



# **Data generation in Excel: Some technical guidelines**

Einar Hjörleifsson

# The simulator in math (almost)

$$N_{a,y} = \begin{cases} R_y = N_{A,y}^{STO} = N_{A,y}^{DET} e^{\varepsilon_R} & \varepsilon \sim N(0, \sigma_R^2) \\ N_{a-1,y-1} e^{-(s_{a-1}F_{y-1} + M_{a-1,y-1})} \\ N_{a-1,y-1} e^{-(s_{a-1}F_{y-1} + M_{a-1,y-1})} + N_{a,y-1} e^{-(s_a F_{y-1} + M_{a,y-1})} \end{cases} \quad \begin{cases} a = A \\ A < a \leq a_{plus} \\ a = a_{plus} \end{cases}$$

$$F_{ay} = s_a F_y$$

$$C_{ay}^{TRUE} = \frac{s_a F_y}{s_a F_y + M_{ay}} \left(1 - e^{-(s_a F_y + M_{ay})}\right) N_{ay}$$

$$C_{ay}^{DET} = (1 - p_{ay}) C_{ay}^{TRUE}$$

$$C_{a,y}^{OBS1} = C_{a,y}^{DET} e^{\varepsilon_{C_a}} \quad \varepsilon_{C_a} \sim N(0, \sigma_{C_a}^2)$$

$$C_{a,y}^{OBS2} = C_{a,y}^{OBS1} \frac{Y_y^{OBS2}}{Y_y^{OBS1}}$$

$$Y_y = \sum_a C_{a,y} w_{a,y}^{catch}$$

$$R_y = f(time)$$

$$F_y = f(time)$$

$$s_{ay} = f(age, time)$$

$$q_a = f(age, time)$$

$$M_{ay} = f(?)$$

$$U_{ay}^{TRUE} = q_a N_{ay}^{\beta_a}$$

$$U_{ay}^{OBS} = U_{ay}^{TRUE} e^{\varepsilon_{U_a}}$$

$$\varepsilon_{U_a} \sim N(0, \sigma_{U_a})$$



# Generating data: Using functions

- Empirical specifications
  - Simplest way to generate data is to punch in some starting numbers for each parameter
  - Cumbersome and inefficient, need to punch in a lot of values. But sometimes valuable.
- Use functions
  - Many commonly observed time trends in fish stocks can be described by 2-5 parametric functions

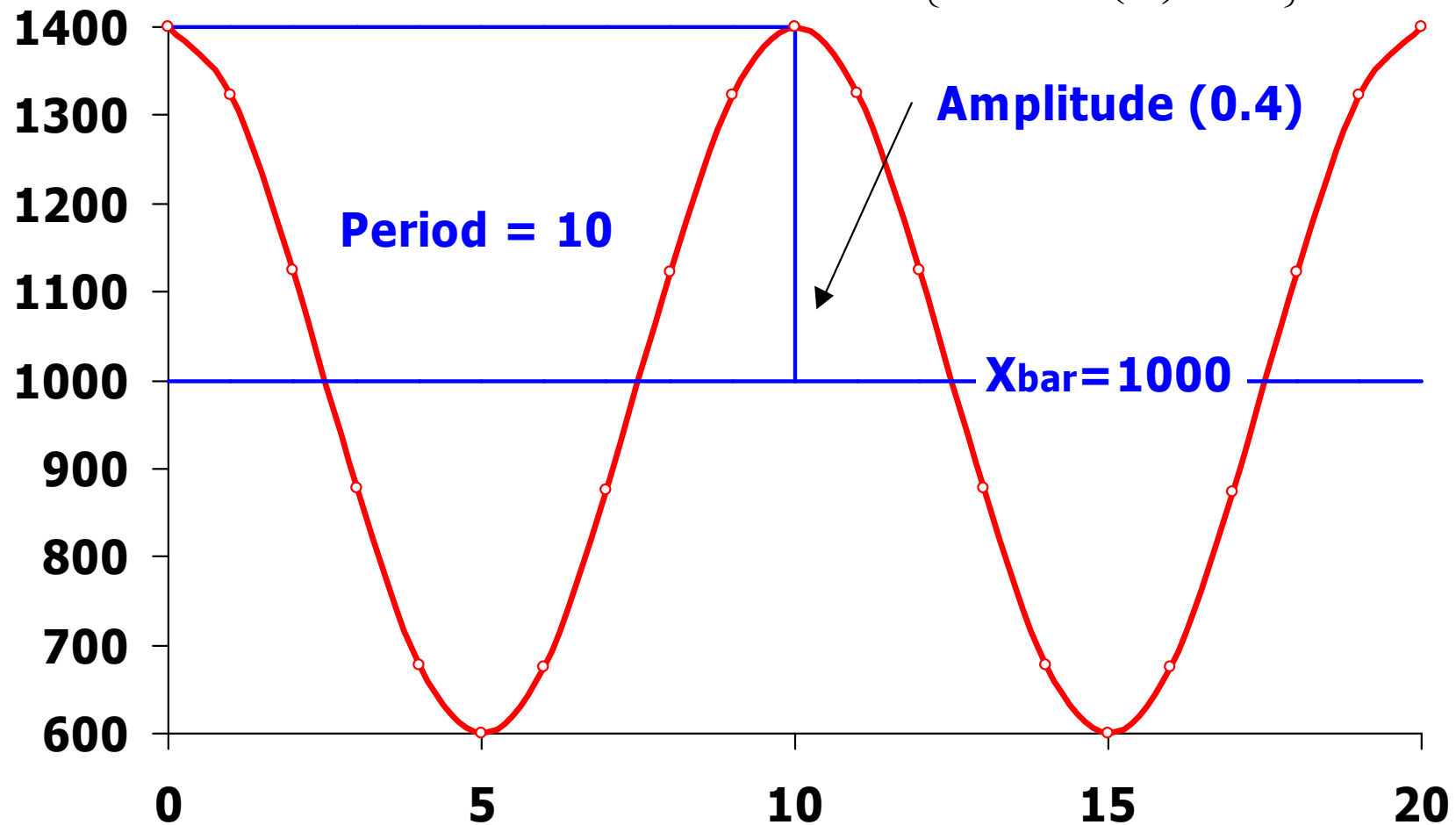
# Periodic function

$$Y_t = \bar{Y} \left\{ 1 + A * \cos\left(\frac{T}{P}\right) * 6.28 \right\}$$

- $\bar{Y}$ : Mean level (over the time of 1 period)
- $A$ : Amplitude
- $T$ : Time
- $P$ : Period
- 6.28 -- conversion factor

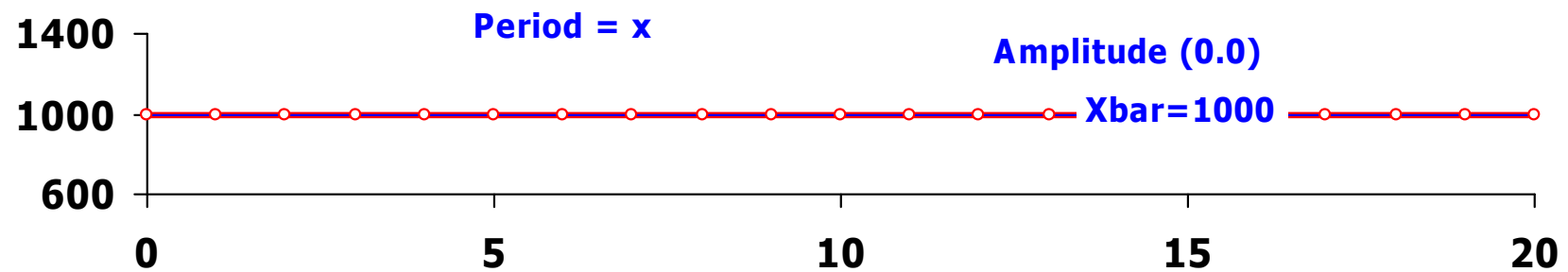
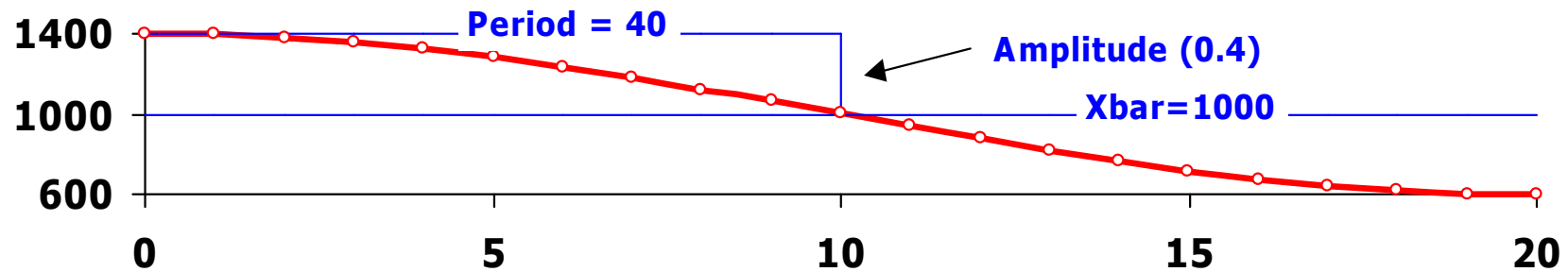
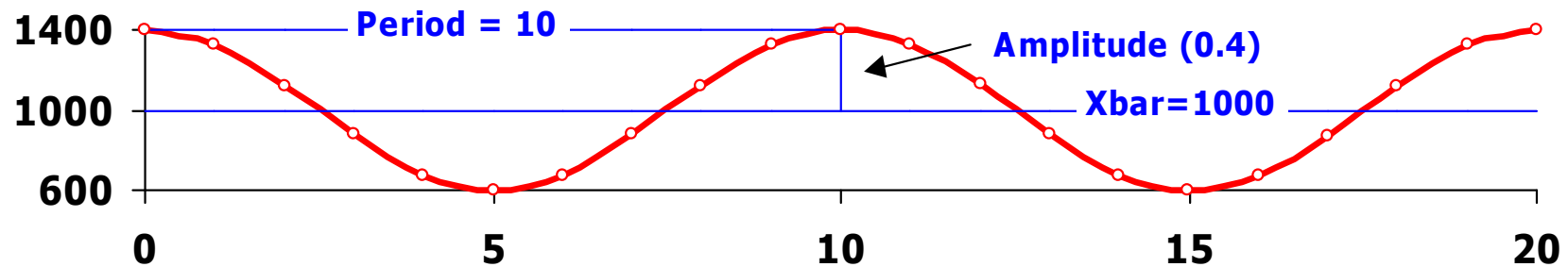
# Periodic function: What it looks like

$$Y_t = \bar{Y} \left\{ 1 + A * \cos\left(\frac{T}{P}\right) * 6.28 \right\}$$



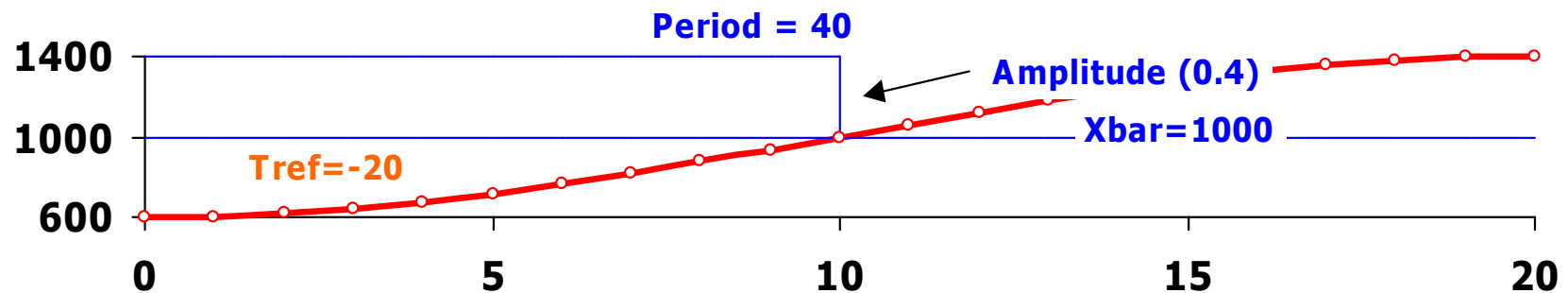
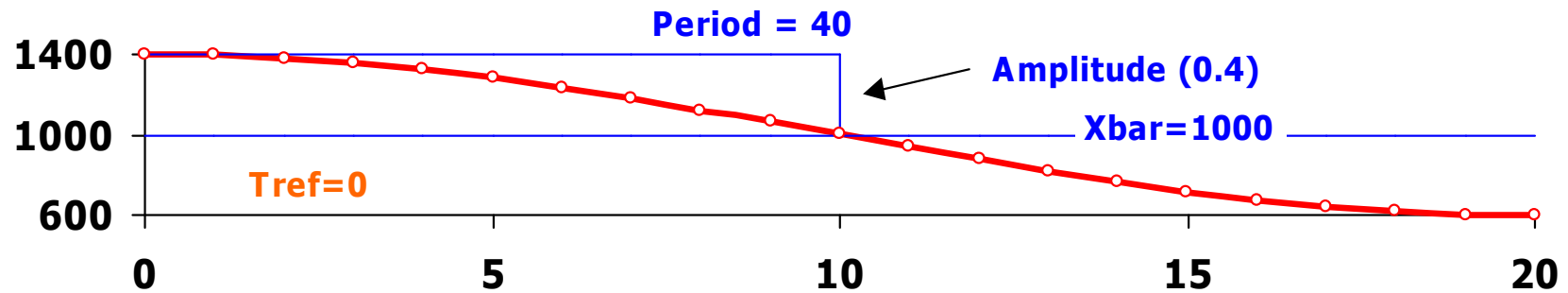
# Periodic function: Very flexible

$$Y_t = \bar{Y} \left\{ 1 + A * \cos\left(\frac{T}{P}\right) * 6.28 \right\}$$



# Periodic function: Phase shift

$$Y_t = \bar{Y} \left\{ 1 + A * \cos \left( \frac{T - T_{ref}}{P} \right) * 6.28 \right\}$$

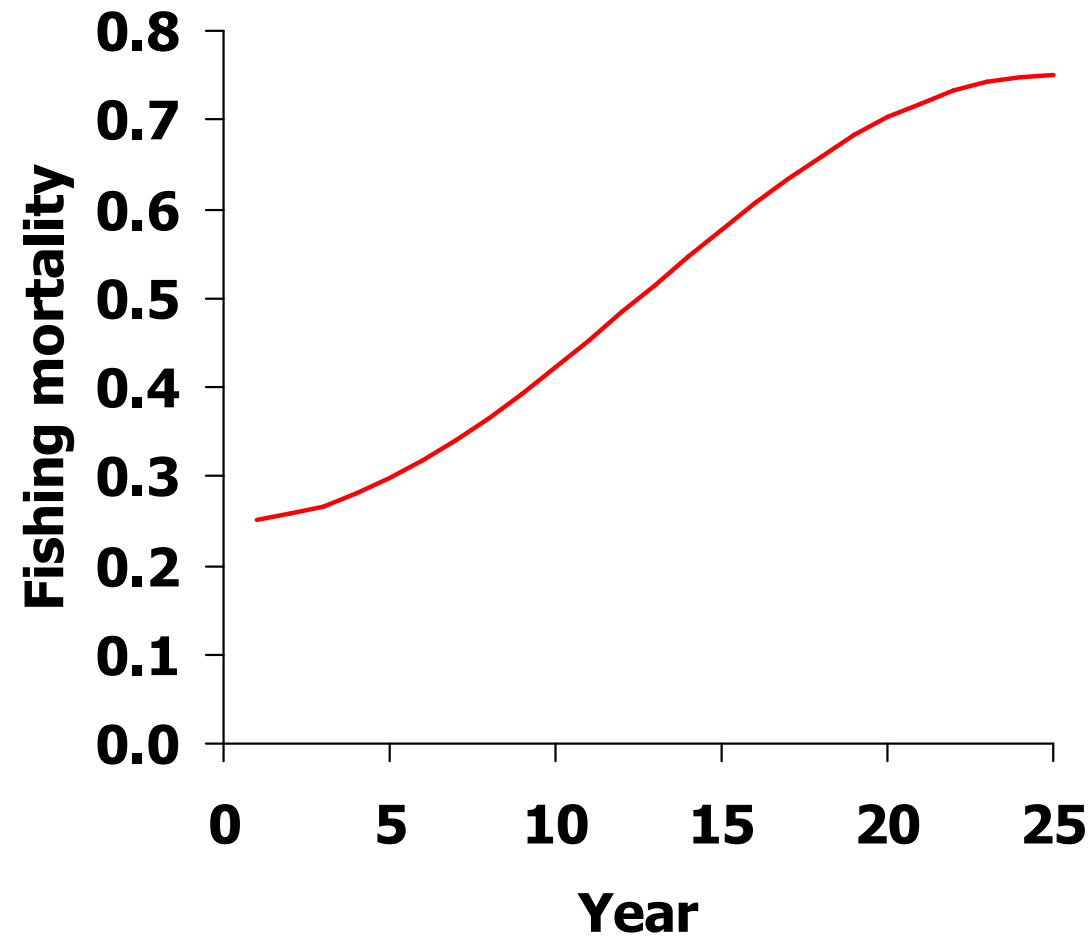




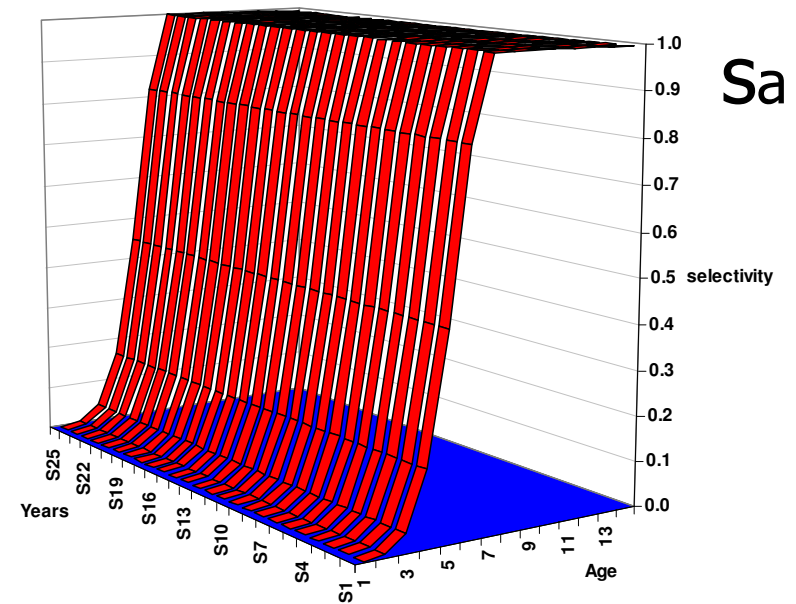
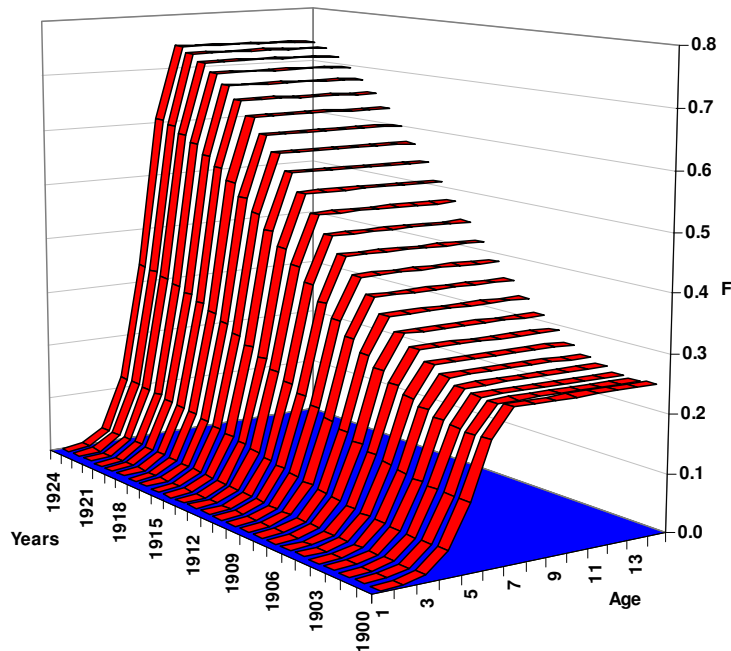
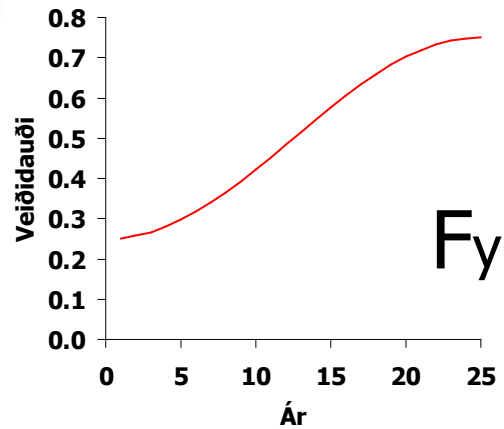
# Development in fishing mortality ( $F_y$ )

Fishing mortality	
Parameters	
A	0.5
P	50
Ybar	0.5
Tshift	-25

Year	$F_y$
1	0.25
2	0.26
3	0.27
4	0.28
5	0.30
6	0.32
7	0.34
8	0.37
9	0.39
10	0.42
11	0.45
12	0.48
13	0.52
14	0.55
15	0.58
16	0.61
17	0.63
18	0.66
19	0.68
20	0.70
21	0.72
22	0.73
23	0.74
24	0.75
25	0.75



# Fishing mortality by age and time: $F_{ay}$



$$F_{ay} = S_a F_y$$

Note: Here selection pattern fixed

# Selectivity - Logistic function

- Various forms describing the same thing.

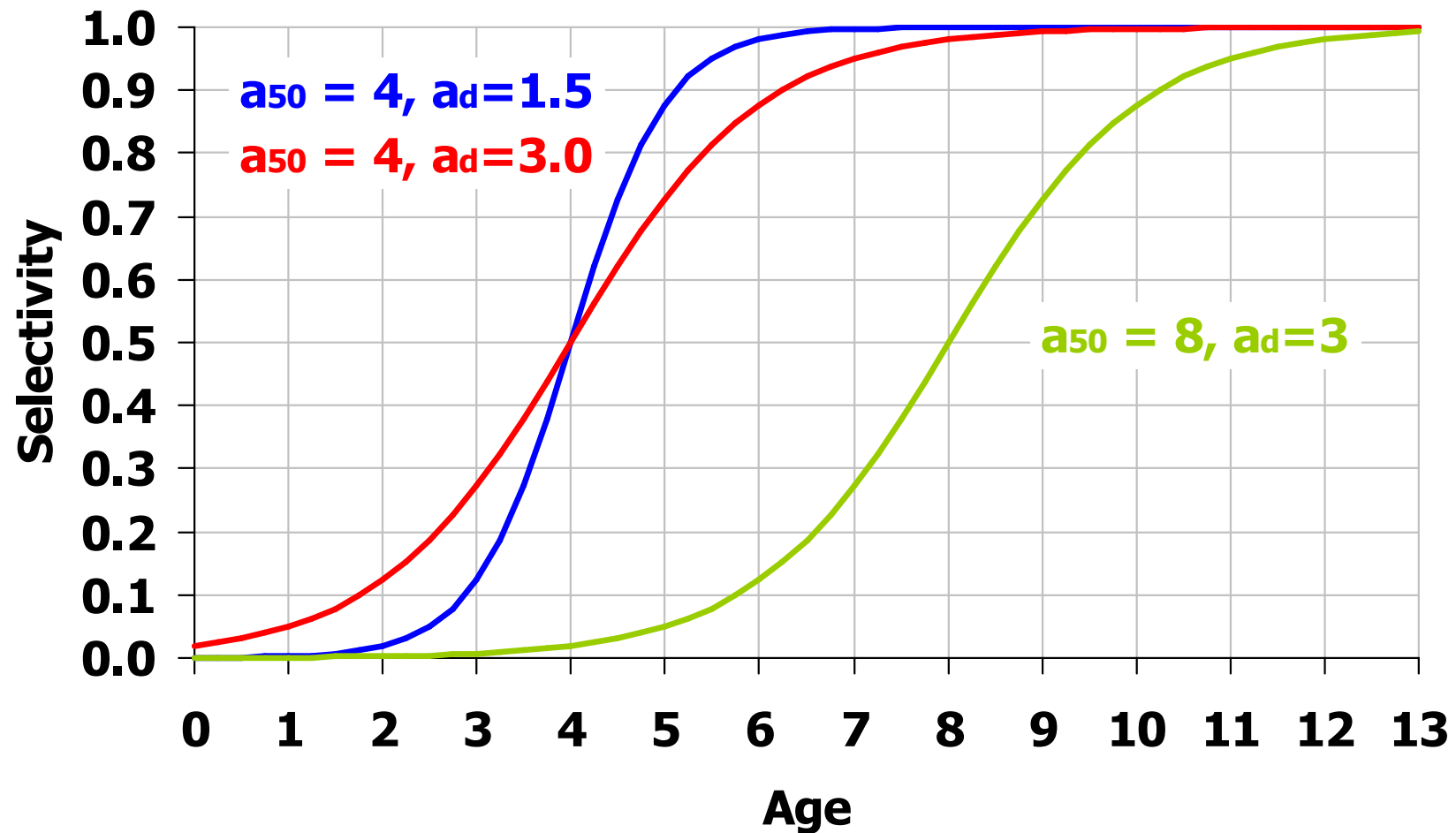
$$s_a = \frac{1}{1 + e^{-k(a - a_{50})}}$$

$$s_a = \frac{1}{1 + e^{-\ln 19 \frac{a - a_{50}}{a_{95} - a_{50}}}}$$

$$s_a = \frac{1}{1 + e^{-\ln 19 \frac{a - a_{50}}{a_d}}}$$

# Selectivity - Logistic

Note symmetry



# Selectivity - double half-Gaussian

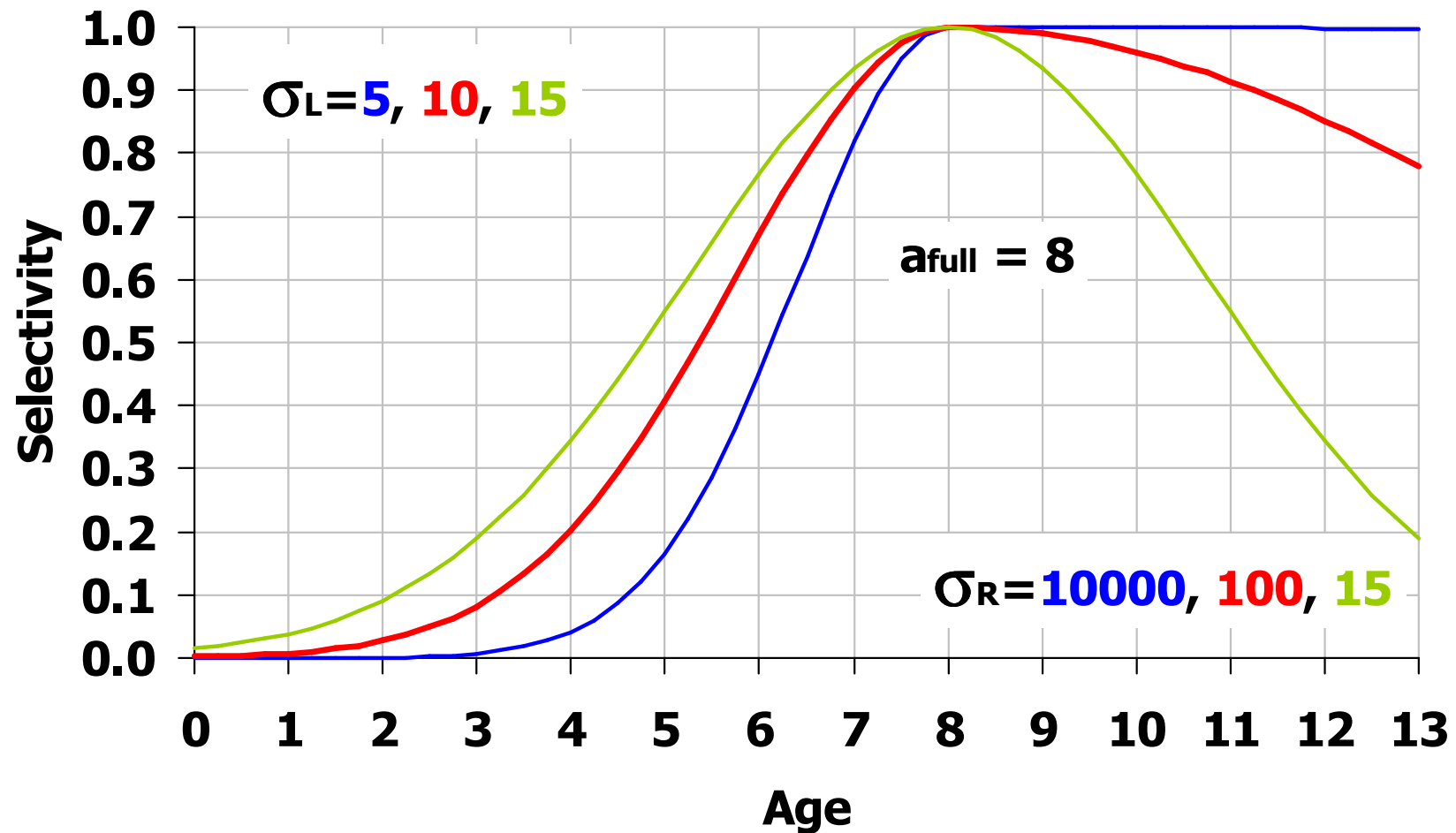
- Are simply two half-normal curves:

$$s_a = \begin{cases} e^{-\frac{(a-a_{full})^2}{\sigma_L}} & \text{for } a \leq a_{full} \\ e^{-\frac{(a-a_{full})^2}{\sigma_R}} & \text{for } a > a_{full} \end{cases}$$

- $a_{full}$ : age at full selectivity
- $\sigma_R$ : Shape factor (variance) for right hand curve
- $\sigma_L$ : Shape factor (variance) for left side of curve

# Selectivity - double half-Gaussian

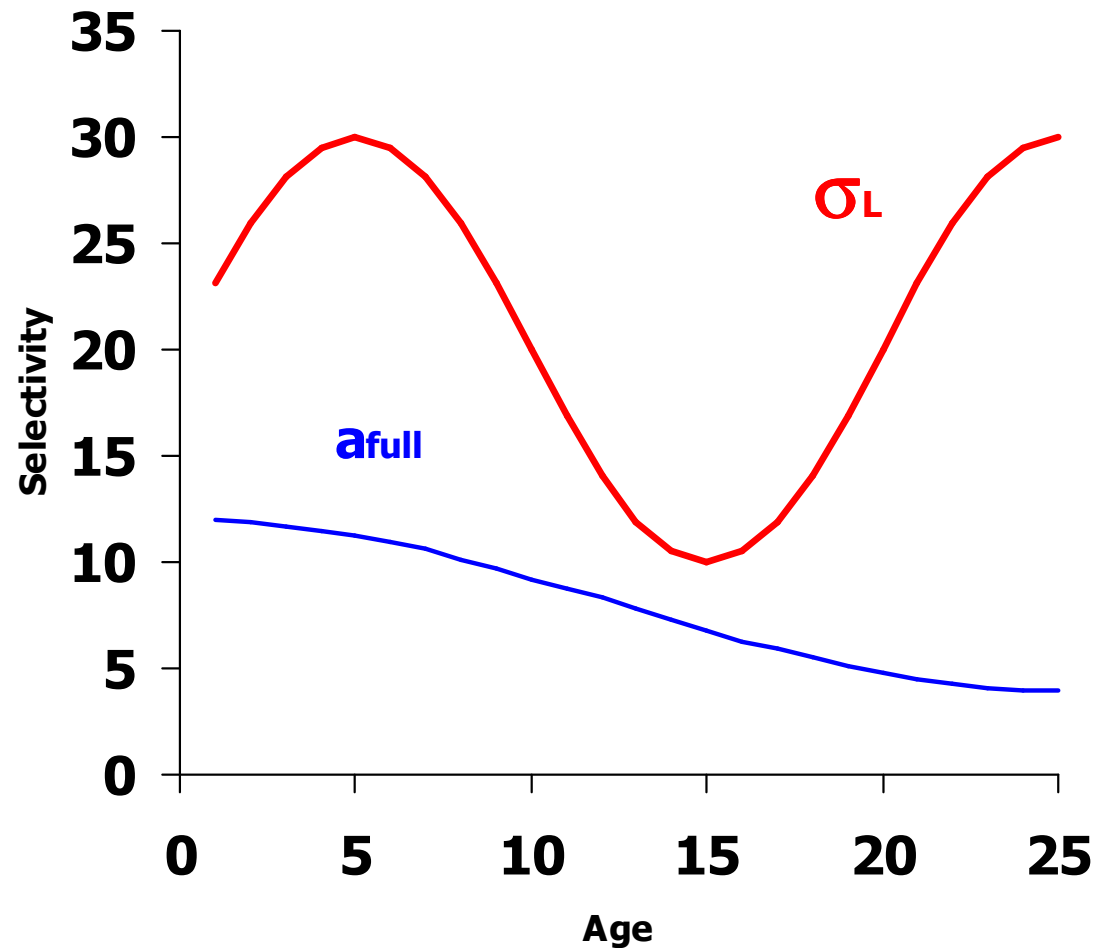
Note asymmetry



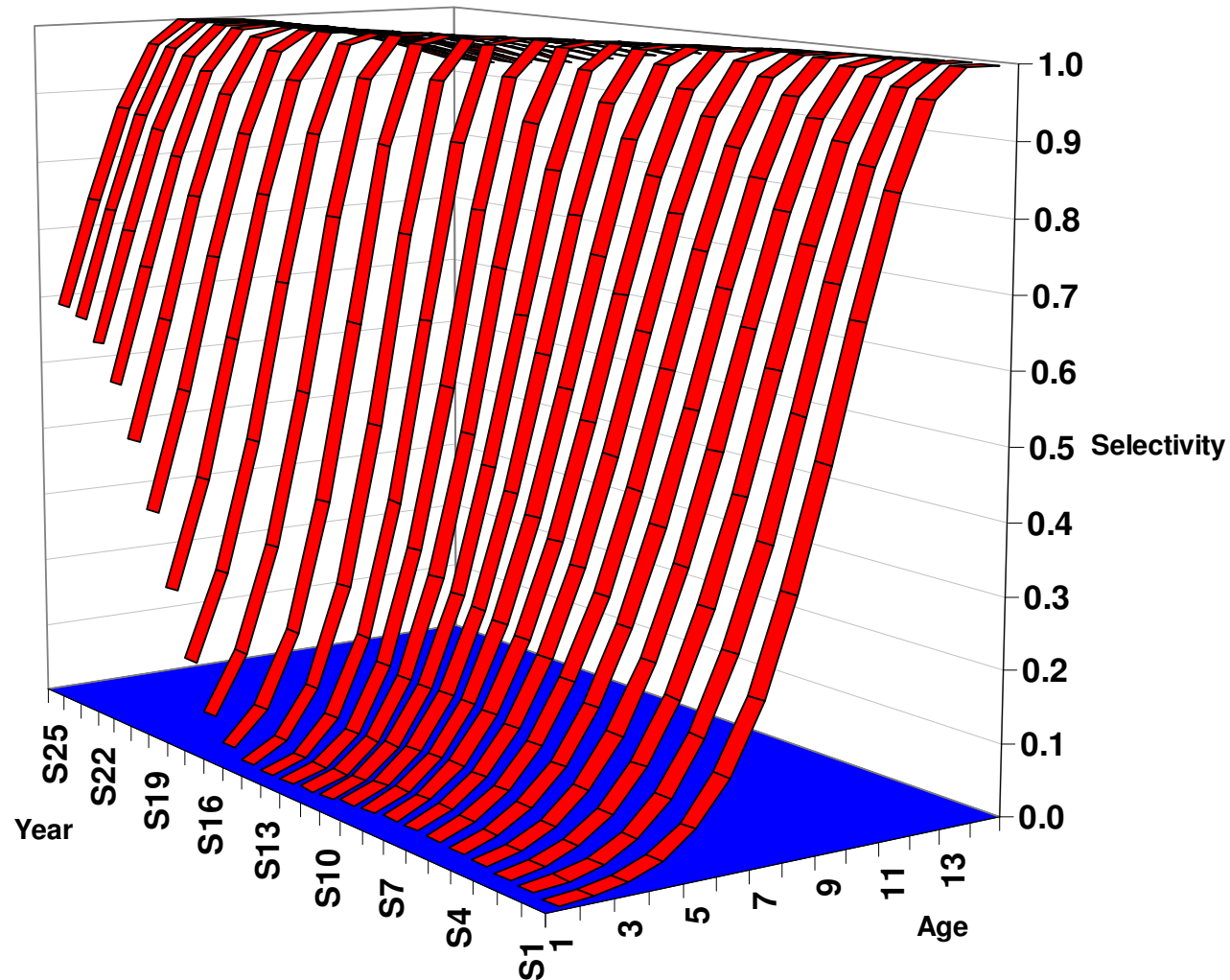
# Selectivity - adding periodic function

Double half Gaussian selectivity			
	Afull	$\sigma_L$	$\sigma_R$
Amplitude	0.5	0.5	0
Period	50	20	50
Level	8	20	1000
Tshift	0	25	25

Year	Afull	$\sigma_L$	$\sigma_R$
1	12.0	23.1	1000
2	11.9	25.9	1000
3	11.7	28.1	1000
4	11.5	29.5	1000
5	11.2	30.0	1000
6	10.9	29.5	1000
7	10.6	28.1	1000
8	10.1	25.9	1000
9	9.7	23.1	1000
10	9.2	20.0	1000
11	8.8	16.9	1000
12	8.3	14.1	1000
13	7.8	11.9	1000
14	7.3	10.5	1000
15	6.8	10.0	1000
16	6.3	10.5	1000
17	5.9	11.9	1000
18	5.5	14.1	1000
19	5.1	16.9	1000
20	4.8	20.0	1000
21	4.5	23.1	1000
22	4.3	25.9	1000
23	4.1	28.1	1000
24	4.0	29.5	1000
25	4.0	30.0	1000

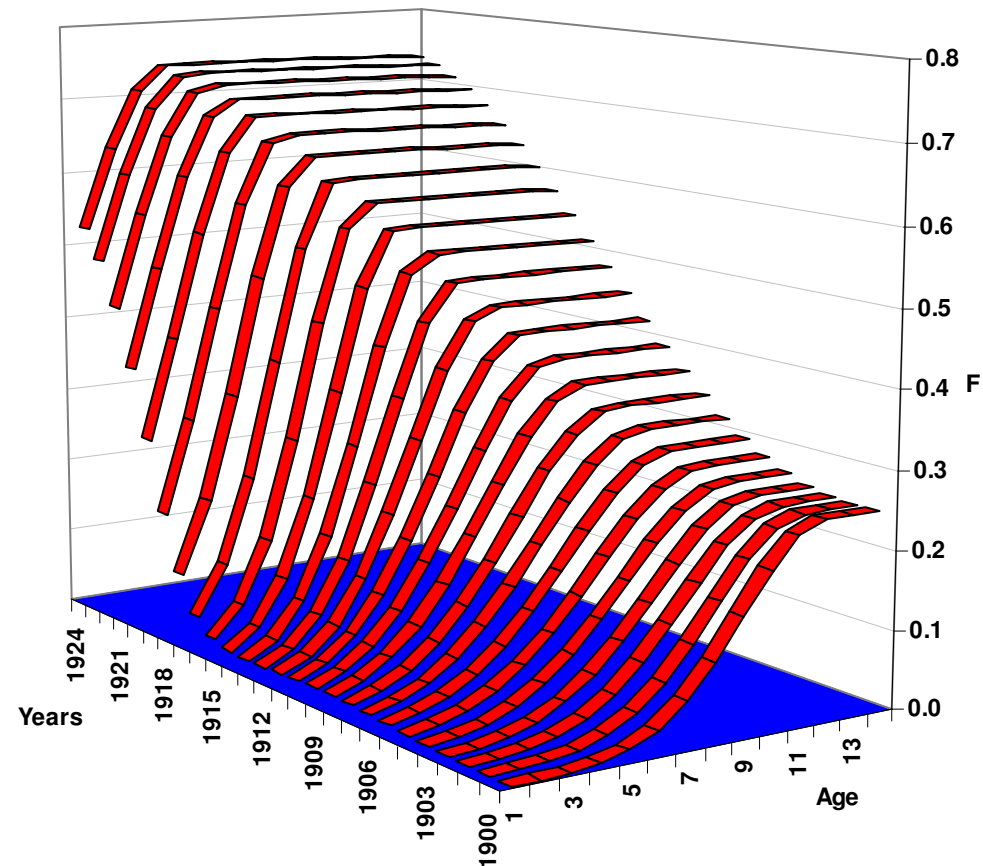


# Resulting selectivity pattern





$$F_{ay} = s_{ay} F_y$$

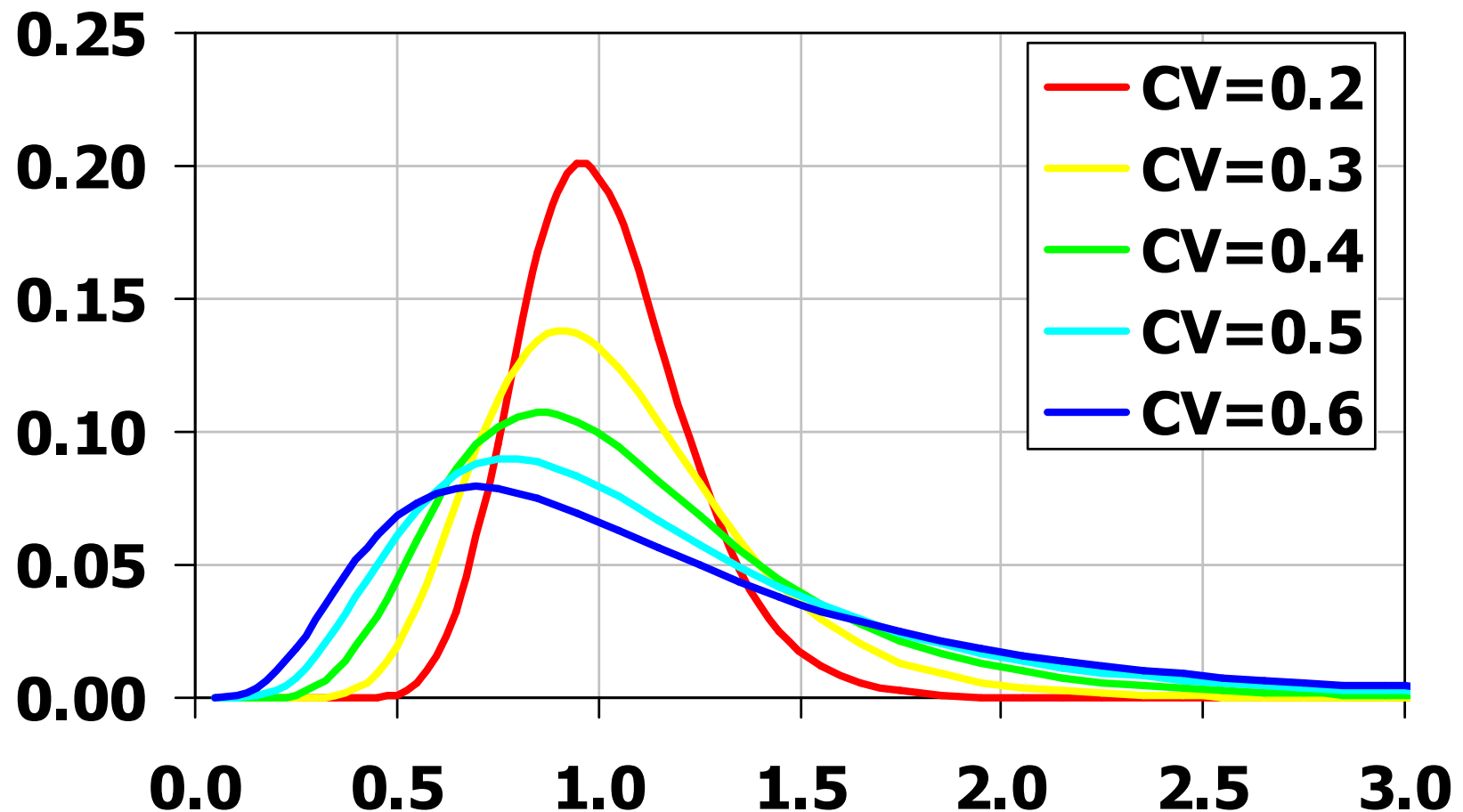


Pattern generated with relatively few parameters

# Adding stochastic / noise

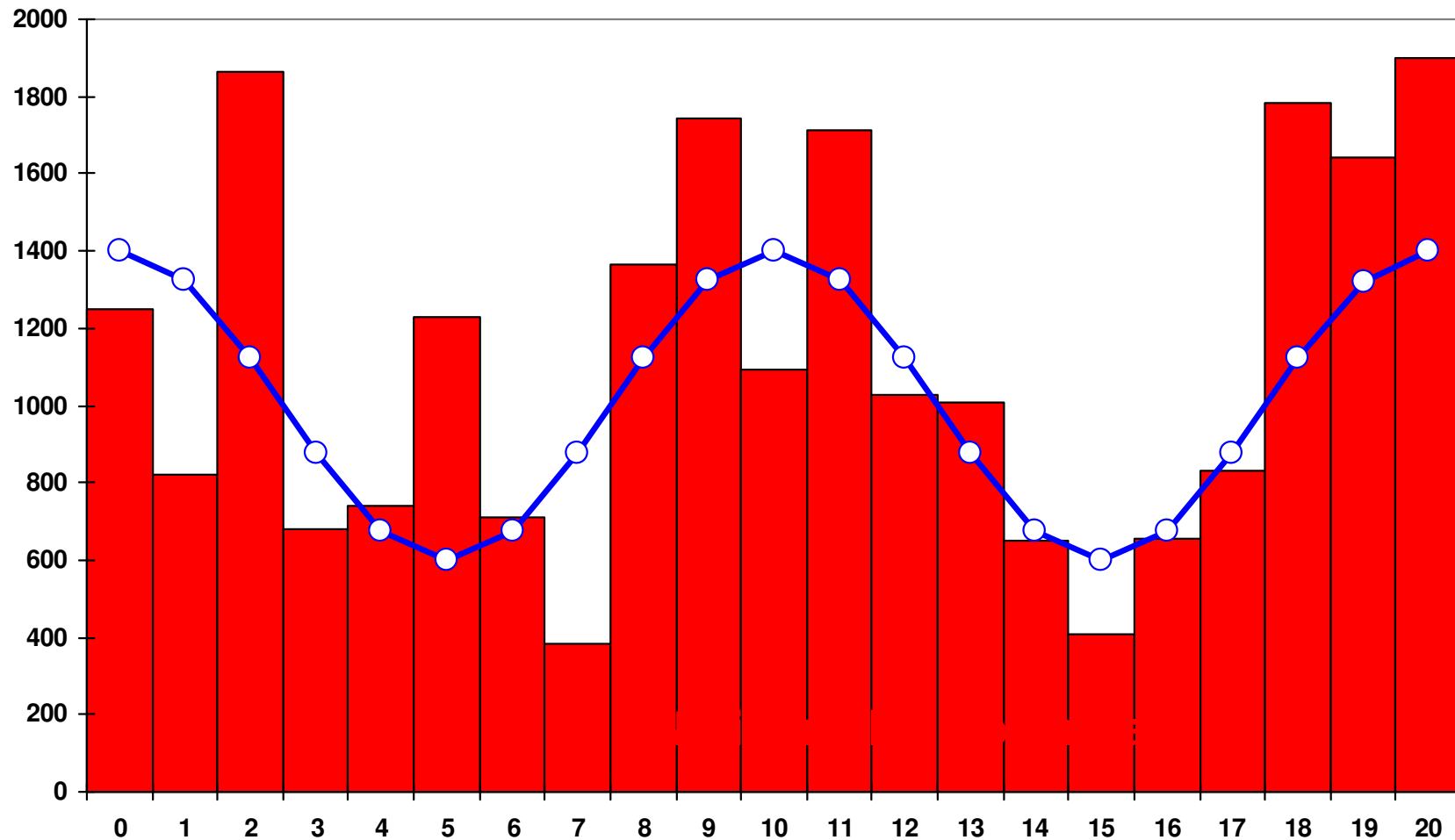
einar

XCEL SPEAK: = LOGINV(RAND();0;CV)



# Recruitment: Deterministic -> stochastic

XCEL SPEAK:  $=Ry*LOGINV(RAND();0;CV)$



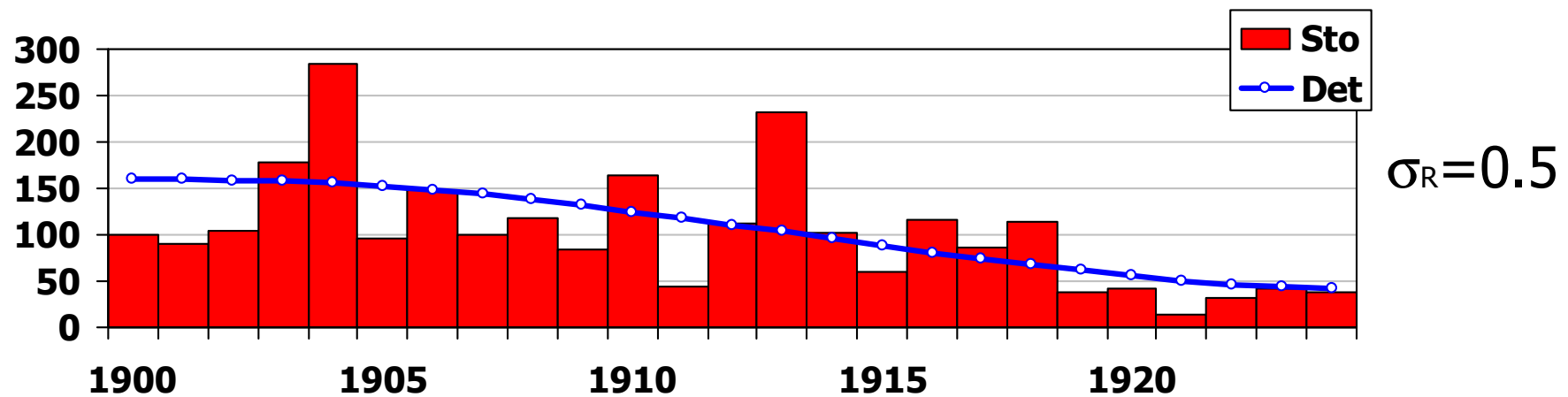
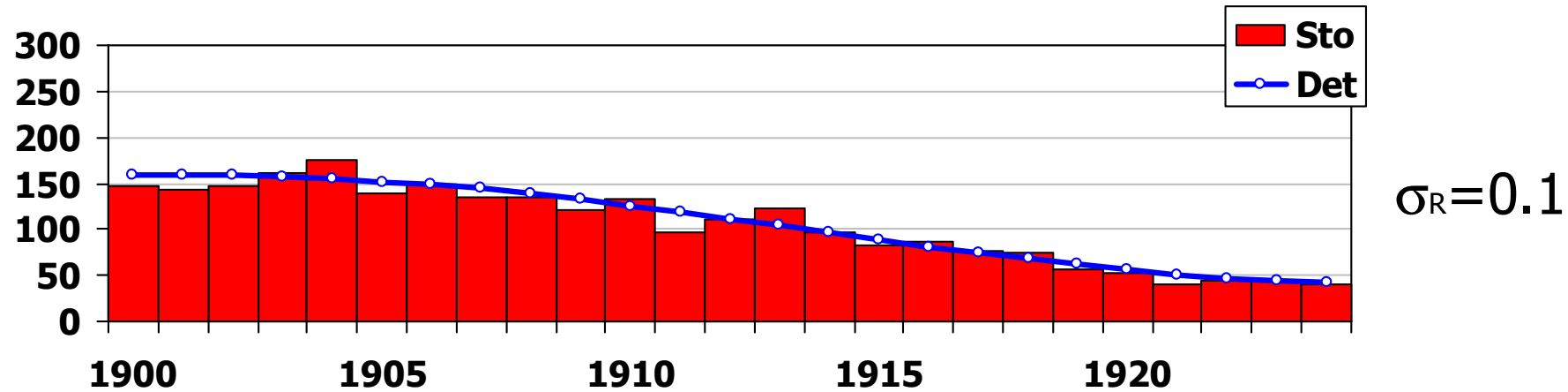
## Stated more formally

$$N_{A,y}^{DET} = \bar{N}_{A,y} \left\{ 1 + A * \cos \frac{y - T_{initY}}{P} * 6.28 \right\}$$

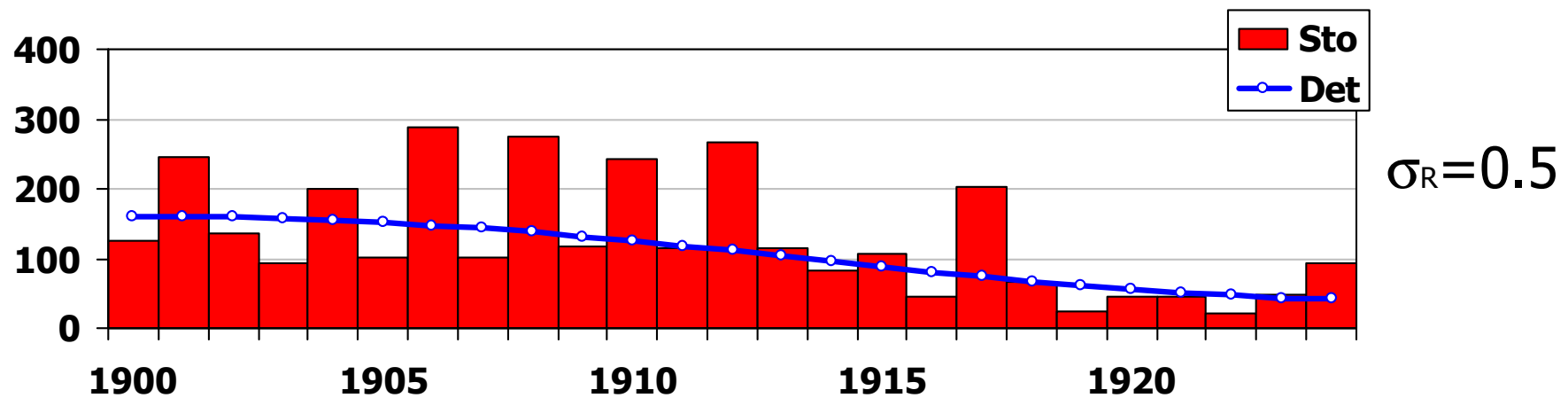
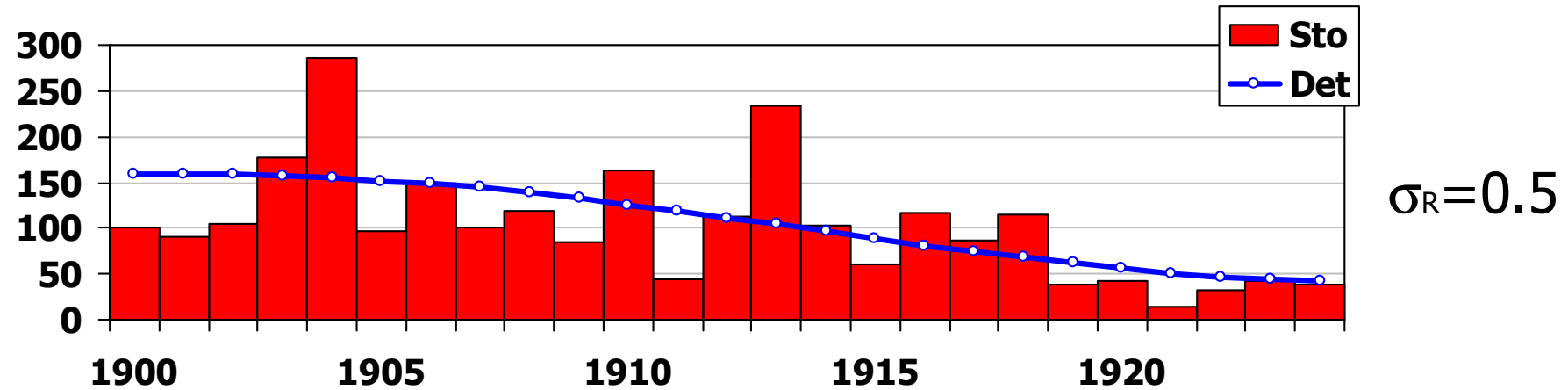
$$N_{A,y}^{STO} = N_{A,y}^{DET} e^{\varepsilon_R} \quad \varepsilon_R \sim N(0, \sigma_R^2)$$

- $\varepsilon_R$  normally distributed random number
  - Mean = 0
  - Standard deviation =  $\sigma_R$
- Normally distributed numbers set to the power of e gives lognormal distribution.
  - Note!  $\text{Loginv}(\text{rand}(), 0, s) = \exp(\text{norminv}(\text{rand}(), 0, s))$
- By changing the value of  $\sigma_R$  we can control the variance in recruitment.
- Often refer to as  $\sigma_R$  CVs.

# Recruitment: Different $\sigma_R$



# Recruitment: Same $\sigma_R$ , new random numbers



- Measurement error for  $C_{ay}$  are often assumed log normally distributed
- In mathematical terms this is described as:

**True catch**  
**Measured catch**

**Lognormal random numbers**

**Normally distributed random numbers with:  
 mean = 0  
 standard deviation =  $\sigma$**

$$C_{a,y}^{OBS} = C_{a,y}^{TRUE} e^{\varepsilon_{C_a}}$$

$$\varepsilon_{C_a} \sim N(0, \sigma_{C_a}^2)$$

- In Excel speak:

$$C_{a,y}^{OBS} = C_{a,y}^{TRUE} \times \text{LOGINV}(\text{RAND}(), 0, \sigma)$$

# Generating error and bias in catches

$$C_{ay}^{TRUE} = \frac{s_a F_y}{s_a F_y + M_{ay}} \left( 1 - e^{-(s_a F_y + M_{ay})} \right) N_{ay}$$

$$C_{ay}^{DET} = (1 - p_{ay}) C_{ay}^{TRUE}$$

$$C_{a,y}^{OBS1} = C_{a,y}^{DET} e^{\varepsilon_{C_a}} \quad \varepsilon_{C_a} \sim N(0, \sigma_{C_a}^2)$$

$$C_{a,y}^{OBS2} = C_{a,y}^{OBS1} \frac{Y_y^{DET}}{Y_y^{OBS1}}$$

$$Y_y^{DET} = \sum_a C_{a,y}^{TRUE} w_{a,y}^{catch}$$

$$Y_y^{OBS1} = \sum_a C_{a,y}^{OBS1} w_{a,y}^{catch}$$

•The actual catch taken

•The actual catch

•Measured catch, lognormal error

•Scaling to total yield

•Formula for total yield

**p<sub>ay</sub>: proportion of unrecorded catch**



- Simple model
- More complex
  - $q_a$ : catchability for age  $a$
  - $\beta_a$ : power parameter for age  $a$ 
    - NB, if  $\beta=1$  it is a simple model
- Simulating measurement errors

$$U_{ay} = q_a N_{a,y}$$

$$U_{ay} = q_a N_{a,y}^{\beta_a}$$

$$U_{ay}^{OBS} = U_{ay}^{TRUE} e^{\varepsilon_{U_a}} \quad \varepsilon_{U_a} \sim N(0, \sigma_{U_a})$$

- In the simulator we can also add year effect and  $q$  –trends.
  - Lets not worry to much about that for now

# The simulator in math (almost)

$$N_{a,y} = \begin{cases} R_y = N_{A,y}^{STO} = N_{A,y}^{DET} e^{\varepsilon_R} & \varepsilon \sim N(0, \sigma_R^2) \\ N_{a-1,y-1} e^{-(s_{a-1}F_{y-1} + M_{a-1,y-1})} \\ N_{a-1,y-1} e^{-(s_{a-1}F_{y-1} + M_{a-1,y-1})} + N_{a,y-1} e^{-(s_a F_{y-1} + M_{a,y-1})} \end{cases} \quad \begin{cases} a = A \\ A < a \leq a_{plus} \\ a = a_{plus} \end{cases}$$

$$F_{ay} = s_a F_y$$

$$C_{ay}^{TRUE} = \frac{s_a F_y}{s_a F_y + M_{ay}} \left( 1 - e^{-(s_a F_y + M_{ay})} \right) N_{ay}$$

$$C_{ay}^{DET} = (1 - p_{ay}) C_{ay}^{TRUE}$$

$$C_{a,y}^{OBS1} = C_{a,y}^{DET} e^{\varepsilon_{C_a}} \quad \varepsilon_{C_a} \sim N(0, \sigma_{C_a}^2)$$

$$C_{a,y}^{OBS2} = C_{a,y}^{OBS1} \frac{Y_y^{OBS2}}{Y_y^{OBS1}}$$

$$Y_y = \sum_a C_{a,y} w_{a,y}^{catch}$$

$$R_y = f(time)$$

$$F_y = f(time)$$

$$s_{ay} = f(age, time)$$

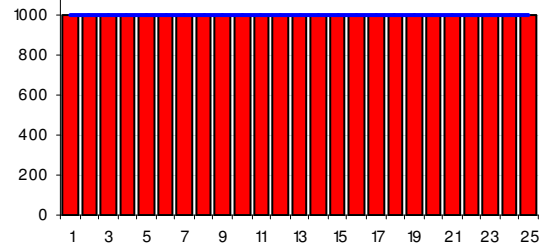
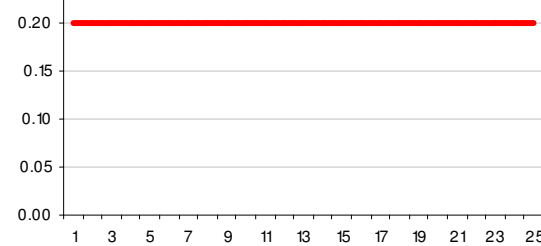
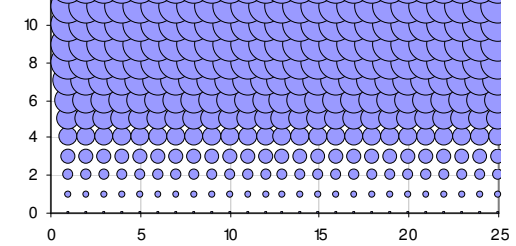
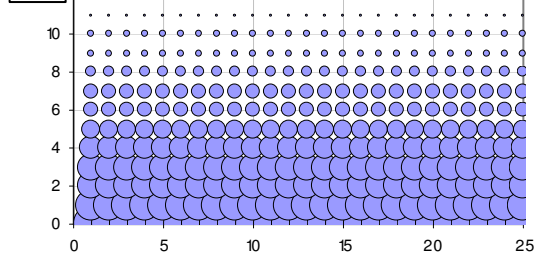
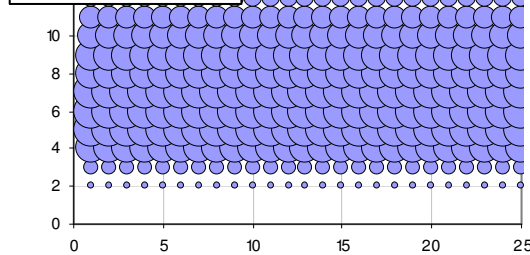
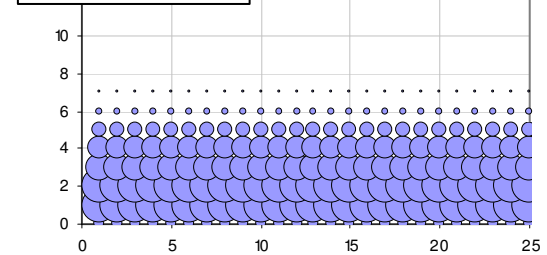
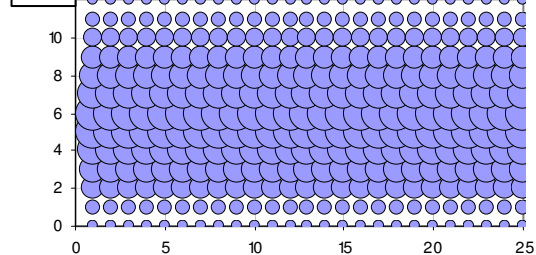
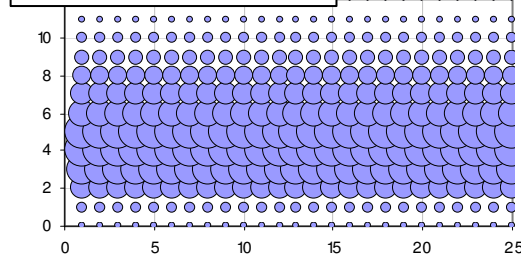
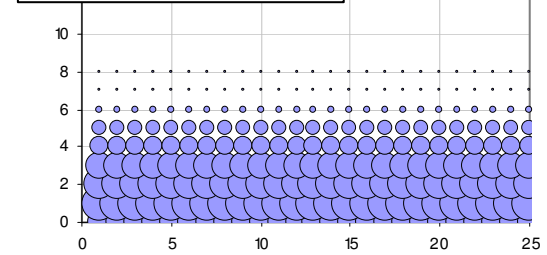
$$q_a = f(age, time)$$

$$M_{ay} = f(?)$$

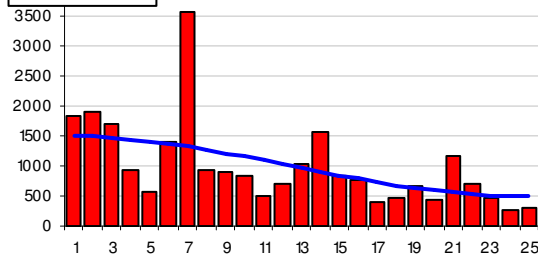
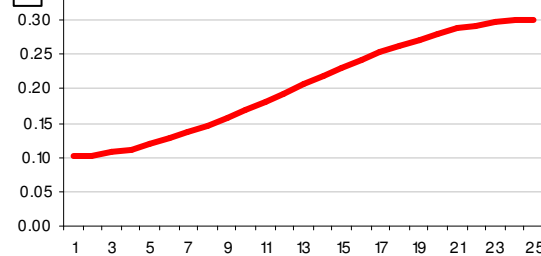
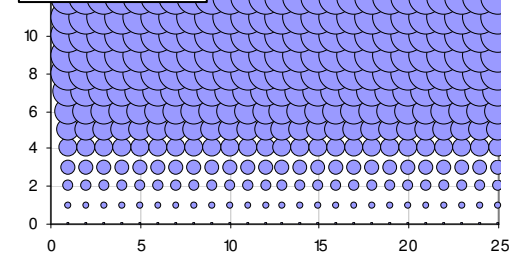
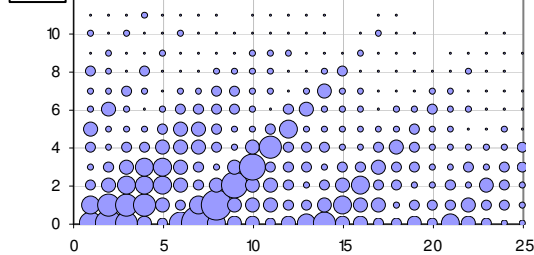
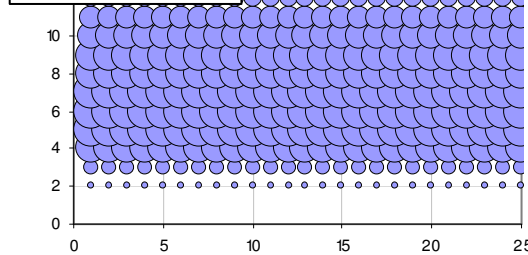
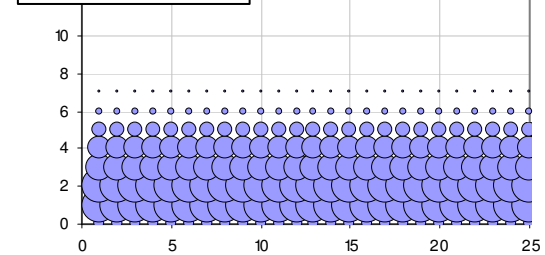
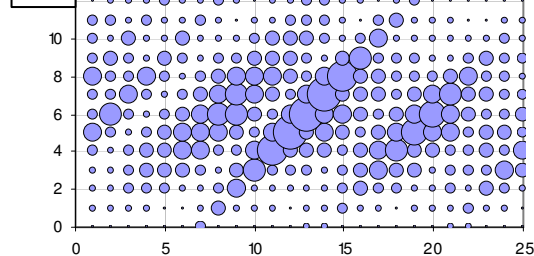
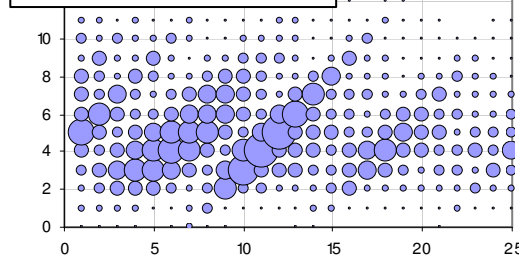
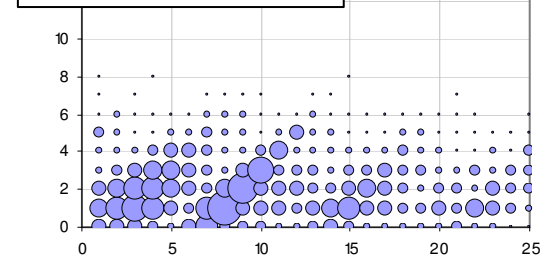
$$U_{ay}^{TRUE} = q_a N_{ay}^{\beta_a}$$

$$U_{ay}^{OBS} = U_{ay}^{TRUE} e^{\varepsilon_{U_a}} \quad \varepsilon_{U_a} \sim N(0, \sigma_{U_a})$$

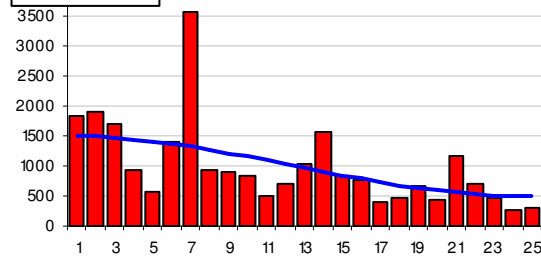
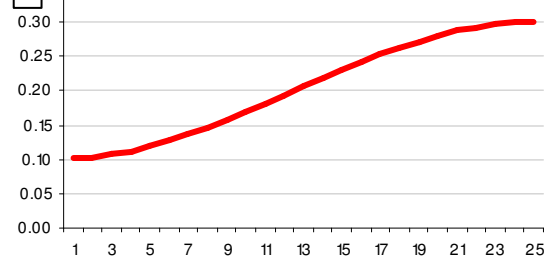
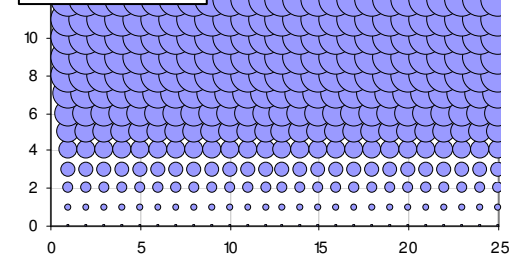
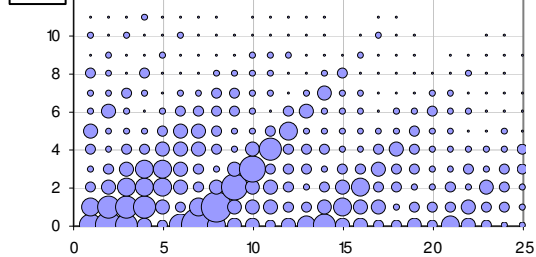
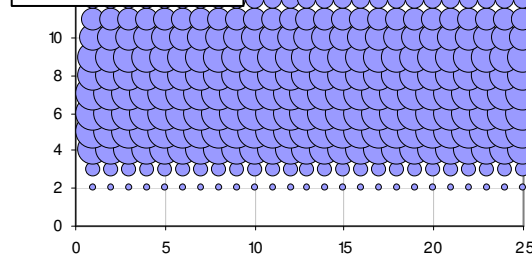
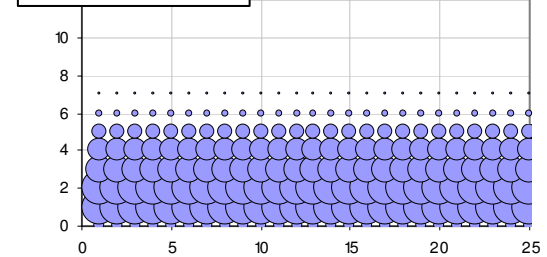
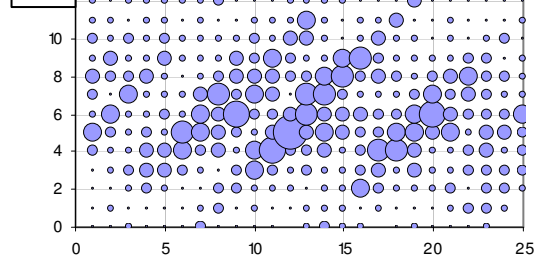
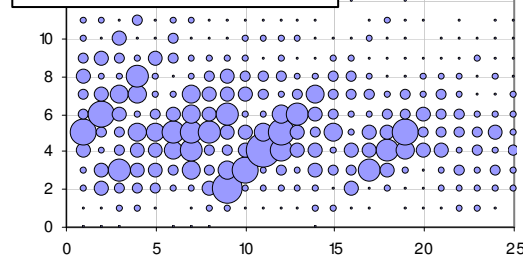
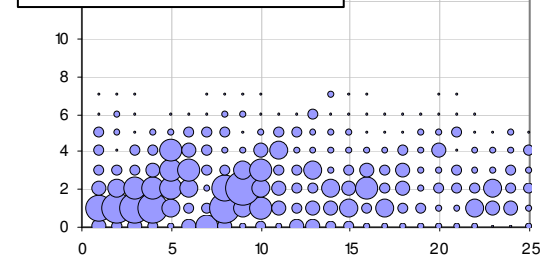
# Example of simulations: Static system

**Recruitment****F****Catch selectivity****Nay****Survey 1 catchability****Survey 2 catchability****c@a****Survey 1, observed numbers****Survey 2, observed numbers**

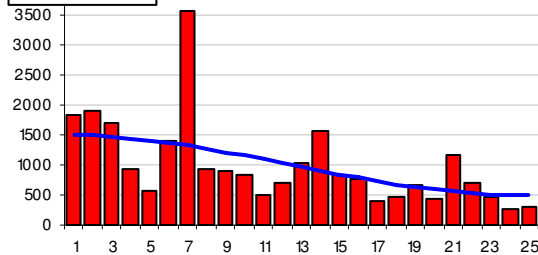
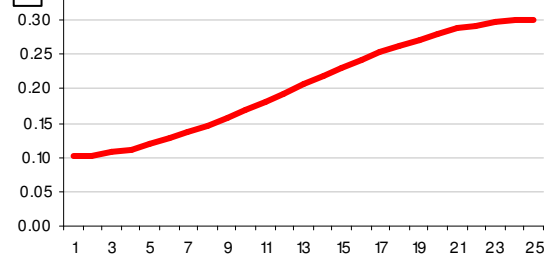
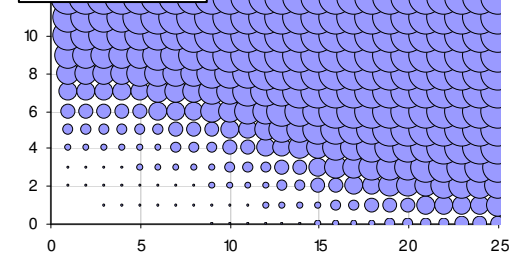
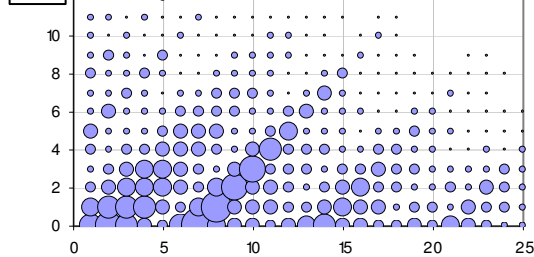
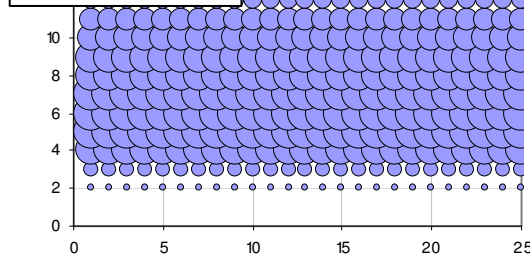
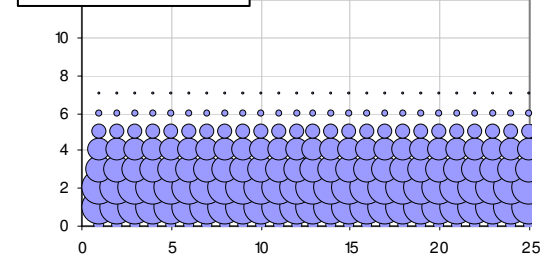
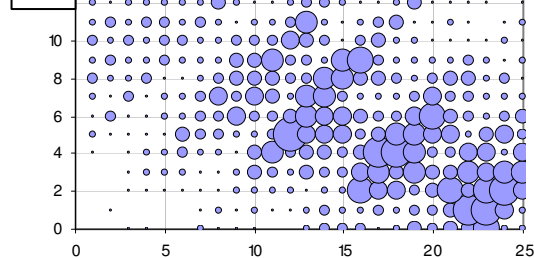
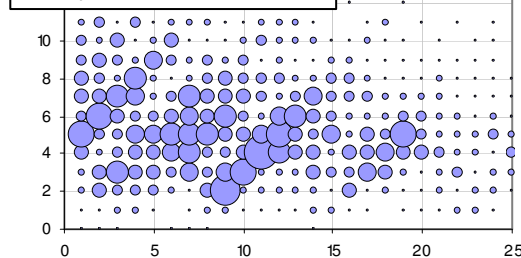
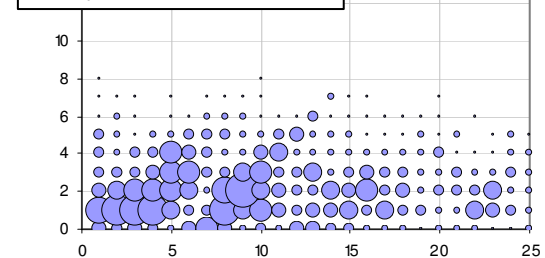
# Adding stochasticity & time trends

**Recruitment****F****Catch selectivity****Nay****Survey 1 catchability****Survey 2 catchability****c@a****Survey 1, observed numbers****Survey 2, observed numbers**

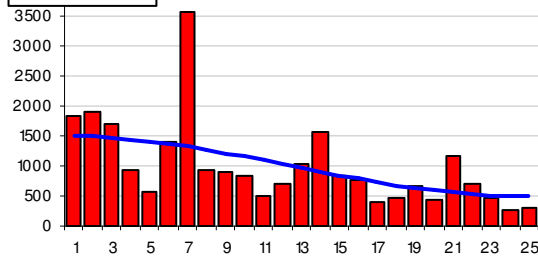
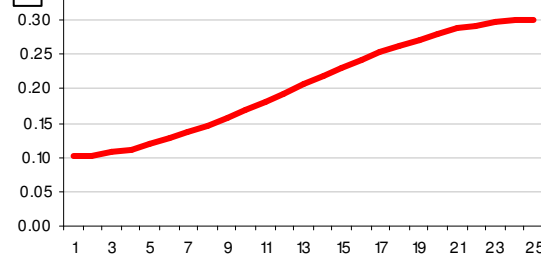
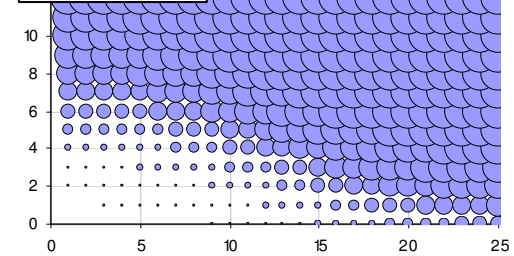
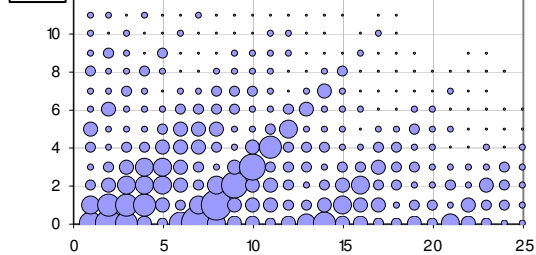
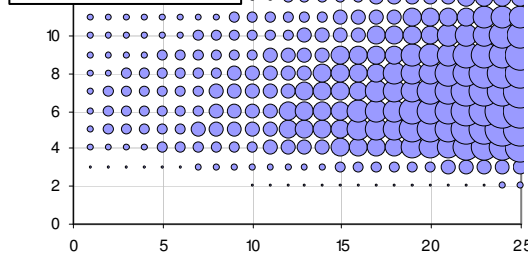
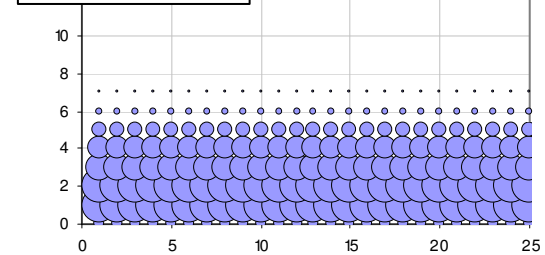
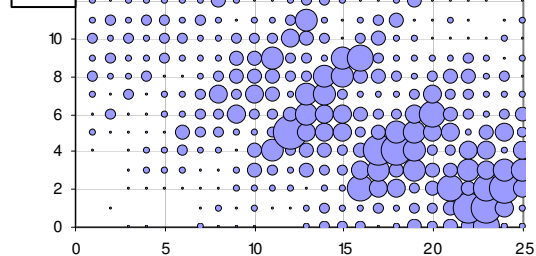
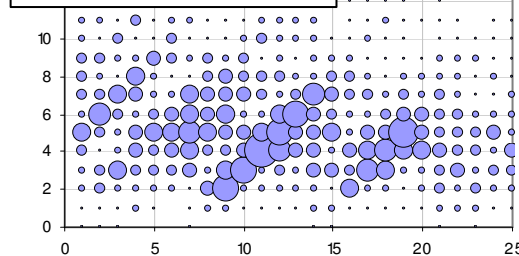
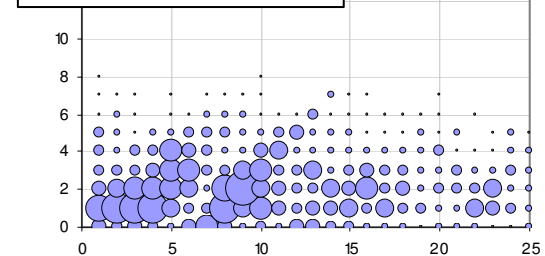
# Adding observation noise

**Recruitment****F****Catch selectivity****Nay****Survey 1 catchability****Survey 2 catchability****c@a****Survey 1, observed numbers****Survey 2, observed numbers**

# Adding change in fleet selectivity

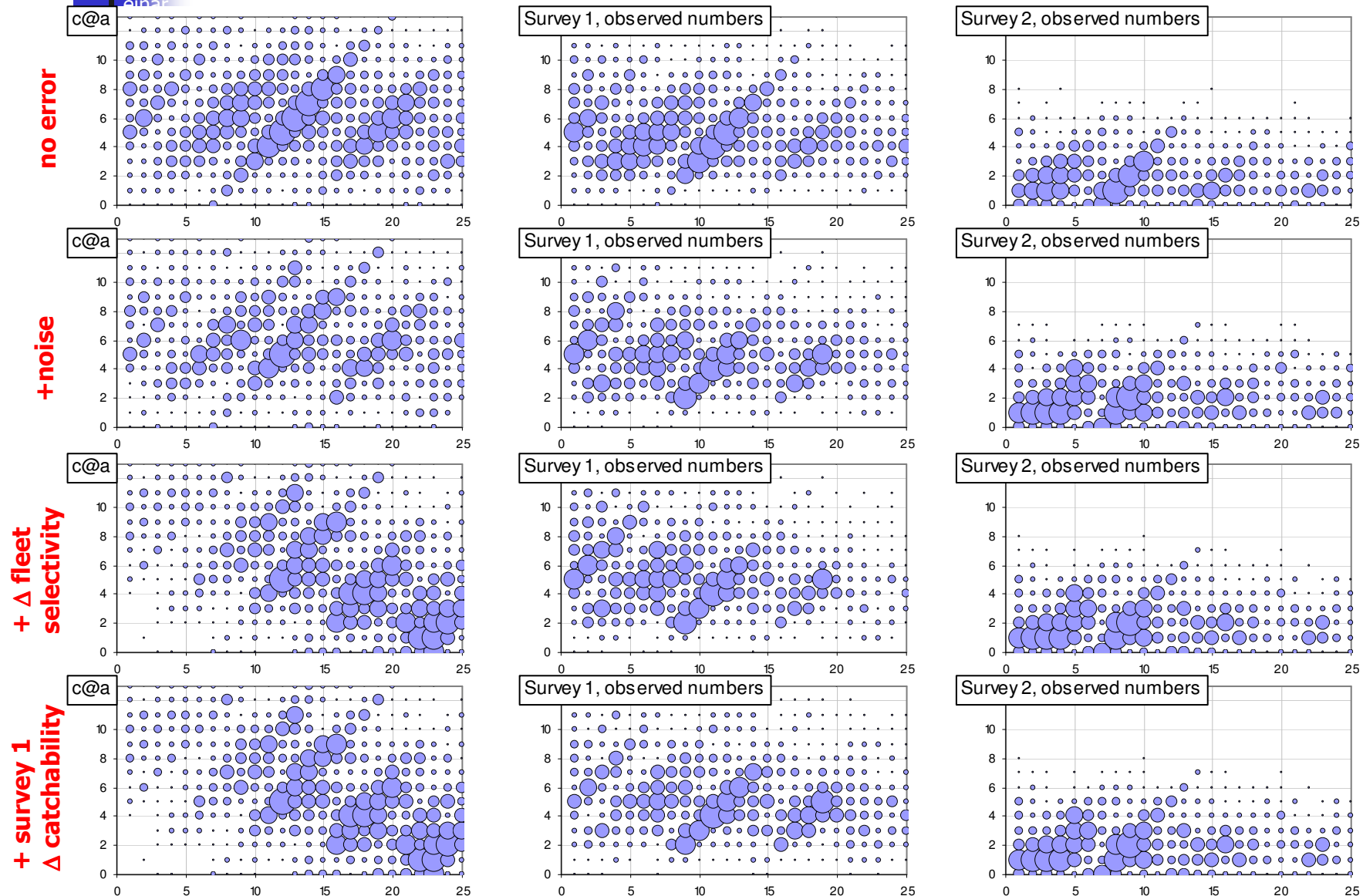
**Recruitment****F****Catch selectivity****Nay****Survey 1 catchability****Survey 2 catchability****c@a****Survey 1, observed numbers****Survey 2, observed numbers**

# Adding trend in survey 1 (q increases)

**Recruitment****F****Catch selectivity****Nay****Survey 1 catchability****Survey 2 catchability****c@a****Survey 1, observed numbers****Survey 2, observed numbers**



# Same R input - different observation







A summary of things done

- Population, fisheries and observation simulation
  - **Operational model**
    - Introduced the concept and the mathematics for simulating population dynamics, fisheries and typical observations in “a data rich” situations.
    - Used a spreadsheet (xGenerator.xls) to get a hands on feel with the intent to improve the understanding of the concept and mathematics in the simulator.
    - Generated observables with given measurement errors:
      - One  $c@a$  matrix – catch at age by years in numbers (Cay)
      - Two  $u@a$  matrices – age based survey indices (Uay)
      - Note: These simulated measurements could be used in any conventional off-the shelf assessment software package, be it XSA, ADAPT, TSA, ADCAM, AMCI, ICA, ...

- Set-up a simple separable **assessment model** (xModel.xls)
  - Use observation ( $C_{ay}$ ,  $U_{ay}$ ), generated with the simulator to estimated stock size ( $N_{ay}$ ) and fishing mortalities ( $F_{ay}$ ), selection pattern, and survey catchability.
    - Note: Could have used actual observations ( $C_{ay}$ ,  $U_{ay}$ ) from any of the stocks where we have  $c@a$  and age based survey indices, e.g. Icelandic cod, Southern hake, ...
    - In our case we are however more interested in understanding how assessment software works, where it goes wrong and how it goes wrong. For that purpose we use the observation generated in the simulator.

- Provided a spreadsheet (xRealityCheck.xls) that made a comparison on some key feature (recruitment, SSB, selection pattern, fishing mortalities) as simulated in the generator ("the truth") with that estimated from the Cay and Uay observation in the assessment model.

# The real world

# What have we done – in short

The Generator is a simplification of the real world and should be taken as such

The BIG question

## Our education world

## Assessment world (and most education worlds)

Measurement/observations

- Landings
- Catch at age
- Survey at age

Output

### xGenerator.xls

Input:

- Recruitment
- Fishing pattern
- Fishing mortality
- Natural mortality
- Survey catchability

Calculated:

- Catch at age
- Survey at age
- Added some noise and biases**

### xRealityCheck.xls

Compared  
known truth  
with estimates

Hopefully learned  
something

Input

Input:

### xModel.xls

- Catch at age
- Survey at age

Estimated:

- Recruitment
- Fishing pattern
- Fishing mortality
- Survey catchability

Assumed:

- M constant
- Landings = catch
- Constant fishing pattern

Diagnostics

- Internal consistency

Uncertainty

# Just to make sure we have this right

- As it currently is set up the only information that is passed from the xGenerator to the xModel are the observed catch at age and observed survey at age, including noise. I.e. any other settings are NOT passed from the xGenerator to the xModel.
  - The xModel thus does not know anything about recruitment, selection patterns, mortalities (including unaccounted ones) or survey catchabilities that were used to generate the observations in the xGenerator.xls
- In our daily world we are effectively limited to the “Assessment world”.
- The generator is a simplified, and possibly a wrong reflection of the real world, the assessment ever more so.

# Assignment work

- You will be given a set of exercises, which you are expected to keep notes on and make a report on, to be evaluated as your assignment.
- How to proceed:
  - Keep a backup of the original sets of Excel spreadsheets
  - Make a separate folder for each exercise
    - xGenerator.xls
    - xModel.xls
    - xRealityCheck.xls
- Keep good notes from classroom discussion, they will help in writing the report

# Exercise 1: The effect of discards

- One of the biggest problems in fisheries science is that one normally has only estimates of actual landings but the true removal is unknown. Discards of a target species is normally assumed to be size related, the rate of discards normally being confined to the smallest fishes.
- Part 1. Start by fitting the model (run solver) without discards.
  - Discarding is controlled in the simulator, worksheet CatchAtAge in area starting in cell B90.
  - For this part you want to make sure that all the values are = 0.0
- Part 2. Set fix discard rate on younger ages through the whole time period.
  - Discarding is controlled in the simulator, worksheet CatchAtAge in area starting in cell B90.
  - Make sure that there is actually some fishing taking place on these younger ages!
  - Suggestion for discard rate:
    - Age 1 & 2: 0.99 – there is no fishing taking place to speak of in these age groups, see the selection pattern (worksheet xx)
    - Age 3: 0.75
    - Age 4: 0.50
    - Age 5: 0.25
    - Age 6: 0.10
- Can one detect any abnormal (nonrandom) patterns in residuals in either fit? Where are the major differences in the parameter estimates of the two fits?
- Compare the fit in Part 2 with that of your simulated population. Which of the key population measures (recruitment, biomass (SSB) and fishing mortality) deviate mostly from the known truth? Explain why this pattern is observed?
- Part 3. Feel free to try any other combination on discard rates.



- Tomorrow we will work further with the tools provided. We will study how:
  - Changing discard rate with time
  - Underreporting
  - Change in the catchability of the survey
  - Values of  $M$
  - , and other factors
- affect our assessment results.
- The brave ones are encourage to proceed right away.