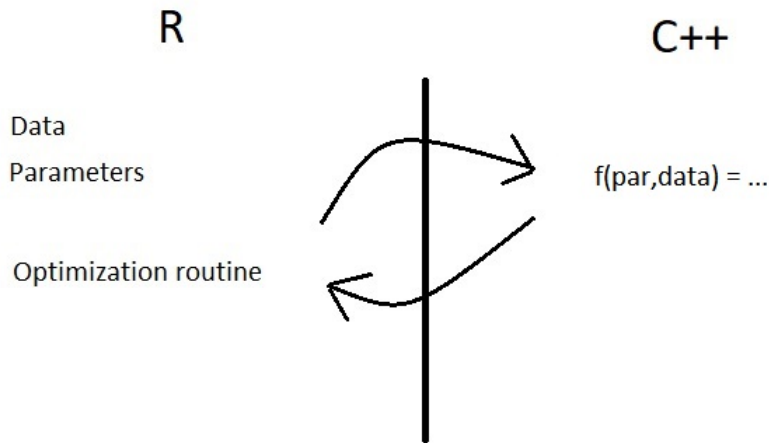# Data and parameters in TMB

Anders Nielsen and Olav Nikolai Breivik

# Data and parameters

# Data and parameters

Simple example:

```
1  library(TMB)
2  compile("scalar.cpp")
3  dyn.load(dynlib("scalar"))
4
5  data = list()
6  data$Y = 5
7
8  par = list()
9  par$mu = 0
10
11 obj = MakeADFun(data,par,DLL = "scalar")
12 opt = nlminb(obj$par,obj$fn,obj$gr)
```

```
1  #include <TMB.hpp>
2  template<class Type>
3  Type objective_function<Type>::operator()()
4  {
5    DATA_SCALAR(Y);
6    PARAMETER(mu);
7    Type nll = pow(Y-mu,2);
8    return nll;
9  }
```

# Transfers basic objects

| What | R side | C++ side |
|---|---|---|
| Number | `1` | `DATA_SCALAR` |
| Vector | `c(1,2,3)` | `DATA_VECTOR` |
| Matrix | `matrix(c(1,2,3,4), nrow=2, ncol=2)` | `DATA_MATRIX` |
| Array | `matrix(c(1,2,3,4), nrow=2, ncol=2)` | `DATA_ARRAY` |
| Integer | `1` | `DATA_INTEGER` |
| Integer Vector | `c(1,2,3)` | `DATA_IVECTOR` |
| Integer Matrix | `matrix(c(1,2,3,4), nrow=2, ncol=2)` | `DATA_IMATRIX` |
| Integer Array | `matrix(c(1,2,3,4), nrow=2, ncol=2)` | `DATA_IARRAY` |
| Factor | `factor(c("a","b"))` | `DATA_FACTOR` |
| String | `"a"` | `DATA_STRING` |

# Checking what is read

- Report values on C-side

```
1  library(TMB)
2  compile("verify.cpp")
3  dyn.load(dynlib("verify"))
4  data = list()
5  data$V = 1:3
6  data$M = matrix(1:6,nrow = 2,ncol = 3)
7  data$A = array(1:6, dim = c(1,2,3))
8
9  par = list()
10 par$mu = 0
11 obj = MakeADFun(data,par, DLL = "verify")
12 out = obj$report()
13 out$M
14 #     [,1] [,2] [,3]
15 #[1,]    1    3    5
16 #[2,]    2    4    6
```

```
1  #include <TMB.hpp>
2  template<class Type>
3  Type objective_function<Type>::operator()()
4  {
5    DATA_VECTOR(V);
6    DATA_MATRIX(M);
7    DATA_ARRAY(A);
8    REPORT(V);
9    REPORT(M);
10   REPORT(A);
11
12   PARAMETER(mu);
13   Type nll = pow(mu,2);
14   return nll;
15 }
```

# Indexing from 0

- In $C^{++}$ the first element is **number 0**
- Different from R, so difficult to remember in the beginning

**R**

```
data <- list()
data$y <- c(1.1, 2.2)
data$z <- myMatrix

y[1] ...  y[n]
z[1,1] ...  z[m,n]
```

$C^{++}$

```
DATA_VECTOR(y)
DATA_ARRAY(z)

y(0) ...  y(n-1)
z(0,0) ...  z(m-1,n-1)
```

# Parameters

- The parameters are what we want to estimate
- Set up list on R side
- The C side evaluates the function (and derivatives)
- Optimization is performed from R, so values need to be passed from and to many times
- Simple example:

```
1  library(TMB)
2  compile("scalar.cpp")
3  dyn.load(dynlib("scalar"))
4
5  data = list()
6
7  par = list()
8  par$mu = 0
9
10 obj = MakeADFun(data,par,DLL = "scalar")
11 opt = nlminb(obj$par,obj$fn,obj$gr)
```

```
1  #include <TMB.hpp>
2  template<class Type>
3  Type objective_function<Type>::operator()()
4  {
5    PARAMETER(mu);
6    Type nll = pow(Type(42)-mu,2);
7    return nll;
8  }
```

# Often we have more than one

- The following parameter types are available:

| Template Syntax | C++ type | R type |
|---|---|---|
| PARAMETER(name) | Type | numeric(1) |
| PARAMETER_VECTOR(name) | vector<Type> | vector |
| PARAMETER_MATRIX(name) | matrix<Type> | matrix |
| PARAMETER_ARRAY(name) | array<Type> | array |

- Just like with data we can specify a list of possibly many parameter objects
- Naturally we need to match each parameter object on the C side

# Simple bounds from R

- Consider the model:

$$X \sim \text{Bin}(100, p)$$

- Let's say we observe X's equal 2, 0 and 1
- We want to estimate our model parameter $p$

```
1  library(TMB)
2  compile("p1.cpp")
3  dyn.load(dynlib("p1"))
4  data = list()
5  data$X = c(2,0,1)
6  par = list()
7  par$p = 0.5
8
9  obj = MakeADFun(data,par,DLL = "p1")
10 opt = nlminb(obj$par,obj$fn,obj$gr,
11 lower=c(0),upper=c(1))
12 rep = sdreport(obj)
13 summary(rep)
14 # Estimate      Std. Error
15 #p     0.01      0.005716054
```

```
1  #include <TMB.hpp>
2  template<class Type>
3  Type objective_function<Type>::operator()()
4  {
5    DATA_VECTOR(X);
6    PARAMETER(p);
7    Type nll = -sum(dbinom(X,Type(100),p,true));
8    return nll;
9  }
```

- Now we want to make a 95% confidence interval - see the problem?

# Bounds from C

- Consider the model parametrized as:

$$X \sim \text{Bin}(100, p), \qquad \text{where } \text{logit}(p) = \alpha$$

- Our code is then:

```
1  library(TMB)
2  compile("p2.cpp")
3  dyn.load(dynlib("p2"))
4  data = list()
5  data$X = c(2,0,1)
6  par = list()
7  par$alpha = 0
8
9  obj = MakeADFun(data,par,DLL = "p2")
10 opt = nlminb(obj$par,obj$fn,obj$gr)
11 rep = sdreport(obj)
12 summary(rep)
13 #      Estimate Std. Error
14 #alpha -4.59512  0.5802588
```

```
1  #include <TMB.hpp>
2  template<class Type>
3  Type trans(Type x){
4    return exp(x)/(Type(1) + exp(x));
5  }
6
7  template<class Type>
8  Type objective_function<Type>::operator()()
9  {
10   DATA_VECTOR(X);
11   PARAMETER(alpha);
12   Type p = trans(alpha);
13   Type nll  = -sum(dbinom(X,Type(100),p,true));
14   return nll;
15 }
```

- Now we can make a 95% confidence interval:

```
1  a = -4.59512  + 0.5802588 * c(-2,2)
2  exp(a)/(1 + exp(a))
3  #[1] 0.003154903 0.031231381
```

# Exercise

Exercise: Suggest how to transform parameter:

1. only positive
2. between -1 and 1
3. Increasing positive vector

# **Exercise**

Exercise: Suggest how to transform parameter:

1. only positive
2. between -1 and 1
3. Increasing positive vector

Solution:

1. $\theta = e^{\alpha}$, where $\alpha \in \mathbb{R}$
2. $\theta = 2/(1 + e^{-\alpha}) - 1$, where $\alpha \in \mathbb{R}$
3. $\theta = (e^{\alpha_1}, e^{\alpha_1} + e^{\alpha_2}, ..., e^{\alpha_1} + \cdots + e^{\alpha_n})$, where $\alpha_1, ..., \alpha_n \in \mathbb{R}$

# Getting results out

- If estimated standard errors are not needed:
    - `REPORT(X)` in the $C^{++}$ file
    - `obj$report()$X` in the $R$ file
- If estimated standard errors are needed:
    - `ADREPORT(X)` in the $C^{++}$ file
    - Reported list:
      `rl <- as.list(sdreport(obj),"Est",report= TRUE)`
    - Reported Sd list:
      `rlsd <- as.list(sdreport(obj),"Std",report= TRUE)`
- Parameter estimates and standard deviations:
    - **Parameter list:** `pl <- as.list(sdreport(obj),"Est")`
    - **Parameter Sd list:** `plsd <- as.list(sdreport(obj),"Std")`

# Exercise: Beverton-Holt stock recruitment model

- The Beverton-Holt model can be written (slightly re-parametrized) as:

$$\log R_i = \log(a) + \log(\text{ssb}_i) - \log(1 + \exp(\log(b))\text{ssb}_i) + \varepsilon_i,$$

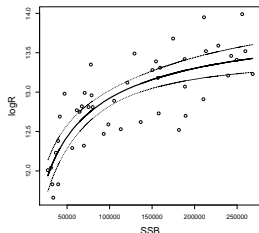where $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$

- A data set of SSB and $\log(R)$ is provided in `bh.dat`
- Code provided in `bh.R` and `bh.cpp`
- Exercise: Estimate the model parameters $\log(a)$ and $\log(b)$ and $\log(\sigma)$.

# Solution, R- and C-side

```
1  dat<-read.table("bh.dat", header=TRUE)
2
3  library(TMB)
4  compile("bh.cpp")
5  dyn.load(dynlib("bh"))
6
7  data <- list(SSB=dat$SSB,logR=dat$logR)
8  parameters <- list(
9    logA=0,
10   logB=0,
11   logSigma=0
12 )
13
14 obj <- MakeADFun(data,parameters,DLL="bh")
15 opt <- nlminb(obj$par,obj$fn,obj$gr)
16 rep <- sdreport(obj)
```

```cpp
1  #include <TMB.hpp>
2
3  template<class Type>
4  Type objective_function<Type>::operator() ()
5  {
6    DATA_VECTOR(SSB)
7    DATA_VECTOR(logR);
8
9    PARAMETER(logA);
10   PARAMETER(logB);
11   PARAMETER(logSigma);
12   vector<Type> pred=logA+log(SSB)-log(Type
         (1)+exp(logB)*SSB);
13   Type ans=-sum(dnorm(logR,pred,exp(logSigma
         ),true));
14   ADREPORT(pred)
15   return ans;
16 }
```



Beverton–Holt

# Collapsing parameters, or fixing them

- The `map` argument of the `MakeADFun` can be used to couple elements in a parameter object
- If we have a parameter vector `alpha` of length 4, then the statement:

```
obj <- MakeADFun(data, param, map=list(alpha=factor(c(1,2,3,3))))
```

  will collapse the last two parameters.

    - Initialized to the mean of the initializations
- If `NA` is included, as in:

```
obj <- MakeADFun(data, param, map=list(alpha=factor(c(1,2,NA,4))))
```

  the optimizer treat that parameter as fixed and equal the initial value.

- Use the map argument to:
    - easily change between different models
    - write the c-side neat

# Map exercise

- Consider the data set `InsectSprays`, which is available in R
- We will use the model

$$\text{count}_i \sim \text{Pois}(\lambda_i), \text{ where } \log \lambda_i = \alpha(\text{spray}_i)$$

- This can be implemented as:

```r
library(TMB)
compile("insect.cpp")
dyn.load(dynlib("insect"))

par <- list()
par$logAlpha=rep(0,nlevels(InsectSprays$
    spray))
obj <- MakeADFun(InsectSprays, par, DLL=
    "insect")
opt <- nlminb(obj$par, obj$fn, obj$gr)
```

```cpp
#include <TMB.hpp>
template<class Type>
Type objective_function<Type>::operator()()
{
  DATA_VECTOR(count);
  DATA_FACTOR(spray);
  PARAMETER_VECTOR(logAlpha);
  Type nll = 0;
  for(int i=0; i<count.size(); ++i){
    Type lambda = exp(logAlpha(spray(i)));
    nll += -dpois(count(i),lambda,true);
  }
  return nll;
}
```

- Use the `map` argument to test if $\alpha(A) = \alpha(B) = \alpha(F)$
- Test the hypothesis $\lambda(A) = \lambda(B) = \lambda(F) = 15$
- Test the hypothesis without modifying the `cpp` file

# Debugging

- Remember that the index starts at 0 in C++ and 1 in R
- Debugging is performed trough
  `TMB::gdbsource("script.R", interactive = TRUE)`
  - A safe version of `source`
- `debug_tutorial.cpp` fails because of index out of bound
- Exercise: Debug `debug_tutorial.cpp`
- Operating system notes:
  - Linux: Works fine
  - Windows: Behaviour of `gdb` depends on versions of R and Rtools.
    See `https://github.com/kaskr/adcomp/wiki/Windows-installation#windows-debugging`
    - Debugging of larger programs may not work in Windows, recommend to debug with Linux
- This exercise is borrowed from `https://github.com/skaug/tmb-case-studies/tree/master/debug_tutorial`

# Some tips

- Check `obj$gr()` if you don't obtain convergence
    - If one element is zero, something is wrong with the implementation.
- Always use same names on the R and C side
- Set `control = list(trace = 1)` to trace the outer fixed parameters in `nlminb()`
- Don't underestimate the importance of writing the C side neat.
- Overview of vector, matrix and array operation in TMB is provided here: `https://kaskr.github.io/adcomp/matrix_arrays_8cpp-example.html`