

UNIVERSITÉ GRENOBLE-ALPES

THÈSE

Pour obtenir le grade de

DOCTEUR DE L'UNIVERSITÉ DE GRENOBLE-ALPES

Spécialité : Modèles, méthodes et algorithmes en biologie, santé et environnement

Arrêté ministériel : ?

Présentée par

Keurcien LUU

Thèse dirigée par **Michael BLUM**

préparée au sein du laboratoire **Techniques de l'Ingénierie Médicale et de la Complexité - Informatique, Mathématiques et Applications de Grenoble (TIMC-IMAG)**

et de l'école doctorale "**Ingénierie de la Santé, de la Cognition et Environnement**" (EDISCE)

Méthodes statistiques en grande dimension pour l'étude de l'adaptation biologique à l'aide de larges bases de données génomiques

Thèse soutenue publiquement le 31 octobre 2017,
devant le jury composé de :

Remerciements

Je tiens à remercier mes collègues Kevin Caye, Thomas Dias-Alves, Thomas Karaouzène et Florian Privé, avec qui j'ai partagé ces trois années de thèse et de qui j'ai beaucoup appris.

Préface

This is an example of a thesis setup to use the reed thesis document class (for LaTeX) and the R bookdown package, in general.

Table des matières

Introduction	1
Chapitre 1 : État de l'art	3
1.1 Analyse en Composantes Principales parcimonieuse	3
1.2 Bootstrap ACP	3
1.3 Contexte	3
1.4 Tests multiples	4
1.5 Contrôle du taux de fausse découverte	4
1.6 R chunks	4
1.7 Inline code	4
1.8 Including plots	5
1.9 Loading and exploring data	5
1.10 Additional resources	9
Chapitre 2 : Adaptation locale	11
2.1 Math	11
2.2 Chemistry 101 : Symbols	11
2.2.1 Typesetting reactions	12
2.2.2 Other examples of reactions	12
2.3 Physics	12
2.4 Biology	12
Chapitre 3 : Introgression adaptative	13
3.1 Coefficients de métissage globaux et locaux	13
3.2 Introgression	13
3.3 Lien entre Analyse en Composantes Principales et métissage global. .	13
3.4 Analyse en Composantes Principales locale	14
3.5 Sensibilité à l'imputation des données manquantes	14
3.6 Simulations	14
3.6.1 Données de peupliers	14
3.6.2 Résultats de la comparaison des logiciels	18
3.7 Figures	22
3.8 Footnotes and Endnotes	24
3.9 Bibliographies	24
3.10 Anything else?	26

Conclusion	27
Annexe A : The First Appendix	29
Annexe B : The Second Appendix, for Fun	31
References	33

Liste des tableaux

1.1 Max Delays by Airline	8
3.1 Correlation of Inheritance Factors for Parents and Child	20

Table des figures

1.1	β	9
3.1	$\lambda = 0.1$	17
3.2	$\lambda = 10$	17
3.3	$\lambda = 50$	17
3.4	Reed logo	22
3.5	Mean Delays by Airline	23
3.6	Subdiv. graph	24
3.7	A Larger Figure, Flipped Upside Down	24

Abstract

The preface pretty much says it all.
Second paragraph of abstract starts here.

Introduction

Chapitre 1

État de l'art

- ACP en génétique des populations
- Partie III : R package pcadapt

1.1 Analyse en Composantes Principales parcimonieuse

1.2 Bootstrap ACP

1.3 Contexte

It's easy to create a list. It can be unordered like

- Item 1
- Item 2

or it can be ordered like

1. Item 1
2. Item 2

Notice that I intentionally mislabeled Item 2 as number 4. *Markdown* automatically figures this out ! You can put any numbers in the list and it will create the list. Check it out below.

To create a sublist, just indent the values a bit (at least four spaces or a tab). (Here's one case where indentation is key !)

1. Item 1
2. Item 2
3. Item 3
 - Item 3a
 - Item 3b

1.4 Tests multiples

1.5 Contrôle du taux de fausse découverte

Le taux de fausse découverte, correspond à la proportion de faux positifs parmi les positifs. En notant FP le nombre de faux positifs, TP le nombre de vrais positifs, on définit le taux de fausse découverte FDR par :

$$FDR = \mathbb{E} \left[\frac{FP}{TP + FP} 1_{FP+TP>0} \right]$$

- Référence cours de Christophe Giraud

q-value, bonferroni, benjamini-hochberg La figure suivante donne les comparaisons entre les différentes procédures de correction :

Now for the correct way :

Here is the first sentence. Here is another sentence. Here is the last sentence to end the paragraph.

This should be a new paragraph.

1.6 R chunks

When you click the **Knit** button above a document will be generated that includes both content as well as the output of any embedded **R** code chunks within the document. You can embed an **R** code chunk like this (**cars** is a built-in **R** dataset) :

```
summary(cars)
```

```
##      speed          dist
##  Min.   : 4.0   Min.   : 2.00
##  1st Qu.:12.0   1st Qu.: 26.00
##  Median :15.0   Median : 36.00
##  Mean   :15.4   Mean   : 42.98
##  3rd Qu.:19.0   3rd Qu.: 56.00
##  Max.   :25.0   Max.   :120.00
```

1.7 Inline code

If you'd like to put the results of your analysis directly into your discussion, add inline code like this :

The `cos` of 2π is 1.

Another example would be the direct calculation of the standard deviation :

The standard deviation of `speed` in `cars` is 5.2876444.

One last neat feature is the use of the `ifelse` conditional statement which can be used to output text depending on the result of an **R** calculation :

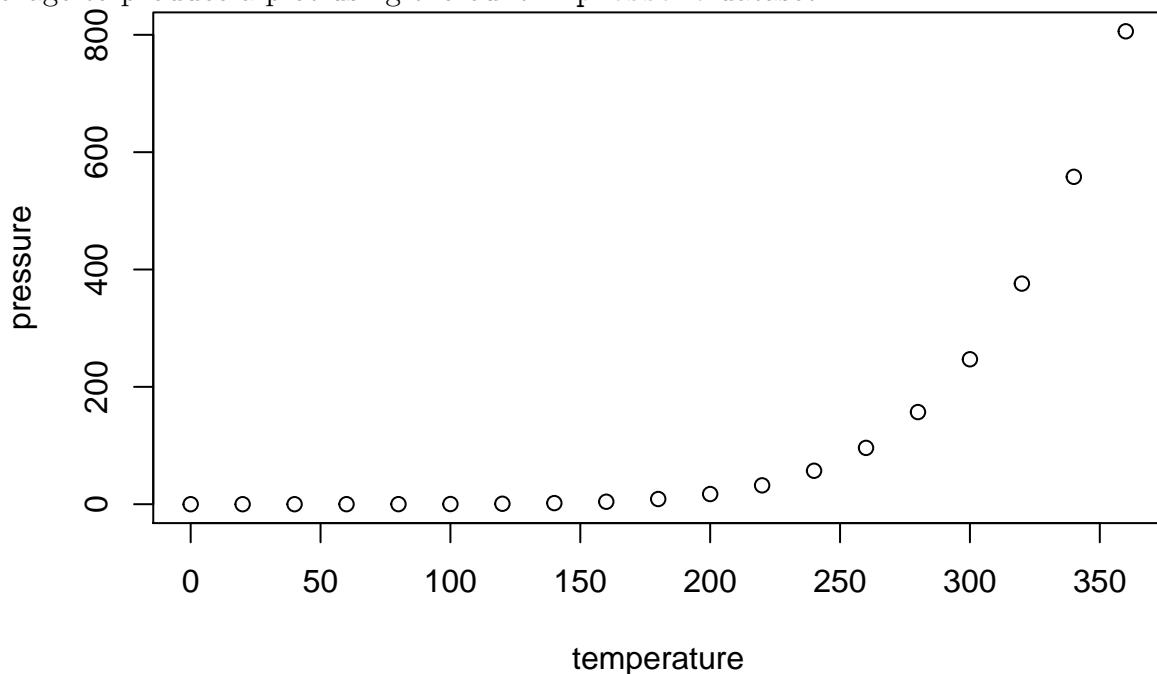
The standard deviation is less than 6.

Note the use of `>` here, which signifies a quotation environment that will be indented.

As you see with `2π` above, mathematics can be added by surrounding the mathematical text with dollar signs. More examples of this are in [Mathematics and Science] if you uncomment the code in Math.

1.8 Including plots

You can also embed plots. For example, here is a way to use the base **R** graphics package to produce a plot using the built-in `pressure` dataset :



Note that the `echo=FALSE` parameter was added to the code chunk to prevent printing of the **R** code that generated the plot. There are plenty of other ways to add chunk options. More information is available at <http://yihui.name/knitr/options/>.

Another useful chunk option is the setting of `cache=TRUE` as you see here. If document rendering becomes time consuming due to long computations or plots that are expensive to generate you can use knitr caching to improve performance. Later in this file, you'll see a way to reference plots created in **R** or external figures.

1.9 Loading and exploring data

Included in this template is a file called `flights.csv`. This file includes a subset of the larger dataset of information about all flights that departed from Seattle and

Portland in 2014. More information about this dataset and its **R** package is available at <http://github.com/ismayc/pnwflights14>. This subset includes only Portland flights and only rows that were complete with no missing values. Merges were also done with the `airports` and `airlines` data sets in the `pnwflights14` package to get more descriptive airport and airline names.

We can load in this data set using the following command :

```
flights <- read.csv("data/flights.csv")
```

The data is now stored in the data frame called `flights` in **R**. To get a better feel for the variables included in this dataset we can use a variety of functions. Here we can see the dimensions (rows by columns) and also the names of the columns.

```
dim(flights)
```

```
## [1] 52808     16
```

```
names(flights)
```

```
##  [1] "month"        "day"          "dep_time"      "dep_delay"
##  [5] "arr_time"      "arr_delay"     "carrier"       "tailnum"
##  [9] "flight"        "dest"         "air_time"      "distance"
## [13] "hour"          "minute"       "carrier_name"  "dest_name"
```

Another good idea is to take a look at the dataset in table form. With this dataset having more than 50,000 rows, we won't explicitly show the results of the command here. I recommend you enter the command into the Console *after* you have run the **R** chunks above to load the data into **R**.

```
View(flights)
```

While not required, it is highly recommended you use the `dplyr` package to manipulate and summarize your data set as needed. It uses a syntax that is easy to understand using chaining operations. Below I've created a few examples of using `dplyr` to get information about the Portland flights in 2014. You will also see the use of the `ggplot2` package, which produces beautiful, high-quality academic visuals.

We begin by checking to ensure that needed packages are installed and then we load them into our current working environment :

```
# List of packages required for this analysis
pkg <- c("dplyr", "ggplot2", "knitr", "bookdown", "devtools", "simulate", "data.table")
# Check if packages are not installed and assign the
# names of the packages not installed to the variable new.pkg
new.pkg <- pkg[!(pkg %in% installed.packages())]
```

```
# If there are any packages in the list that aren't installed,  
# install them  
if (length(new.pkg))  
  install.packages(new.pkg, repos = "http://cran.rstudio.com")  
# Load packages (thesisdown will load all of the packages as well)  
library(thesisdown)
```

The example we show here does the following :

- Selects only the `carrier_name` and `arr_delay` from the `flights` dataset and then assigns this subset to a new variable called `flights2`.
- Using `flights2`, we determine the largest arrival delay for each of the carriers.

```
flights2 <- flights %>%
  select(carrier_name, arr_delay)
max_delays <- flights2 %>%
  group_by(carrier_name) %>%
  summarize(max_arr_delay = max(arr_delay, na.rm = TRUE))
```

A useful function in the `knitr` package for making nice tables in *R Markdown* is called `kable`. It is much easier to use than manually entering values into a table by copying and pasting values into Excel or LaTeX. This again goes to show how nice reproducible documents can be! (Note the use of `results="asis"`, which will produce the table instead of the code to create the table.) The `caption.short` argument is used to include a shorter title to appear in the List of Tables.

```
kable(max_delays,
  col.names = c("Airline", "Max Arrival Delay"),
  caption = "Maximum Delays by Airline",
  caption.short = "Max Delays by Airline",
  longtable = TRUE,
  booktabs = TRUE)
```

TABLE 1.1 – Maximum Delays by Airline

Airline	Max Arrival Delay
Alaska Airlines Inc.	338
American Airlines Inc.	1539
Delta Air Lines Inc.	651
Frontier Airlines Inc.	575
Hawaiian Airlines Inc.	407
JetBlue Airways	273
SkyWest Airlines Inc.	421
Southwest Airlines Co.	694
United Air Lines Inc.	472
US Airways Inc.	347
Virgin America	366

The last two options make the table a little easier-to-read.

We can further look into the properties of the largest value here for American Airlines Inc. To do so, we can isolate the row corresponding to the arrival delay of

1539 minutes for American in our original `flights` dataset.

```
flights %>% filter(arr_delay == 1539,
                     carrier_name == "American Airlines Inc.") %>%
  select(-c(month, day, carrier, dest_name, hour,
           minute, carrier_name, arr_delay))

##   dep_time dep_delay arr_time tailnum flight dest air_time distance
## 1     1403        1539    1553    1934 N595AA   DFW      182     1616
```

We see that the flight occurred on March 3rd and departed a little after 2 PM on its way to Dallas/Fort Worth. Lastly, we show how we can visualize the arrival delay of all departing flights from Portland on March 3rd against time of departure.

```
flights %>% filter(month == 3, day == 3) %>%
  ggplot(aes(x = dep_time, y = arr_delay)) + geom_point()
```

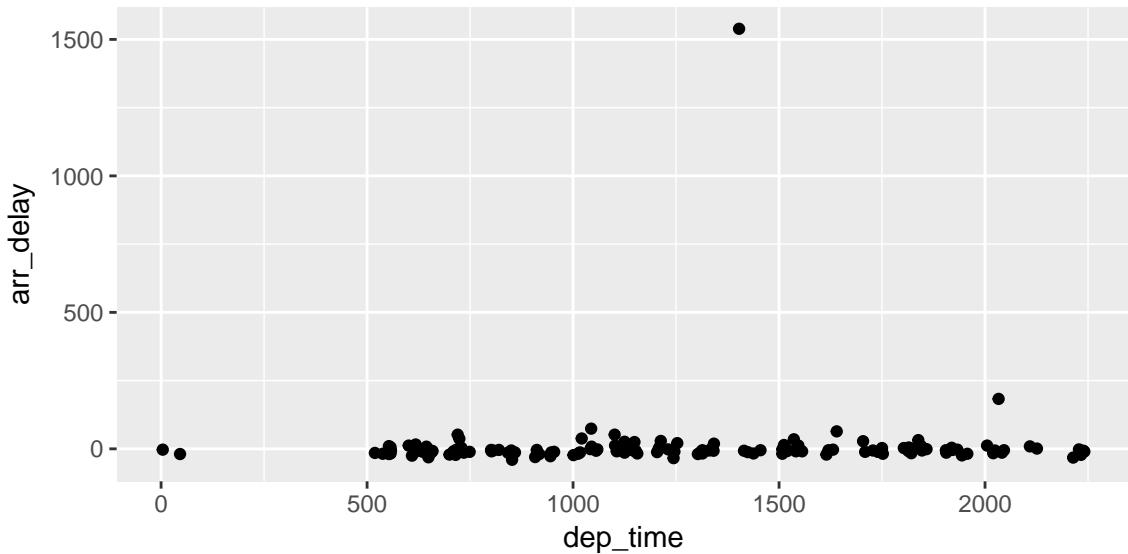


FIGURE 1.1 – β

1.10 Additional resources

- *Markdown Cheatsheet* - <https://github.com/adam-p/markdown-here/wiki/Markdown-Cheatsheet>
- *R Markdown Reference Guide* - <https://www.rstudio.com/wp-content/uploads/2015/03/rmarkdown-reference.pdf>
- Introduction to `dplyr` - <https://cran.rstudio.com/web/packages/dplyr/vignettes/introduction.html>
- `ggplot2` Documentation - <http://docs.ggplot2.org/current/>

Chapitre 2

Adaptation locale

Une population est dite localement adaptée à son environnement si elle a connu une évolution différente de celles qu'ont connu les autres populations de la même espèce, et ce, en réponse des pressions sélectives

2.1 Math

TeX is the best way to typeset mathematics. Donald Knuth designed TeX when he got frustrated at how long it was taking the typesetters to finish his book, which contained a lot of mathematics. One nice feature of *R Markdown* is its ability to read LaTeX code directly.

If you are doing a thesis that will involve lots of math, you will want to read the following section which has been commented out. If you're not going to use math, skip over or delete this next commented section.

2.2 Chemistry 101 : Symbols

Chemical formulas will look best if they are not italicized. Get around math mode's automatic italicizing in LaTeX by using the argument `$\mathsf{formula here}$`, with your formula inside the curly brackets. (Notice the use of the backticks here which enclose text that acts as code.)

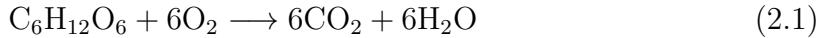
So, $\text{Fe}_2^{2+}\text{Cr}_2\text{O}_4$ is written `$\mathsf{Fe_2^{2+}Cr_20_4}$`.
Exponent or Superscript : O^-
Subscript : CH_4

To stack numbers or letters as in Fe_2^{2+} , the subscript is defined first, and then the superscript is defined.

Bullet : $\text{CuCl} \bullet 7\text{H}_2\text{O}$
Delta : Δ
Reaction Arrows : \longrightarrow or $\xrightarrow{\text{solution}}$
Resonance Arrows : \leftrightarrow
Reversible Reaction Arrows : \rightleftharpoons

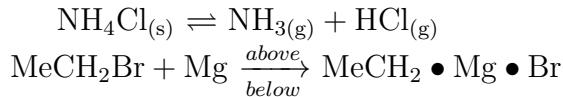
2.2.1 Typesetting reactions

You may wish to put your reaction in an equation environment, which means that LaTeX will place the reaction where it fits and will number the equations for you.



We can reference this combustion of glucose reaction via Equation (2.1).

2.2.2 Other examples of reactions



2.3 Physics

Many of the symbols you will need can be found on the math page <http://web.reed.edu/cis/help/latex/math.html> and the Comprehensive LaTeX Symbol Guide (<http://mirror.utexas.edu/ctan/info/symbols/comprehensive/symbols-letter.pdf>).

2.4 Biology

You will probably find the resources at <http://www.lecb.ncifcrf.gov/~toms/latex.html> helpful, particularly the links to bsts for various journals. You may also be interested in TeXShade for nucleotide typesetting (<http://homepages.uni-tuebingen.de/beitz/txe.html>). Be sure to read the proceeding chapter on graphics and tables.

Chapitre 3

Introgression adaptative

3.1 Coefficients de métissage globaux et locaux

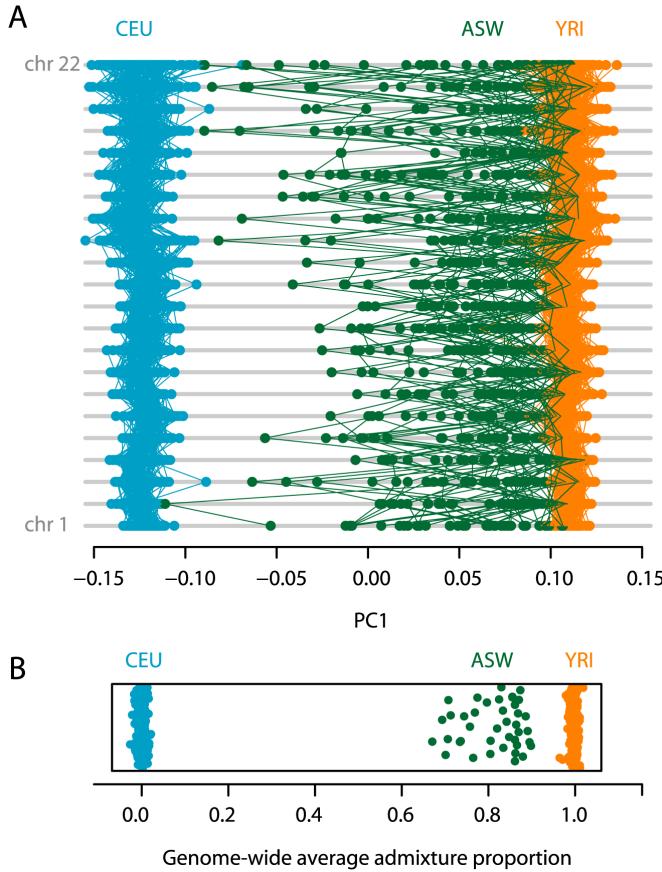
Étant données des populations ancestrales, il est possible d'estimer pour un individu donné, la proportion de son génome provenant de chacune des populations ancestrales. Ces proportions sont connues plus communément sous le nom de *coefficients de métissage globaux*. De nombreux logiciels existent pour l'estimation de ces coefficients : STRUCTURE, ADMIXTURE (Alexander, 2009), LEA (Frichot, 2015), tess3r (Caye, 2016). En complément à cette information globale, il peut être intéressant de déterminer sur des portions plus petites du génome, de la même manière que dans le cas global, les proportions venant de telle ou telle population ancestrale pour chacune de ces portions. Nous parlons dans ce cas de *coefficients de métissage locaux*. Encore une fois, plusieurs logiciels ont été proposés dans le but d'estimer ces coefficients : Hapmix (Price, 2009), EILA (Yang, 2013), LAMP (Thornton, 2014), loter ou encore RFmix (Maples, 2013).

3.2 Introgression

L'introgression peut être détectée de différentes façons. Une première approche consiste à utiliser les *coefficients de métissage locaux*. Les méthodes mentionnées plus haut estiment ces coefficients pour chaque individu, permettant de calculer à partir de ceux-ci des coefficients de métissage locaux pour chaque population.

3.3 Lien entre Analyse en Composantes Principales et métissage global.

L'un des premiers articles à établir un lien entre l'ACP et les coefficients de métissage global fut sur l'interprétation généalogique de l'ACP de Gil McVean (McVean, 2009) :



Pour chacun des 22 chromosomes,

3.4 Analyse en Composantes Principales locale

Notant p le nombre de marqueurs génétiques, i un entier compris entre 1 et p , et x_i la position génétique (en Morgans) ou la position physique (en paires de bases) du i -ème marqueur génétique. Nous définissons pour cet entier i la fenêtre W_i^T de taille T et centrée en i :

$$W_i^T = \{j \in [1, p], |x_i - x_j| \leq T/2\}$$

3.5 Sensibilité à l'imputation des données manquantes

3.6 Simulations

3.6.1 Données de peupliers

Le premier jeu de données est issu d'une étude d'introgression adaptative chez les peupliers d'Amérique du Nord (Suarez-Gonzalez, 2016). La simulation d'haplotypes

d'individus admixés est effectuée à partir des deux populations ancestrales qui y sont présentes. La première, *Populus Balsamifera*, est une espèce de peupliers qui peuple le nord du continent nord-américain, d'Est en Ouest, et se trouve exposée à des conditions climatiques peu clémentes. La seconde, *Populus Trichocarpa*, est principalement localisée en Californie, et bénéficie d'un climat continental.

Chacune des simulations est constituée de 50 haplotypes de la souche continentale, de 50 haplotypes de la souche boréale, ainsi que de 50 haplotypes d'individus hybrides générés à partir des haplotypes ancestraux. Ces haplotypes ancestraux ont été estimés à l'aide du logiciel Beagle. À partir des positions en paires de base, une carte de recombinaison génétique est générée en utilisant le taux de recombinaison moyen chez le peuplier. Le taux de recombinaison, noté τ_r , correspond au nombre moyen de paires de bases à parcourir pour qu'ait lieu un épisode de recombinaison génétique, *i.e.*, notant L la longueur du chromosome en Morgans (M), et N_{bp} le nombre de paires de bases le constituant, le taux de recombinaison génétique pour ce chromosome est donné par la relation :

$$\tau_r = \frac{L}{N_{bp}}$$

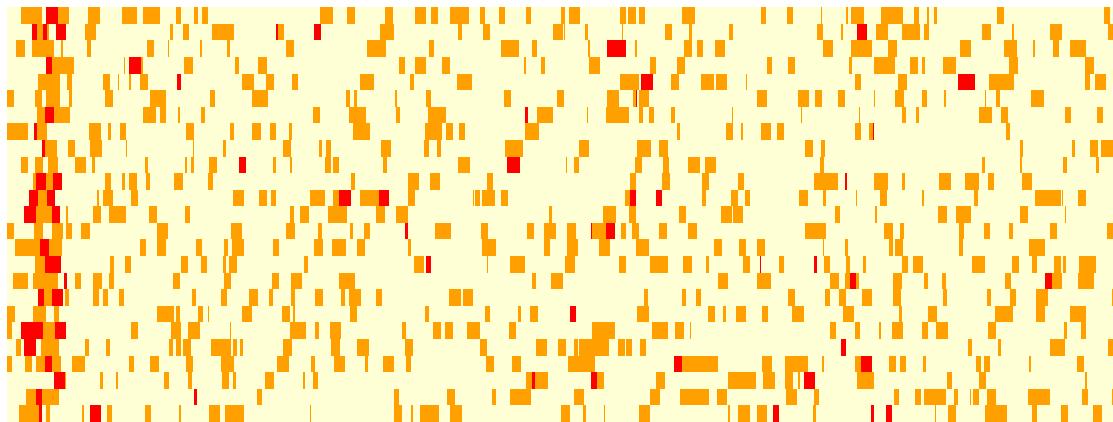
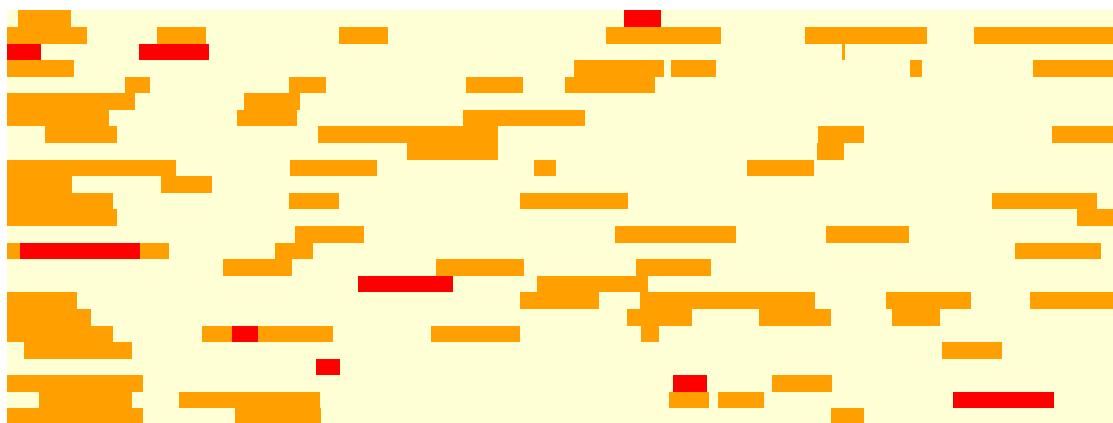
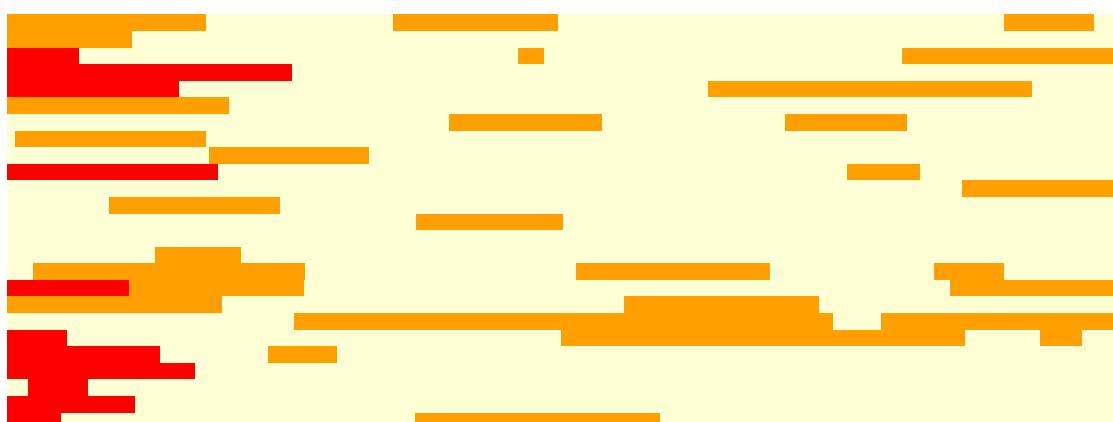
Dans ce scénario, les simulations ont été produites en utilisant un taux de recombinaison génétique moyen τ_r de 0.05 centiMorgans par million de paire de bases, correspondant à la valeur utilisée par les auteurs de l'étude avec le logiciel RASPberry (*Recombination via Ancestry Switch Probability*).

```
path <- "~/Documents/thesis/git/simulations/introgression/"
output.name <- "populus"
recombinationRate <- 0.05 # in Morgans per Megabase
nSNP <- 50000
ancctrl.1 <- 1
ancctrl.2 <- 3
hyb <- 4
intro.size <- 500
global.ancestry <- 0.1
inverted.ancestry <- 0.5

info.map <- as.matrix(data.table::fread(paste0(path, output.name, ".map"),
                                         data.table = FALSE))
H1 <- as.matrix(data.table::fread(paste0(path, output.name, "_H1"),
                                   data.table = FALSE))
H2 <- as.matrix(data.table::fread(paste0(path, output.name, "_H2"),
                                   data.table = FALSE))
n.hyb <- ncol(H1) / 2

### Introgression region
idx <- sample(1:nSNP, size = 1)
beg.reg <- max(1, idx - intro.size)
```

```
end.reg <- min(nSNP, idx + intro.size)
intro.reg <- beg.reg:end.reg
```

FIGURE 3.1 – $\lambda = 0.1$ FIGURE 3.2 – $\lambda = 10$ FIGURE 3.3 – $\lambda = 50$

3.6.2 Résultats de la comparaison des logiciels

```

setwd("~/Documents/thesis/git/simulations/introgression/")
output.name <- "populus"
recombinationRate <- 0.05 # in Morgans per Megabase
nSNP <- 50000
ancctrl.1 <- 1
ancctrl.2 <- 3
hyb <- 4
intro.size <- 500
global.ancestry <- 0.1
inverted.ancestry <- 0.5
N <- 10
pop <- c(rep(ancctrl.1, ncol(H1) / 2),
         rep(ancctrl.2, ncol(H2) / 2),
         rep(4, n.hyb))
results <- data.frame(Software = c("pcadapt", "eila", "RFMix"), Power = c(0, 0, 0), FDR
info.map <- as.matrix(data.table::fread(paste0(path, output.name, ".map")),
                       data.table = FALSE))

compute.fdr = function(list, ground.truth){
  if (length(list) == 0){
    x <- 0
  } else {
    x <- sum(!(list %in% ground.truth)) / length(list)
  }
  return(x)
}

compute.power = function(list, ground.truth){
  if (length(ground.truth) == 0){
    warning("The list of true positives is empty.")
  } else {
    x <- sum(list %in% ground.truth) / length(ground.truth)
  }
  return(x)
}

for (n.simu in 21:30){
  dir.name <- paste0("RFMix_v1.5.4/simu", n.simu, "/")

  input.pcadapt <- as.matrix(data.table::fread(paste0(dir.name, "simu.pcadapt")), data.table = FALSE)
  input.eila <- simulate::eila_from_pcadapt(input.pcadapt, pop, anc1 = ancctrl.1, anc2 = ancctrl.2)
  param <- read.table(paste0(dir.name, "/parameters.txt"))
}

```

```

gt <- (param$begin):(param$end)

#### run pcadapt
wsize <- 1000
mmaf <- 0.01
nomap <- 1:nSNP
maf <- pcadapt::cmpt_minor_af(input.pcadapt, 2)
proxy.map <- info.map[1:nSNP]
filtered.map <- nomap[maf >= mmaf]
stat.pcadapt <- pcadapt::scan.intro(input.pcadapt, K = 1, pop = pop,
                                      ancstrl.1 = ancstrl.2,
                                      ancstrl.2 = ancstrl.1,
                                      admxd = hyb,
                                      min.maf = mmaf,
                                      window.size = wsize,
                                      ploidy = 2,
                                      side = "middle",
                                      map = nomap)

#### run eila
obj.eila <- EILA::eila(admixed = input.eila$admixed, anc1 = input.eila$anc1,
                        anc2 = input.eila$anc2, position = info.map[, 1], lambda
loc.anc.eila <- simulate::haplo_to_ancestry(obj.eila$local.ancestry, 1)

#### run rfmix
allele <- paste0("./simu", n.simu, "/rfmix_alleles.txt")
classes <- paste0("./simu", n.simu, "/rfmix_classes.txt")
markerLocation <- paste0("./simu", n.simu, "/rfmix_markerLocation.txt")
output <- paste0("simu", n.simu, "/output_simu", n.simu)
window.rfmix <- 0.00002
command <- paste("python2.7 RunRFMix.py PopPhased", allele, classes, markerLoc
setwd("~/Documents/thesis/git/simulations/introgression/RFMix_v1.5.4/")
system(command = command)
setwd("~/Documents/thesis/git/simulations/introgression/")
aux.rfmix <- simulate::rfmix.local.ancestry(paste0("RFMix_v1.5.4/simu", n.simu,
loc.anc.rfmix <- simulate::haplo_to_ancestry(aux.rfmix, 1)

#### FDR
interp <- approx(filtered.map, stat.pcadapt[[1]], 1:nSNP)
sd.pcadapt <- sd(interp$y, na.rm = TRUE)
list.pcadapt <- which(interp$y > 3)
results[1, 3] <- results[1, 3] + compute.fdr(list.pcadapt, gt) / N
results[1, 2] <- results[1, 2] + compute.power(list.pcadapt, gt) / N

```

```

sd.eila <- sd(loc.anc.eila, na.rm = TRUE)
stat.eila <- (loc.anc.eila - mean(loc.anc.eila)) / sd.eila
list.eila <- which(stat.eila > 3)
results[2, 3] <- results[2, 3] + compute.fdr(list.eila, gt) / N
results[2, 2] <- results[2, 2] + compute.power(list.eila, gt) / N

sd.rfmix <- sd(loc.anc.rfmix, na.rm = TRUE)
stat.rfmix <- (loc.anc.rfmix - mean(loc.anc.rfmix)) / sd.rfmix
list.rfmix <- which(stat.rfmix > 3)
results[3, 3] <- results[3, 3] + compute.fdr(list.rfmix, gt) / N
results[3, 2] <- results[3, 2] + compute.power(list.rfmix, gt) / N
}

ggres <- data.frame(Software = rep(c("pcadapt", "eila", "RFMix"), 2), Stat = rep(0, 6),
                      Percent = rep(0, 6))
ggres$Stat[1:3] <- results$Power * 100
ggres$Stat[4:6] <- results$FDR * 100
ggres$Percent <- as.numeric(format(ggres$Stat, digits = 2))
p0 <- ggplot(ggres, aes(x = Software, y = Stat, fill = as.factor(Type)))
p0 <- p0 + ggtile(expression(lambda == 1)) + ylab("")
p0 <- p0 + geom_bar(stat = "identity", position = position_dodge(width = 0.9))
p0 <- p0 + guides(fill = guide_legend(title = ""))
p0 <- p0 + geom_text(aes(label = Percent), position = position_dodge(width = 0.9),
                      color = "white", vjust = 1.4, size = 5)
p0 <- p0 + theme_bw() + theme(axis.text = element_text(size = 15),
                                axis.title = element_text(size = 15, face = "bold"),
                                title = element_text(size = 15, face = "bold"),
                                legend.text = element_text(size = 15),
                                legend.key.height = unit(1, "line"),
                                legend.key.width = unit(3, "line"))
)
print(p0)

```

TABLE 3.1 – Correlation of Inheritance Factors for Parents and Child

Factors	Correlation between Parents & Child	Inherited
Education	-0.49	Yes
Socio-Economic Status	0.28	Slight
Income	0.08	No
Family Size	0.18	Slight
Occupational Prestige	0.21	Slight

We can also create a link to the table by doing the following : Table 3.1. If you go back to Loading and exploring data and look at the **kable** table, we can create a reference to this max delays table too : Table 1.1. The addition of the (`\#tab:inher`) option to the end of the table caption allows us to then make a reference to Table `\@ref(tab:label)`. Note that this reference could appear anywhere throughout the document after the table has appeared.

3.7 Figures

If your thesis has a lot of figures, *R Markdown* might behave better for you than that other word processor. One perk is that it will automatically number the figures accordingly in each chapter. You'll also be able to create a label for each figure, add a caption, and then reference the figure in a way similar to what we saw with tables earlier. If you label your figures, you can move the figures around and *R Markdown* will automatically adjust the numbering for you. No need for you to remember! So that you don't have to get too far into LaTeX to do this, a couple **R** functions have been created for you to assist. You'll see their use below.

In the **R** chunk below, we will load in a picture stored as `reed.jpg` in our main directory. We then give it the caption of “Reed logo”, the label of “reedlogo”, and specify that this is a figure. Make note of the different **R** chunk options that are given in the R Markdown file (not shown in the knitted document).

```
include_graphics(path = "figure/reed.jpg")
```



FIGURE 3.4 – Reed logo

Here is a reference to the Reed logo : Figure 3.4. Note the use of the `fig:` code here. By naming the **R** chunk that contains the figure, we can then reference that figure later as done in the first sentence here. We can also specify the caption for the figure via the R chunk option `fig.cap`.

Below we will investigate how to save the output of an **R** plot and label it in a way similar to that done above. Recall the `flights` dataset from Chapter ???. (Note that we've shown a different way to reference a section or chapter here.) We will next explore a bar graph with the mean flight departure delays by airline from Portland for 2014. Note also the use of the `scale` parameter which is discussed on the next page.

```
flights %>% group_by(carrier) %>%
  summarize(mean_dep_delay = mean(dep_delay)) %>%
  ggplot(aes(x = carrier, y = mean_dep_delay)) +
  geom_bar(position = "identity", stat = "identity", fill = "red")
```

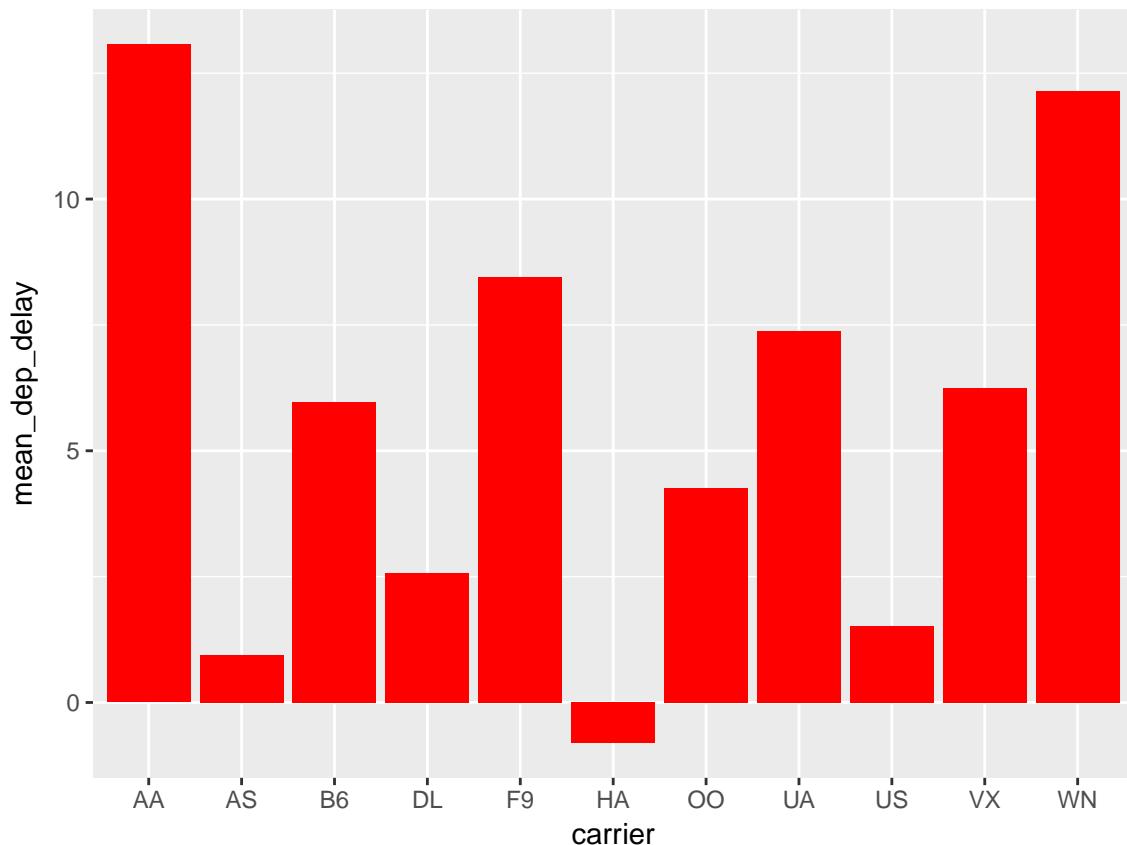


FIGURE 3.5 – Mean Delays by Airline

Here is a reference to this image : Figure 3.5.

A table linking these carrier codes to airline names is available at <https://github.com/ismayc/pnwflights14/blob/master/data/airlines.csv>.

Next, we will explore the use of the `out.extra` chunk option, which can be used to shrink or expand an image loaded from a file by specifying "`scale=`". Here we use the mathematical graph stored in the "subdivision.pdf" file.

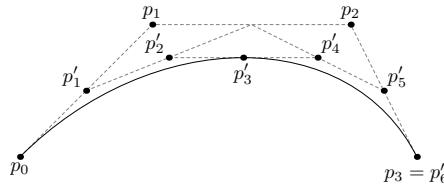


FIGURE 3.6 – Subdiv. graph

Here is a reference to this image : Figure 3.6. Note that `echo=FALSE` is specified so that the **R** code is hidden in the document.

More Figure Stuff

Lastly, we will explore how to rotate and enlarge figures using the `out.extra` chunk option. (Currently this only works in the PDF version of the book.)

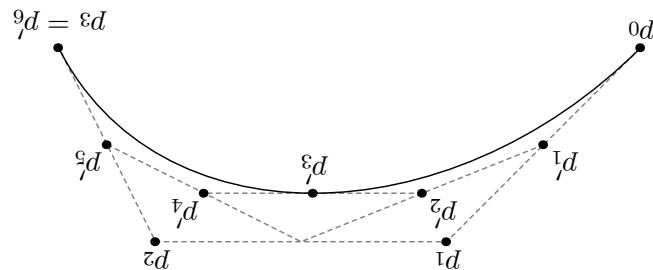


FIGURE 3.7 – A Larger Figure, Flipped Upside Down

As another example, here is a reference : Figure 3.7.

3.8 Footnotes and Endnotes

You might want to footnote something.¹ The footnote will be in a smaller font and placed appropriately. Endnotes work in much the same way. More information can be found about both on the CUS site or feel free to reach out to `data@reed.edu`.

3.9 Bibliographies

Of course you will need to cite things, and you will probably accumulate an armful of sources. There are a variety of tools available for creating a bibliography database (stored with the `.bib` extension). In addition to BibTeX suggested below, you may want to consider using the free and easy-to-use tool called Zotero. The Reed librarians have created Zotero documentation at <http://libguides.reed.edu/>

1. footnote text

`citation/zotero`. In addition, a tutorial is available from Middlebury College at <http://sites.middlebury.edu/zoteromiddlebury/>.

R Markdown uses *pandoc* (<http://pandoc.org/>) to build its bibliographies. One nice caveat of this is that you won't have to do a second compile to load in references as standard LaTeX requires. To cite references in your thesis (after creating your bibliography database), place the reference name inside square brackets and precede it by the “at” symbol. For example, here's a reference to a book about worrying : (Molina & Borkovec, 1994). This Molina1994 entry appears in a file called `thesis.bib` in the `bib` folder. This bibliography database file was created by a program called BibTeX. You can call this file something else if you like (look at the YAML header in the main .Rmd file) and, by default, is placed in the `bib` folder.

For more information about BibTeX and bibliographies, see our CUS site (<http://web.reed.edu/cis/help/latex/index.html>)². There are three pages on this topic : `bibtex` (which talks about using BibTeX, at <http://web.reed.edu/cis/help/latex/bibtex.html>), `bibtextstyles` (about how to find and use the bibliography style that best suits your needs, at <http://web.reed.edu/cis/help/latex/bibtextstyles.html>) and `bibman` (which covers how to make and maintain a bibliography by hand, without BibTeX, at <http://web.reed.edu/cis/help/latex/bibman.html>). The last page will not be useful unless you have only a few sources.

If you look at the YAML header at the top of the main .Rmd file you can see that we can specify the style of the bibliography by referencing the appropriate csl file. You can download a variety of different style files at <https://www.zotero.org/styles>. Make sure to download the file into the csl folder.

Tips for Bibliographies

- Like with thesis formatting, the sooner you start compiling your bibliography for something as large as thesis, the better. Typing in source after source is mind-numbing enough ; do you really want to do it for hours on end in late April? Think of it as procrastination.
- The cite key (a citation's label) needs to be unique from the other entries.
- When you have more than one author or editor, you need to separate each author's name by the word “and” e.g. `Author = {Noble, Sam and Youngberg, Jessica},`,
- Bibliographies made using BibTeX (whether manually or using a manager) accept LaTeX markup, so you can italicize and add symbols as necessary.
- To force capitalization in an article title or where all lowercase is generally used, bracket the capital letter in curly braces.
- You can add a Reed Thesis citation³ option. The best way to do this is to use the `phdthesis` type of citation, and use the optional “type” field to enter “Reed thesis” or “Undergraduate thesis.”

2. Reed College (2007)

3. Noble (2002)

3.10 Anything else ?

If you'd like to see examples of other things in this template, please contact the Data @ Reed team (email `data@reed.edu`) with your suggestions. We love to see people using *R Markdown* for their theses, and are happy to help.

Conclusion

If we don't want Conclusion to have a chapter number next to it, we can add the `{-}` attribute.

More info

And here's some other random info : the first paragraph after a chapter title or section head *shouldn't be* indented, because indents are to tell the reader that you're starting a new paragraph. Since that's obvious after a chapter or section title, proper typesetting doesn't add an indent there.

Annexe A

The First Appendix

This first appendix includes all of the R chunks of code that were hidden throughout the document (using the `include = FALSE` chunk tag) to help with readability and/or setup.

In the main Rmd file

```
# This chunk ensures that the thesisdown package is
# installed and loaded. This thesisdown package includes
# the template files for the thesis.
if(!require(devtools))
  install.packages("devtools", repos = "http://cran.rstudio.com")
if(!require(thesisdown))
  devtools::install_github("ismayc/thesisdown")
library(thesisdown)
```

In Chapter ?? :

```
# This chunk ensures that the thesisdown package is
# installed and loaded. This thesisdown package includes
# the template files for the thesis and also two functions
# used for labeling and referencing
opts_chunk$set(cache=TRUE)

if(!require(devtools))
  install.packages("devtools", repos = "http://cran.rstudio.com")
if(!require(dplyr))
  install.packages("dplyr", repos = "http://cran.rstudio.com")
if(!require(ggplot2))
  install.packages("ggplot2", repos = "http://cran.rstudio.com")
if(!require(bookdown))
  install.packages("bookdown", repos = "http://cran.rstudio.com")
if(!require(thesisdown)){
  library(devtools)
```

```
devtools::install_github("ismayc/thesisdown")
}
library(thesisdown)
library(pcadapt)
library(EILA)
library(simulate)
flights <- read.csv("data/flights.csv")
```

Annexe B

The Second Appendix, for Fun

References

- Alexander, D. (2009). *Fast model-based estimation of ancestry in unrelated individuals*.
- Angel, E. (2000). *Interactive computer graphics : A top-down approach with opengl*. Boston, MA : Addison Wesley Longman.
- Angel, E. (2001a). *Batch-file computer graphics : A bottom-up approach with quicktime*. Boston, MA : Wesley Addison Longman.
- Angel, E. (2001b). *Test second book by angel*. Boston, MA : Wesley Addison Longman.
- Caye, K. (2016). *TESS3 : Fast inference of spatial population structure and genome scans for selection*.
- Frichot, É. (2015). *LEA : An r package for landscape and ecological association studies*.
- Maples, B. K. (2013). *RFMix : A discriminative modeling approach for rapid and robust local-ancestry inference*.
- McVean, G. (2009). A genealogical interpretation of principal components analysis.
- Molina, S. T., & Borkovec, T. D. (1994). The Penn State worry questionnaire : Psychometric properties and associated characteristics. In G. C. L. Davey & F. Tallis (Eds.), *Worrying : Perspectives on theory, assessment and treatment* (pp. 265–283). New York : Wiley.
- Noble, S. G. (2002). *Turning images into simple line-art* (Undergraduate thesis). Reed College.
- Price, A. L. (2009). *Sensitive detection of chromosomal segments of distinct ancestry in admixed populations*.
- Reed College. (2007, march). LaTeX your document. Retrieved from <http://web.reed.edu/cis/help/LaTeX/index.html>
- Suarez-Gonzalez, et a., Adriana. (2016). Genomic and functional approaches reveal a case of adaptive introgression from *populus balsamifera* (balsam poplar) in *p. trichocarpa* (black cottonwood). *Molecular Ecology*, 2427–2442.
- Thornton, T. (2014). *Local and global ancestry inference, and applications to genetic*

association analysis for admixed populations.

Yang, J. J. (2013). *Efficient inference of local ancestry.*