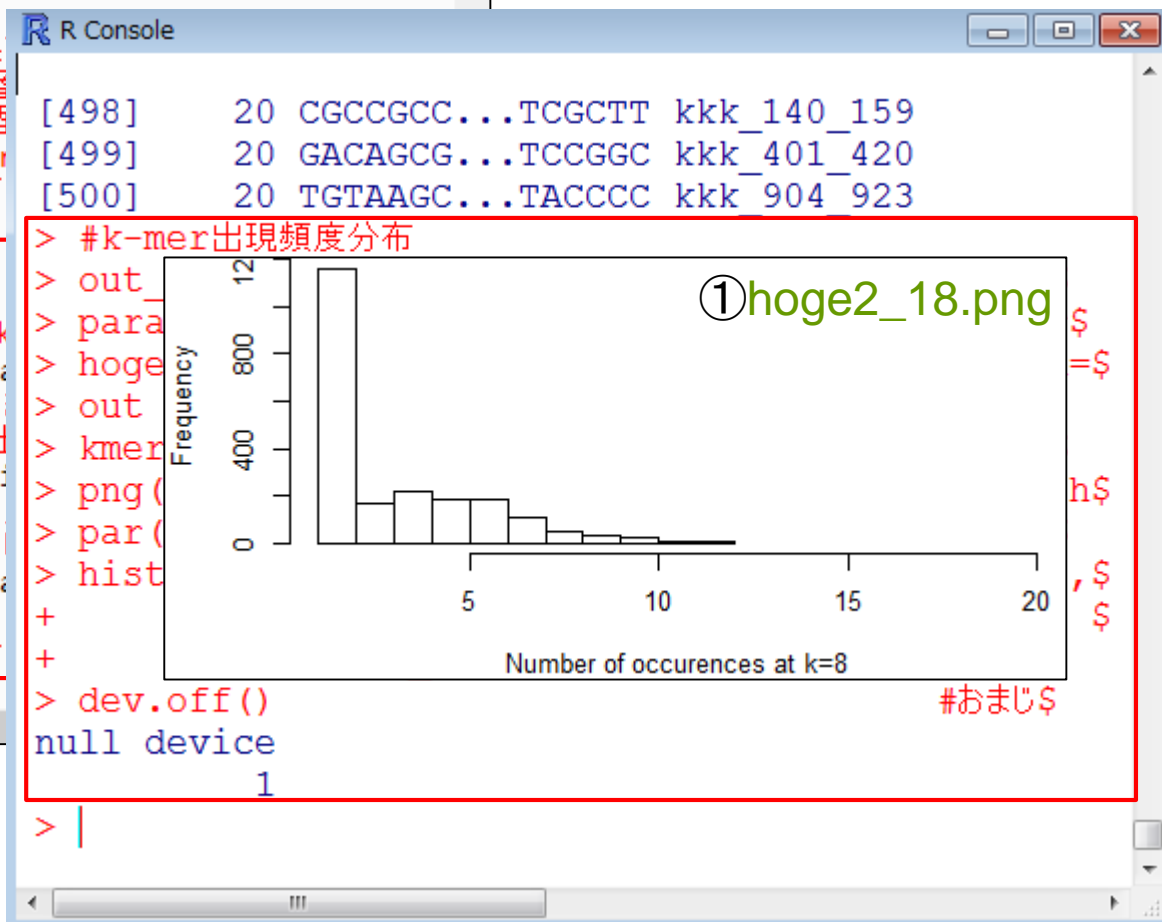
```
enkichikan(as.character(fasta[500]), 18) #18番目の塩基を置換
```

#塩基置換本番

```
fasta_org <- fasta #塩基置換  
hoge <- enkichikan(as.character(fasta), 18) #18番  
fasta <- DNAStringSet(hoge) #DNA塩基配列  
names(fasta) <- names(fasta_org) #fasta_or  
fasta #確認して
```

#k-mer出現頻度分布

```
out_f <- "hoge2_18.png" ① #出力ファイル名  
param_kmer <- 8 #k-merの幅  
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #k-merの出現頻度  
out <- colSums(hoge) #列ごとの出現頻度  
kmer <- out[out > 0] #1回以上出現するk-mer  
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2]) #下、左、右、上  
par(mar=c(4, 4, 0, 0)) #下、左、右、上  
hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム  
      xlab=paste("Number of occurrences at k=", param_kmer), #おまじな  
      xlim=param_range, main="")  
dev.off()
```



塩基置換後の

①塩基置換による仮想シーケンスエラー導入後と②導入前(同じk=8)の主な違いは、①のシーケンスエラーを含むデータにおいて、③1-2回しか出現しないk-merが1,200個弱存在する点。これがリアルデータでもよく見られるk-mer出現頻度分布

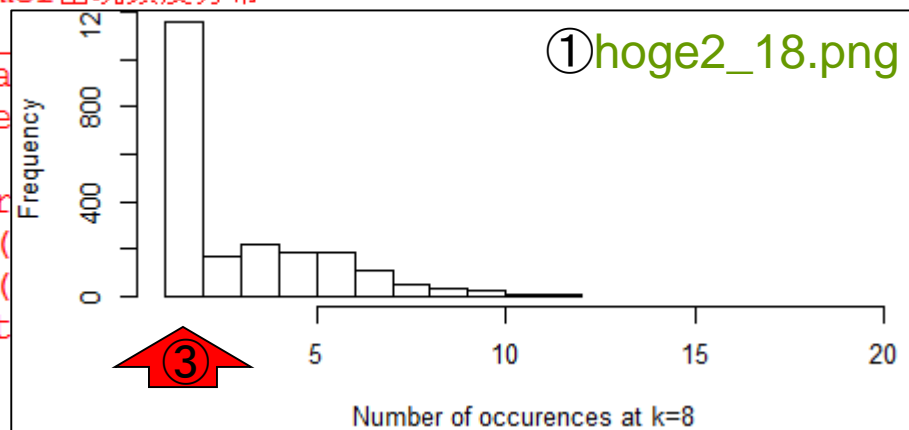
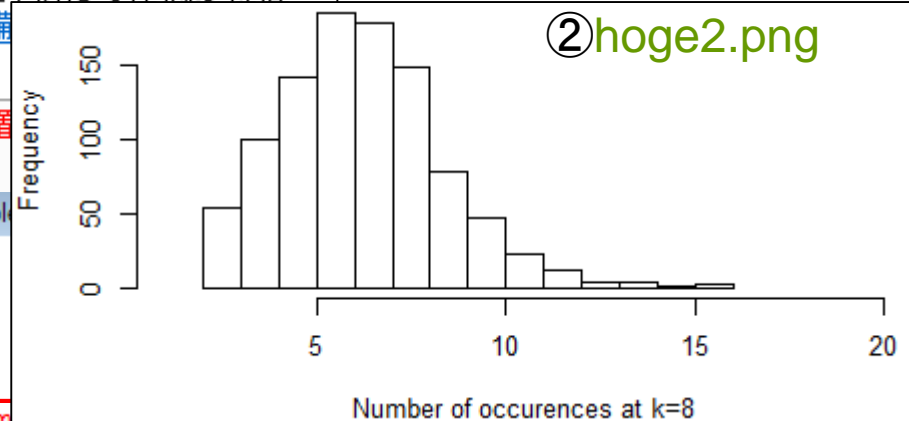
塩基置換

GならC、AならTのような塩基置換を行うenkichikan関数(関数)の塩基置換を行う。「書籍 | トランスクリプトーム解析 | 2.3.4 マッピング(準備)」

```
enkichikan(as.character(fasta[500]), 18) #18番目の塩基を置換  
  
#塩基置換本番  
fasta_org <- fasta  
hoge <- enkichikan(as.character(fasta), 18)#18番目の塩基を置換  
fasta <- DNASTringSet(hoge) #DNA塩基置換  
names(fasta) <- names(fasta_org) #fasta_orgの名称を保持  
fasta #確認して
```

```
#k-mer出現頻度分布  
out_f <- "hoge2_18.png" #出力ファイル名  
param_kmer <- 8 #k-merの幅  
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #oligonucleotideFrequency関数  
out <- colSums(hoge) #列ごとの出現頻度  
kmer <- out[out > 0] #1回以上出現するk-mer  
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2]) #png関数  
par(mar=c(4, 4, 0, 0)) #下、左、右、上の余白を調整  
hist(kmer, breaks=max(kmer), ylab="Frequency", #hist関数  
      xlab=paste("Number of occurrences at k=", param_kmer), #x軸ラベル  
      xlim=param_range, main="") #ヒストグラムの範囲とタイトル  
dev.off() #おまじな
```

①



```
> #k-mer出現頻度分布  
> out_f <- "hoge2_18.png"  
> param_kmer <- 8  
> hoge <- oligonucleotideFrequency(fasta, width=param_kmer)  
> out <- colSums(hoge)  
> kmer <- out[out > 0]  
> png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2])  
> par(mar=c(4, 4, 0, 0))  
> hist(kmer, breaks=max(kmer), ylab="Frequency",  
      xlab=paste("Number of occurrences at k=", param_kmer),  
      xlim=param_range, main="")  
+  
+  
> dev.off()  
null device  
1  
  
#おまじな
```

縦軸の範囲が違いすぎるので、①0-250に統一して眺める

縦軸の範囲を揃える

- 縦軸の範囲を揃える

ついでに、x軸の範囲指定のところの名前をparam_rangeからparam_xrangeに変更しています。

```
in_f <- "sample34_ngs.fasta"
param_fig <- c(450, 210)
param_xrange <- c(1, 20)
param_yrange <- c(0, 250)
param_kmer <- 8
```



#入力ファイル名を指
#ヒストグラム描画時
#表示したい横軸の範
#表示したい縦軸の範
#k-merのkの値を指定

```
library(Biostrings)
fasta <- readDNAStringSet(in_f, format="fasta")
fasta
```

#パッケージの読み込
#fastaで指定
#確認してるだけです

#塩基置換関数の作成

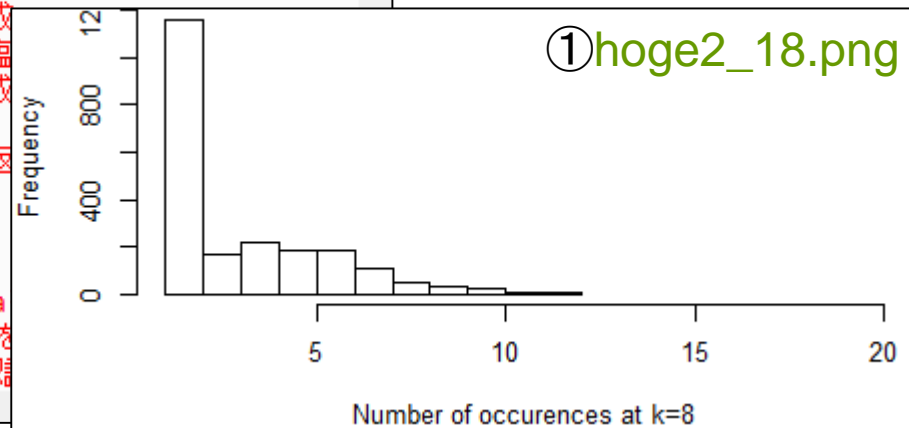
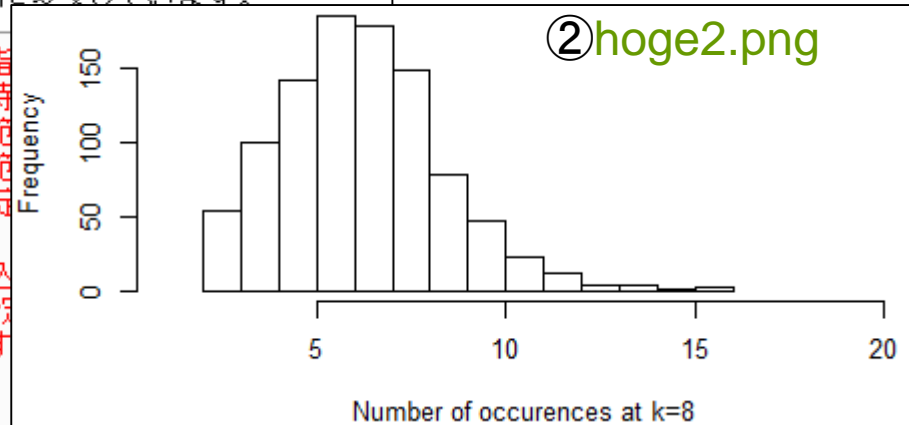
```
enkichikan <- function(fa, p) {
  t <- substring(fa, p, p)
  t_c <- chartr("CGAT", "GCTA", t)
  substring(fa, p, p) <- t_c
  return(fa)
}
```

#関数名や引数の作成
#置換したい位置の塩
#置換後の塩基を作成
#置換
#置換後のデータを返

#塩基置換本番

```
fasta_org <- fasta
hoge <- enkichikan(as.character(fasta), 18)
fasta <- DNAStringSet(hoge)
```

#塩基置換前のfasta
#18番目の塩基を
#DNA塩基配列だと認



縦軸の範囲を揃える

①コード下部はこんな感じ。②塩基置換後と③置換前のpngファイルを一気に作成して眺める。コード全体をコピー

縦軸の範囲を揃える

ついでに、x軸の範囲指定のところの名前をparam_rangeからparam_xrangeに変更しています。

```
#k-mer出現頻度分布(塩基置換後)
out_f <- "hoge2_18_250.png"
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)#k連続塩基の出現頻度
out <- colSums(hoge)
kmer <- out[out > 0]
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2])#出力ファイル
par(mar=c(4, 4, 0, 0))
hist(kmer, breaks=max(kmer), ylab="Frequency",#ヒストグラムを描画
      xlab=paste("Number of occurrences at k=", param_kmer, sep=""),#ヒストグラ
      xlim=param_xrange, ylim=param_yrange, main="")#ヒストグラムを描画
dev.off()
```

```
#k-mer出現頻度分布(塩基置換前)
fasta <- fasta_org
out_f <- "hoge2_250.png"
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)#k連続塩基の出現頻度
out <- colSums(hoge)
kmer <- out[out > 0]
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2])#出力ファイル
par(mar=c(4, 4, 0, 0))
hist(kmer, breaks=max(kmer), ylab="Frequency",#ヒストグラムを描画
      xlab=paste("Number of occurrences at k=", param_kmer, sep=""),#ヒストグラ
```

縦軸の範囲を揃える

①塩基置換後と②置換前。シーケンスエラーに相当する置換塩基を含むk-merは③のあたりに位置する。それ以外は置換塩基を含まないものたち

• 縦軸の範囲を揃える

ついでに、x軸の範囲指定のところの名前をparam_rangeからparam_xrangeに変更しています。

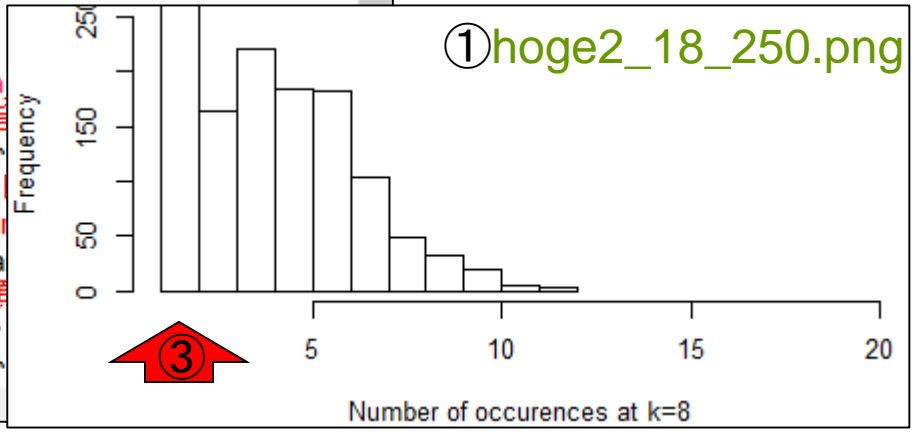
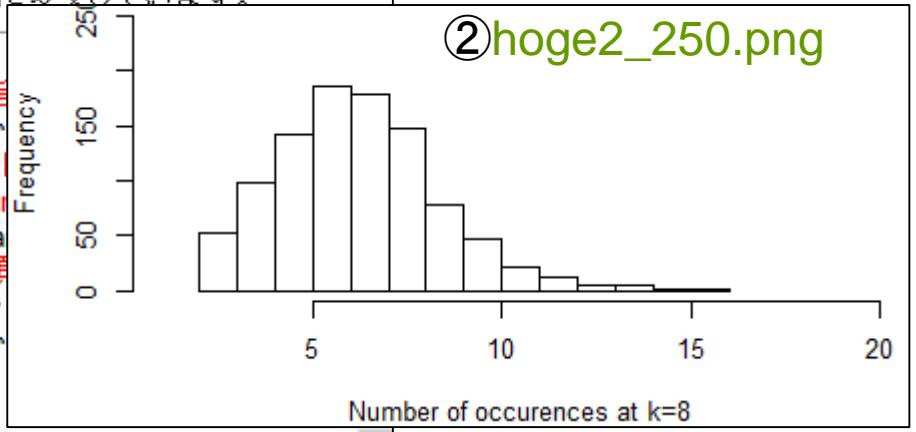
#k-mer出現頻度分布(塩基置換後)

```
out_f <- "hoge2_18_250.png"
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)
kmer <- out[out > 0]
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],
     par(mar=c(4, 4, 0, 0)))
hist(kmer, breaks=max(kmer), ylab="Frequency",
     xlab=paste("Number of occurrences at k=", param_kmer),
     xlim=param_xrange, ylim=param_yrange, main="")
dev.off()
```



#k-mer出現頻度分布(塩基置換前)

```
fasta <- fasta_org
out_f <- "hoge2_250.png"
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)
kmer <- out[out > 0]
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],
     par(mar=c(4, 4, 0, 0)))
hist(kmer, breaks=max(kmer), ylab="Frequency",
     xlab=paste("Number of occurrences at k=", param_kmer),
     xlim=param_xrange, ylim=param_yrange, main="")
dev.off()
```



おさらい

今は、長さ20 bpのリードが500個からなる仮想データに対して、18番目の塩基を置換してk-mer分布がどう変化するかを調べている

k=8の場合 : 13個のk-mer

GGTACGGTTCCGGTTGCCGA

GGTACGGT
GTACGGTT
TACGGTTC
ACGGTTCC
CGGTTCCG
GGTTCCGG
GTTCCGGT
TTCCGGTT
TCCGGTTG
CCGGTTGC
CGGTTGCC
GGTTGCCG
GTTGCCGA



k=8の場合 : 13個のk-mer

GGTACGGTTCCGGTTGC**G**GA

GGTACGGT
GTACGGTT
TACGGTTC
ACGGTTCC
CGGTTCCG
GGTTCCGG
GTTCCGGT
TTCCGGTT
TCCGGTTG
CCGGTTGC
CGGTTGC**G**
GGTTGC**GG**
GTTGC**G**GA



おさらい

k=8の場合は、リードあたり13個のk-merが生成できる。500リードのトータルで、 $500 \times 13 = 6,500$ 個のk-merとなる。そのうち、置換を行った18番目の塩基を含む、仮想シーケンスエラー由来k-merは、リードあたり3個。トータルで $500 \times 3 = 1,500$ 個

k=8の場合 : 13個のk-mer

GGTACGGTTCCGGTTGCCGA

GGTACGGT
GTACGGTT
TACGGTTC
ACGGTTCC
CGGTTCCG
GGTTCCGG
GTTCCGGT
TTCCGGTT
TCCGGTTG
CCGGTTGC
CGGTTGCC
GGTTGCCG
GTTGCCGA



k=8の場合 : 13個のk-mer

GGTACGGTTCCGGTTGC**G**GA

GGTACGGT
GTACGGTT
TACGGTTC
ACGGTTCC
CGGTTCCG
GGTTCCGG
GTTCCGGT
TTCCGGTT
TCCGGTTG
CCGGTTGC
CGGTTGC**G**
GGTTGC**GG**
GTTGC**GGA**



縦軸の範囲を揃え

①の総数が1,200弱程度だったのは、シークエンスエラー由来k-merが1,500個程度だったことに起因する。②塩基置換後の分布が③置換前に比べて左にずれているのは、エラーなしが6,500→5,000個に減っているため

縦軸の範囲を揃える

ついでに、x軸の範囲指定のところの名前をparam_rangeからparam_xrange

```
#k-mer出現頻度分布(塩基置換後)
out_f <- "hoge2_18_250.png"
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)
kmer <- out[out > 0]
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],
    par(mar=c(4, 4, 0, 0)))
hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム
     xlab=paste("Number of occurrences at k=", param_kmer),
     xlim=param_xrange, ylim=param_yrange, main="") #ヒストグラム
dev.off()
```

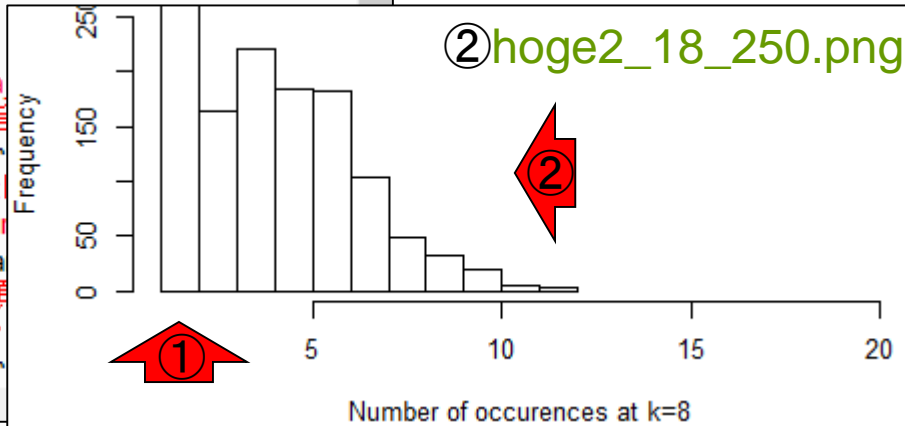
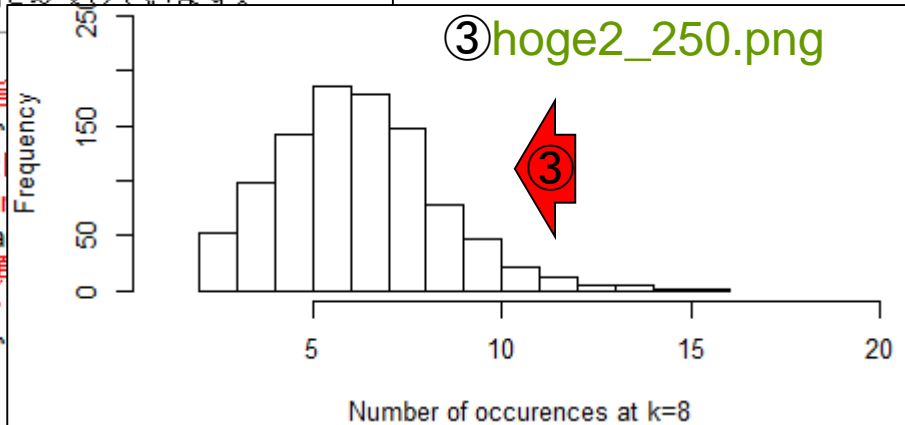
②

#出力ファイル名を指定
#列ごとの総和を出力
#1回以上出現したk-mer
#下、左、上、右の順番
#おまじない

```
#k-mer出現頻度分布(塩基置換前)
fasta <- fasta_org
out_f <- "hoge2_250.png"
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)
kmer <- out[out > 0]
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],
    par(mar=c(4, 4, 0, 0)))
hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム
     xlab=paste("Number of occurrences at k=", param_kmer),
     xlim=param_xrange, ylim=param_yrange, main="") #ヒストグラム
```

③

#塩基置換前のfasta
#出力ファイル名を指定
#列ごとの総和を出力
#1回以上出現したk-mer
#下、左、上、右の順番



Contents

- NGS解析手段、ウェブツール(DDBJ Pipeline)との連携
- DDBJ PipelineでPlatanusを実行
- k-mer解析(k個の連続塩基に基づく各種解析)の基礎
 - 塩基ごとの出現頻度解析(k=1)、一気に計算
 - 2連続塩基の出現頻度解析(k=2)、基本スキルの復習や作図
- de novoアセンブリ時のエラー補正やゲノムサイズ推定の基本的な考え方
 - ランダムな塩基配列(仮想ゲノムおよび仮想NGSデータ)の作成
 - k-mer解析の応用、ゲノムサイズ推定の基礎
 - ゲノムサイズ推定(1,000 bpの仮想ゲノムの場合)
 - ゲノムサイズ推定(1,000 bpの仮想ゲノム; 4X → 10X coverageの場合)
 - k-mer出現頻度分布
 - シークエンスエラーを含む場合
 - 最終確認



最終確認

k=8の場合は、(8番目から13番目の範囲であればどこでも同じであるが)もし②13番目の塩基に置換を入れれば、仮想シーケンスエラー由来k-merは、リードあたり8個。最後の50リード(451-500番目のリード)に対してこの置換を導入すると、トータルで $50 \times 8 = 400$ 個の仮想シーケンスエラー由来k-merを発生させることになる

k=8の場合: 13個のk-mer

GGTACGGTTCCGGTTGCCGA

```
GGTACGGT
GTACGGTT
TACGGTTC
ACGGTTCC
CGGTTCCG
GGTTCCGG
GTTCCGGT
TTCCGGTT
TCCGGTTG
CCGGTTGC
CGGTTGCC
GGTTGCCG
GTTGCCGA
```



k=8の場合: 13個のk-mer

GGTACGGTTCCGCTTGCCGA

```
GGTACGGT
GTACGGTT
TACGGTTC
ACGGTTCC
CGGTTCCG
GGTTCCGC
GTTCCGCT
TTCCGCTT
TCCGCTTG
CCGCTTGC
CGCTTGCC
GCTTGCCG
CTTGCCGA
```



狙い

エラーを全く含まない1-450番目のリードのみのk-mer分布と、エラーを多く含む451-500番目のリードのk-mer出現頻度分布を独立に作成して眺め、エラーを含まないk-merのみ取り扱える(フィルタリングできる)ことを学ぶ

k=8の場合: 13個のk-mer

GGTACGGTTCCGGTTGCCGA

```
GGTACGGT
GTACGGTT
TACGGTTC
ACGGTTCC
CGGTTCCG
GGTTCCGG
GTTCCGGT
TTCCGGTT
TCCGGTTG
CCGGTTGC
CGGTTGCC
GGTTGCCG
GTTGCCGA
```



k=8の場合: 13個のk-mer

GGTACGGTTCCGGCTTGCCGA

```
GGTACGGT
GTACGGTT
TACGGTTC
ACGGTTCC
CGGTTCCG
GGTTCCGGC
GTTCCGGCT
TTCCGGCTT
TCCGGCTTG
CCGGCTTGC
CGCTTGCC
GCTTGCCG
CTTGCCGA
```



最終確認

塩基置換を行わない最初の450リード分の情報は①fasta_first450としている

最終確認

500リードのうち、451:500番目のリードについてののみ13番目の塩基を置換。

```
in_f <- "sample34_ngs.fasta"
param_fig <- c(450, 210)
param_xrange <- c(1, 20)
param_yrange <- c(0, 350)
param_kmer <- 8
```

```
#入力ファイル名を指定してin_fに格納
#ヒストグラム描画時の横幅と縦幅を指定(単
#表示したい横軸の範囲を指定
#表示したい縦軸の範囲を指定
#k-merのkの値を指定
```

```
library(Biostrings)
fasta <- readDNAStringSet(in_f,
fasta
fasta_first450 <- fasta[1:450]
fasta_first450
```

#塩基置換関数の作成

```
enkichikan <- function(fa, p) {
  t <- substring(fa, p, p)
  t_c <- chartr("CGAT", "GCTA",
  substring(fa, p, p) <- t_c
  return(fa)
}
```

#塩基置換本番(451:500のみ)

```
fasta_org <- fasta
```

R Console

```
[499]    20 GACAGCGGGAGTGATCCGGC      kkk_401_420
[500]    20 TGTAAGCCGTACCCTACCCC      kkk_904_923
> fasta_first450 <- fasta[1:450]      #塩基置換を行わない$
> fasta_first450                       #確認してるだけです
A DNAStringSet instance of length 450
      width seq
[1]    20 GGTACGGTTCGGTTGCCGA      kkk_658_677
[2]    20 TTATCCGGCAGTCCTATATG      kkk_197_216
[3]    20 CAATAGACACCACGCGCCAC      kkk_559_578
[4]    20 ACCGAAAGTTCTGGAGGTCA      kkk_55_74
[5]    20 TTGGGCCCGATATCTATAAA      kkk_681_700
...    ...
[446]   20 ATATGATATGGGCCGCTGCG      kkk_212_231
[447]   20 ATTTTCGTTTCCGCATGACAC      kkk_118_137
[448]   20 GCTCGCTTTTACTTCCGGA      kkk_152_171
[449]   20 TTTTGAAC TCACTACACCAG      kkk_10_29
[450]   20 TGAAGACGCGGACGGCAGTC      kkk_34_53
> |
```

最終確認

①ちょっと下に移動。②赤枠部分で451-500番目のリードのみに対して13番目の塩基を置換した結果を③fasta_last50としている

- 最終確認
500リードのうち、451:500番目のリードについてののみ13番目の塩基を置換。

```
#塩基置換関数の作成
enkichikan <- function(fa, p) {
  t <- substring(fa, p, p)
  t_c <- chartr("CGAT", "GCTA", t)
  substring(fa, p, p) <- t_c
  return(fa)
}

#塩基置換本番(451:500のみ)
fasta_org <- fasta
hoge <- enkichikan(as.character(fasta[451:500]), 13)
fasta <- DNASTringSet(hoge)
names(fasta) <- names(fasta_org)
fasta_last50 <- fasta
fasta_last50

#k-mer出現頻度分布(original)
fasta <- fasta_org
out_f <- "hoge2_original.png"
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)
kmer <- out[out > 0]
```



最終確認

①13番目の塩基置換後のfasta_last50の中身。②置換された塩基ポジション

最終確認

500リードのうち、451:500番目のリードについてのみ13番目の塩基を置換。

#塩基置換関数の作成

```
enkichikan <- function(fa, p) {  
  t <- substring(fa, p, p)           #関数名や引数の作成  
  t_c <- chartr("CGAT", "GCTA", t)  #置換したい位置の塩基を取り出す  
  substring(fa, p, p) <- t_c        #置換後の塩基を作成  
  return(fa)                         #置換
```

#塩基置換本番(451:500のみ)

```
fasta_org <- fasta  
hoge <- enkichikan(as.character(fasta_org[451:500]), 13)  
fasta <- DNASTringSet(hoge)  
names(fasta) <- names(fasta_org[451:500])  
fasta_last50 <- fasta
```

#k-mer出現頻度分布(original)

```
fasta <- fasta_org  
out_f <- "hoge2_original.png"  
hoge <- oligonucleotideFrequency(fasta, k=13)  
out <- colSums(hoge)  
kmer <- out[out > 0]
```

```
R Console  
> fasta <- DNASTringSet(hoge)           #DNA塩基配列だと認識$  
> names(fasta) <- names(fasta_org[451:500]) #fasta_orgのdescri$  
> fasta_last50 <- fasta                #fasta_last50として$  
> fasta_last50                         #確認してるだけです  
A DNASTringSet instance of length 50  
  width seq                               names  
[1]    20 CTGCGACCGTTCCTCTAAATC          kkk_746_765  
[2]    20 TGGGCCTCCGGAACCCAGTT          kkk_956_975  
[3]    20 GGACGCTTTTCATGTCCAGC          kkk_933_952  
[4]    20 CAGTTTGTCTAACTGGTAGG          kkk_265_284  
[5]    20 GTACGGTTCGGATGCCGAT            kkk_659_678  
...  
[46]   20 GAGCTATCAATGAACCCAG            kkk_433_452  
[47]   20 AACCTAATCGCCGCGTGGCG            kkk_805_824  
[48]   20 CGCCGCCTTCGACCTCGCTT            kkk_140_159  
[49]   20 GACAGCGGGAGTCATCCGGC            kkk_401_420  
[50]   20 TGTAAGCCGTACGCTACCCC            kkk_904_923  
> |
```



最終確認

①置換前のfasta_orgオブジェクトの同じ部分を表示。②確かに13番目の塩基を意図通りに置換できていることがわかる

- 最終確認
500リードのうち、451:500番目のリードについてのみ13番目の塩基を置換。

```
#塩基置換関数の作成
enkichikan <- function(fa, p) {
  t <- substring(fa, p, p)
  t_c <- chartr("CGAT", "GCTA", t)
  substring(fa, p, p) <- t_c
  return(fa)
}
#関数名や引数の作成
#置換したい位置の塩基を取り出す
#置換後の塩基を作成
#置換
```

```
#塩基置換本番(451:500のみ)
fasta_org <- fasta
hoge <- enkichikan(as.character(fasta))
fasta <- DNASTringSet(hoge)
names(fasta) <- names(fasta_org)
fasta_last50 <- fasta
fasta_last50
```

```
#k-mer出現頻度分布(original)
fasta <- fasta_org
out_f <- "hoge2_original.png"
hoge <- oligonucleotideFrequency(fasta)
out <- colSums(hoge)
kmer <- out[out > 0]
```

```
R Console
[48] 20 CGCCGCCTTCGACCTCGCTT kkk_140_159
[49] 20 GACAGCGGGAGTCATCCGGC kkk_401_420
[50] 20 TGTAAGCCGTACGCTACCCC kkk_904_923
> fasta_org[451:500]
A DNASTringSet instance of length 50
width seq names
[1] 20 CTGCGACCGTTCACTAAATC kkk_746_765
[2] 20 TGGGCCTCCGGATCCCAGTT kkk_956_975
[3] 20 GGACGCTTTTCAAGTCCAGC kkk_933_952
[4] 20 CAGTTTGTCTAAGTGGTAGG kkk_265_284
[5] 20 GTACGGTTCGGTTGCCGAT kkk_659_678
... ..
[46] 20 GAGCTATCAATGTACCCAG kkk_433_452
[47] 20 AACCTAATCGCCCGTGGCG kkk_805_824
[48] 20 CGCCGCCTTCGAGCTCGCTT kkk_140_159
[49] 20 GACAGCGGGAGTGATCCGGC kkk_401_420
[50] 20 TGTAAGCCGTACCCCTACCCC kkk_904_923
> |
```

塩基置換前の全リード

まずはおさらいとして、①オリジナルの全500リードからなる塩基置換前のfasta_orgのk-mer出現頻度分布や周辺情報を見ておく

- 最終確認
500リードのうち、451:500番目のリードについてのみ13番目の塩基を置換。

```
#k-mer出現頻度分布(original) ①
fasta <- fasta_org                #fasta_orgをfastaに格納
out_f <- "hoge2_original.png"    #出力ファイル名を指定してout_fに格納
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #k連続塩基の出現頻度
out <- colSums(hoge)              #列ごとの総和をoutに格納
kmer <- out[out > 0]             #1回以上出現したk-merのみをkmerに格納
length(kmer)                     #1回以上出現したk-merの種類数を表示
table(kmer)                      #(k-merの種類は問わずに)k-merの出現回数
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2]) #出力ファイル
par(mar=c(4, 4, 0, 0))          #下、左、上、右の順で余白(行)を指定
hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラムを描画
      xlab=paste("Number of occurrences at k=", param_kmer, sep=""), #ヒストグラム
      xlim=param_xrange, ylim=param_yrange, main="") #ヒストグラムを描画
dev.off()                        #おまじない
kmer_org <- kmer                 #kmerをkmer_orgに格納

#k-mer出現頻度分布(最初の450リードのみ; 塩基置換なし)
fasta <- fasta_first450          #fasta_first450をfastaに格納
out_f <- "hoge2_first450.png"   #出力ファイル名を指定してout_fに格納
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #k連続塩基の出現頻度
out <- colSums(hoge)            #列ごとの総和をoutに格納
kmer <- out[out > 0]            #1回以上出現したk-merのみをkmerに格納
```

塩基置換前の全リード

①の結果は、②k-merの種類数は975個であり、これが推定ゲノムサイズ(正解は1,000 bp)。③ヒストグラムはbin次第で見栄えが異なりうる。④table(kmer)で元データも眺める

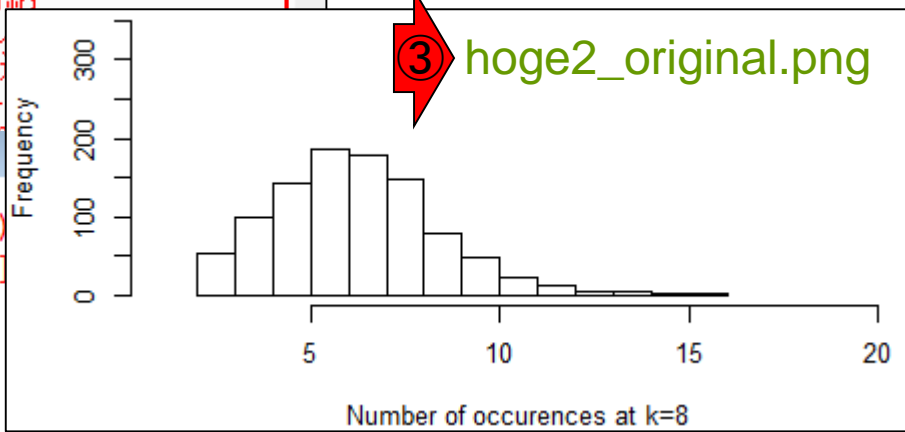
- 最終確認
500リードのうち、451:500番目のリードについてのみ13番目の塩基を置換。

#k-mer出現頻度分布(original) ①

```
fasta <- fasta_org  
out_f <- "hoge2_original.png"  
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #k-mer出現頻度分布  
out <- colSums(hoge) #出力ファイル名を指定し、列ごとの総和をoutに格納  
kmer <- out[out > 0] #k-merの種類を抽出  
length(kmer) #k-merの種類数を取得  
table(kmer) #k-merの種類数を取得  
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],  
par(mar=c(4, 4, 0, 0)) #下、左、上、右の順  
hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム  
xlab=paste("Number of occurrences at k=", param_kmer, s$  
xlim=param_xrange, ylim=param_yrange, main="") #ヒストグラム  
dev.off() #おまじない  
kmer_org <- kmer #kmerをkmer_orgに格納
```

#k-mer出現頻度分布(最初の450リード)
fasta <- fasta_first450
out_f <- "hoge2_first450.png"
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)
kmer <- out[out > 0]

```
R Console  
> out <- colSums(hoge)  
> kmer <- out[out > 0]  
> length(kmer)  
[1] 975  
> table(kmer)  
kmer  
 2  3  4  5  6  7  8  9 10 11 12 13 14 15 16  
18 35 99 142 185 178 148 78 47 22 12 4 4 1 2  
> png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],  
> par(mar=c(4, 4, 0, 0)) #下、左、上、右の順  
> hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム  
+ xlab=paste("Number of occurrences at k=", param_kmer, s$  
+ xlim=param_xrange, ylim=param_yrange, main="") #ヒストグラム  
> dev.off() #おまじない  
null device  
 1  
> kmer_org <- kmer #kmerをkmer_orgに格納  
> |
```



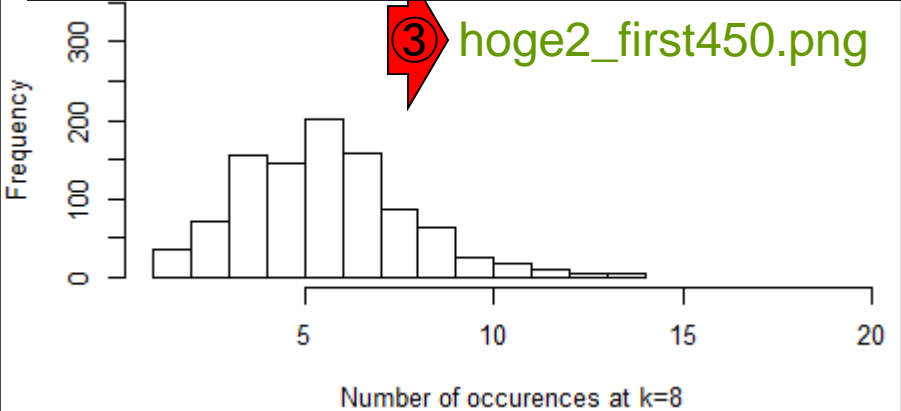
最初の450リードのみ

- 最終確認
500リードのうち、451:500番目のリードについてのみ13番目の塩基を置換。

```
#k-mer出現頻度分布(最初の450リードのみ; 塩基置換なし) ①  
fasta <- fasta_first450 #fasta_first50をfasta  
out_f <- "hoge2_first450.png" #出力ファイル名を指定し  
hoge <- oligonucleotideFrequency(fasta, width=param_kmer) #k  
out <- colSums(hoge) #列ごとの総和をoutに格  
kmer <- out[out > 0]  
length(kmer)  
table(kmer) ②  
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],  
par(mar=c(4, 4, 0, 0)) #下、左、上、右の順$  
hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム$  
xlab=paste("Number of occurrences at k=", param_kmer, s$  
xlim=param_xrange, ylim=param_yrange, main="") #ヒスト$  
dev.off() #おまじない  
kmer_first450 <- kmer
```

```
#k-mer出現頻度分布(最後の50リードのみ)  
fasta <- fasta_last50  
out_f <- "hoge2_last50.png"  
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)  
out <- colSums(hoge)  
kmer <- out[out > 0]
```

```
R Console  
> out <- colSums(hoge)  
> kmer <- out[out > 0]  
> length(kmer)  
[1] 975  
> table(kmer)  
kmer  
 1  2  3  4  5  6  7  8  9 10 11 12 13 14  
 5 31 71 155 145 201 158 87 64 24 17 10 4 3  
> png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],  
> par(mar=c(4, 4, 0, 0)) #下、左、上、右の順$  
> hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム$  
+ xlab=paste("Number of occurrences at k=", param_kmer, s$  
+ xlim=param_xrange, ylim=param_yrange, main="") #ヒスト$  
> dev.off() #おまじない  
null device  
 1  
> kmer_first450 <- kmer #kmerをkmer_first450$  
> |
```



①の結果は、②k-merの種類数は975個であり、これが推定ゲノムサイズ(正解は1,000 bp)。③ヒストグラムはbin次第で見栄えが異なりうる。④table(kmer)で元データも眺める。シーケンスエラーはないが、リード数が減っているため、分布も左に少しずれている

最後の50リードのみ

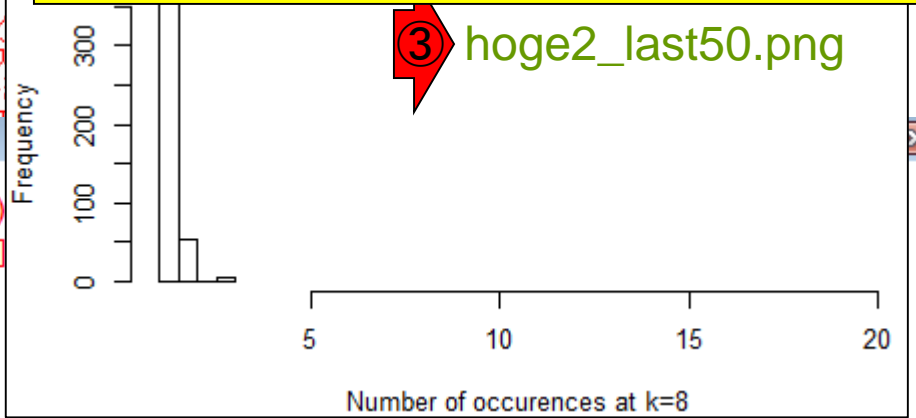
①の結果は、②k-merの種類数は590個であり、これが推定ゲノムサイズ(正解は1,000 bp)。③ヒストグラムはbin次第で見栄えが異なりうる。④table(kmer)で元データも眺める。シークエンスエラー由来k-merがほとんどなので、出現回数が最大でも3回になっている

- 最終確認
500リードのうち、451:500番目のリードについてのみ13番目の塩基を置換。

```
#k-mer出現頻度分布(最後の50リードのみ;塩基置換あり) ①
fasta <- fasta_last50 #fasta_last50をfastaに代入
out_f <- "hoge2_last50.png" #出力ファイル名を指定し
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)#k
out <- colSums(hoge) #列ごとの総和をoutに格
kmer <- out[out > 0] #0以上出現した
length(kmer)
table(kmer)
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],
par(mar=c(4, 4, 0, 0))
hist(kmer, breaks=max(kmer), #ヒストグラム$
xlab=paste("Number of occurrences at k=", param_kmer, s$
xlim=param_xrange, ylim=param_yrange, main="") #ヒスト$
dev.off()
kmer_last50 <- kmer
```

```
#k-mer出現頻度分布(merged)
fasta <- c(fasta_first450, fasta_last50)
out_f <- "hoge2_merged.png"
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)
out <- colSums(hoge)
kmer <- out[out > 0]
```

```
R Console
> out <- colSums(hoge)
> kmer <- out[out > 0]
> length(kmer)
[1] 590
> table(kmer)
kmer
 1  2  3
533 54  3
> png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],
par(mar=c(4, 4, 0, 0))
> hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム$
+ xlab=paste("Number of occurrences at k=", param_kmer, s$
+ xlim=param_xrange, ylim=param_yrange, main="") #ヒスト$
> dev.off()
null device
      1
> kmer_last50 <- kmer #kmerをkmer_last50に代入
> |
```



マージ後のデータ

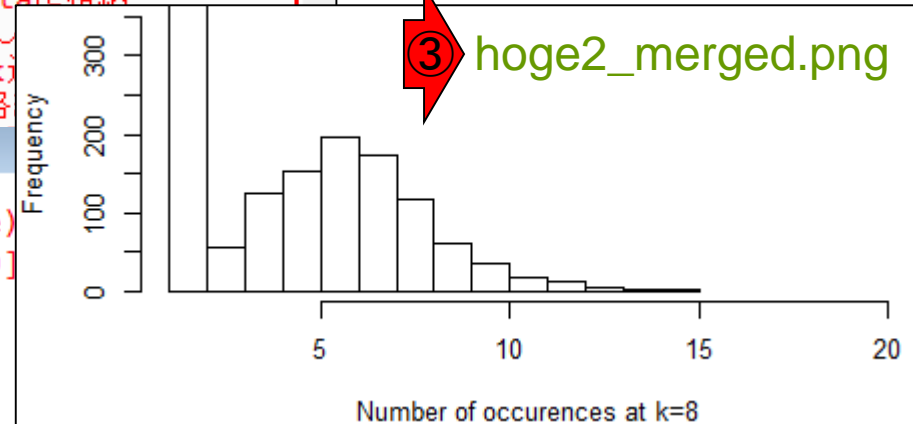
①の結果は、②k-merの種類数は1,341個であり、これが推定ゲノムサイズ(正解は1,000 bp)。③ヒストグラムはbin次第で見栄えが異なりうる。④table(kmer)で元データも眺める

- 最終確認
500リードのうち、451:500番目のリードについてのみ13番目の塩基を置換。

```
#k-mer出現頻度分布(merged)
fasta <- c(fasta_first450, fasta_last50)#マージしたものをfastaに格納
out_f <- "hoge2_merged.png" #出力ファイル名を指定し
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)#k
out <- colSums(hoge) #列ごとの総和をoutに格納
kmer <- out[out > 0]
length(kmer)
table(kmer)
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],
par(mar=c(4, 4, 0, 0))
hist(kmer, breaks=max(kmer),
xlab=paste("Number of occurrences at k=", param_kmer, "bp"),
xlim=param_xrange, ylim=param_yrange, main="")
dev.off()
kmer_merged <- kmer
```

R Console

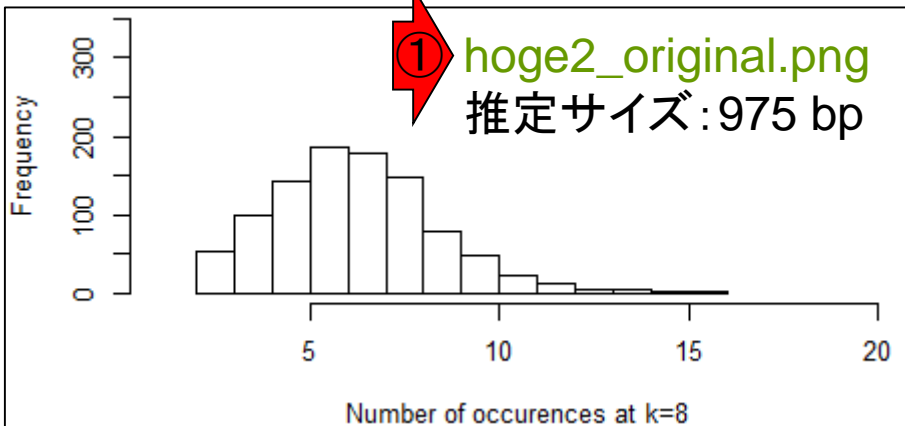
```
> out <- colSums(hoge)
> kmer <- out[out > 0]
> length(kmer)
[1] 1341
> table(kmer)
kmer
 1  2  3  4  5  6  7  8  9 10 11 12 13 14 15
342 51 56 125 153 195 172 116 60 35 16 11 5 2 2
> png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],
> par(mar=c(4, 4, 0, 0)) #下、左、上、右の順$
> hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム$
+ xlab=paste("Number of occurrences at k=", param_kmer, "bp"), s$
+ xlim=param_xrange, ylim=param_yrange, main="") #ヒスト$
> dev.off() #おまじない
null device
      1
> kmer_merged <- kmer #kmerをkmer_mergedに$
> |
```



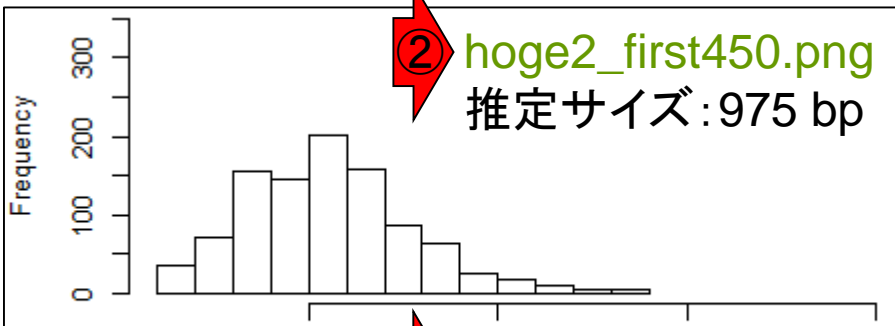
比較

①オリジナルの塩基置換前の全500リードの結果。②最初の450リードの結果、③13番目の塩基置換後の最後の50リードの結果。④は、②と③を合わせた結果。③の結果はシーケンスエラー由来k-merが多くを占めていたことを思い出そう

① hoge2_original.png
推定サイズ: 975 bp



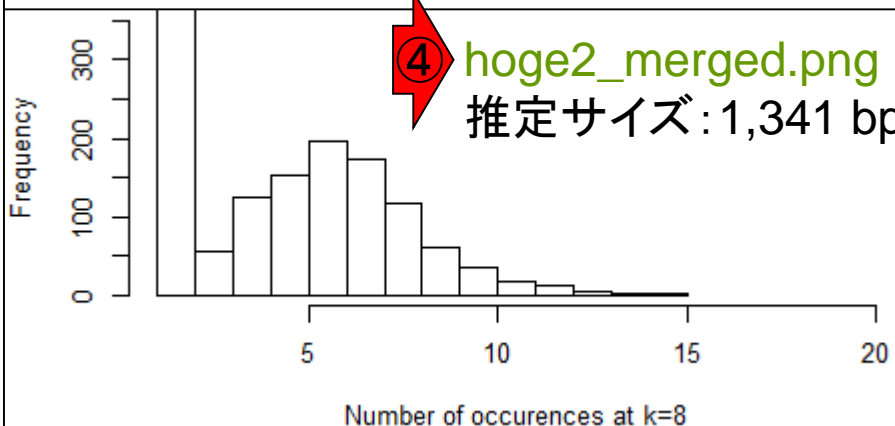
② hoge2_first450.png
推定サイズ: 975 bp



③ hoge2_last50.png
推定サイズ: 590 bp



④ hoge2_merged.png
推定サイズ: 1,341 bp



マージ後のデータ

①シークエンスエラー由来k-merが多くを占めていた部分まで含めて算出した②推定ゲノムサイズは1,341 bp。正解(1,000 bp)からのずれは、③シークエンスエラー由来k-merの結果部分に起因するので除いちゃえばいい

最終確認

500リードのうち、451:500番目のリードについてのみ13番目の塩基を置換。

#k-mer出現頻度分布(merged)

```
fasta <- c(fasta_first450, fasta_last50)#マージしたものをfastaに格納
out_f <- "hoge2_merged.png" #出力ファイル名を指定し
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)#k
out <- colSums(hoge) #列ごとの総和をoutに格納
kmer <- out[out > 0]
length(kmer)
table(kmer)
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],
par(mar=c(4, 4, 0, 0))
hist(kmer, breaks=max(kmer),
xlab=paste("Number of occurrences at k=", param_kmer, s$),
xlim=param_xrange, ylim=param_yrange, main="")#ヒストグラム$
dev.off()
kmer_merged <- kmer
```



R Console

```
> out <- colSums(hoge)
> kmer <- out[out > 0]
> length(kmer)
[1] 1341
> table(kmer)
kmer
 1  2  3  4  5  6  7  8  9 10 11 12 13 14 15
342 51 56 125 153 195 172 116 60 35 16 11 5 2 2
> png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2],
> par(mar=c(4, 4, 0, 0)) #下、左、上、右の順$
> hist(kmer, breaks=max(kmer), ylab="Frequency", #ヒストグラム$
+ xlab=paste("Number of occurrences at k=", param_kmer, s$),
+ xlim=param_xrange, ylim=param_yrange, main="")#ヒスト$
> dev.off() #おまじない
null device
      1
> kmer_merged <- kmer #kmerをkmer_mergedに$
> |
```

k	Frequency
1	342
2	51
3	56
4	125
5	153
6	195
7	172
8	116
9	60
10	35
11	16
12	11
13	5
14	2
15	2

マージ後のデータ

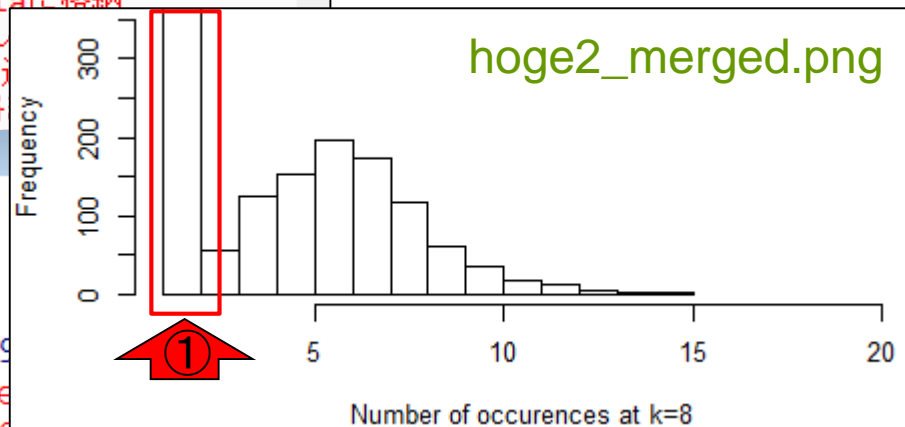
①シークエンスエラーを多く含むk-mer部分を除いて考えるのは、この場合は②2回以上、あるいは③3回以上出現するk-merの種類数のみをカウントすることに相当

最終確認

500リードのうち、451:500番目のリードについてのみ13番目の塩基を置換。

```
#k-mer出現頻度分布(merged)
fasta <- c(fasta_first450, fasta_last50)#マージしたものをfastaに格納
out_f <- "hoge2_merged.png" #出力ファイル名を指定し
hoge <- oligonucleotideFrequency(fasta, width=param_kmer)#k
out <- colSums(hoge) #列ごとの総和をoutに格納
kmer <- out[out > 0]
length(kmer)
table(kmer)
png(out_f, pointsize=13, width=
par(mar=c(4, 4, 0, 0))
hist(kmer, breaks=max(kmer), ylab=
xlab=paste("Number of occurrences at k=", param_kmer, s$
xlim=param_xrange, ylim=param_yrange, main="")#ヒストグラム$
dev.off()
kmer_merged <- kmer
```

```
R Console
> table(kmer)
kmer
 1  2  3  4  5
342 51 56 125 153 19
> sum(kmer >= 2)
[1] 999
> sum(kmer >= 3)
[1] 948
> |
```



#下、左、上、右の順番
#ヒストグラム\$
#ヒスト\$
#おまじない
#kmerをkmer_mergedに\$