

Advanced Bayesian modeling

Joachim Vandekerckhove

Order of operations

$$X = \begin{bmatrix} 2 & 1 & 2 \\ 1 & 1 & 5 \end{bmatrix}$$

Order of operations

$$X = \begin{bmatrix} 2 & 1 & 2 \\ 1 & 1 & 5 \end{bmatrix}$$

Sum of sums:

$$\sum_{r=1}^2 \left(\sum_{c=1}^3 x_{rc} \right) = \sum_{c=1}^3 \left(\sum_{r=1}^2 x_{rc} \right)$$

Order of operations

$$X = \begin{bmatrix} 2 & 1 & 2 \\ 1 & 1 & 5 \end{bmatrix}$$

Sum of sums:

$$\sum_{r=1}^2 \left(\sum_{c=1}^3 x_{rc} \right) = \sum_{c=1}^3 \left(\sum_{r=1}^2 x_{rc} \right)$$

Sum of products vs. product of sums:

$$\sum_{r=1}^2 \left(\prod_{c=1}^3 x_{rc} \right) \neq \prod_{c=1}^3 \left(\sum_{r=1}^2 x_{rc} \right)$$

Order of operations

In general, order of operations matters:

$$f \circ g(x) \neq g \circ f(x)$$

Estimating parameters is an operation

In science, we often want to make statements about averages.

Estimating parameters is an operation

In science, we often want to make statements about averages.

Estimating model parameters from data is an operation:

$$\hat{\theta} = f(x)$$

Estimating parameters is an operation

In science, we often want to make statements about averages.

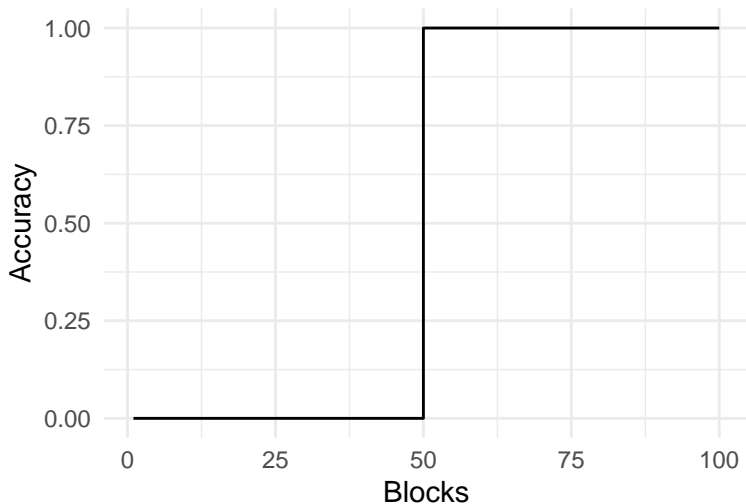
Estimating model parameters from data is an operation:

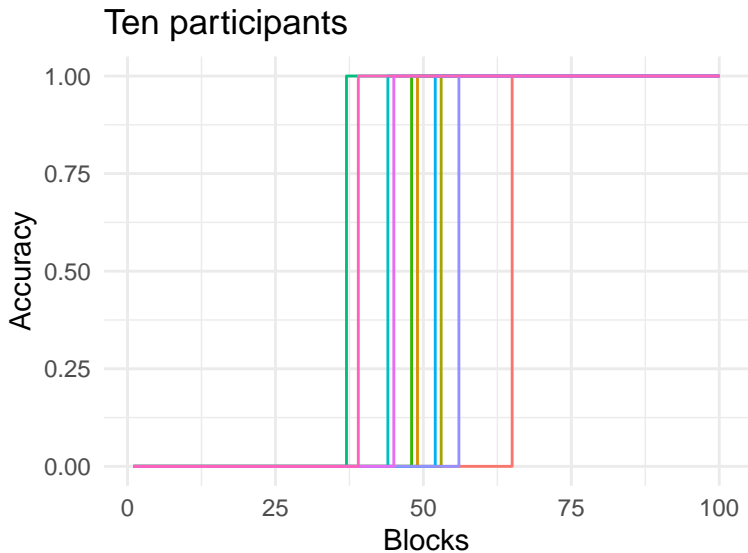
$$\hat{\theta} = f(x)$$

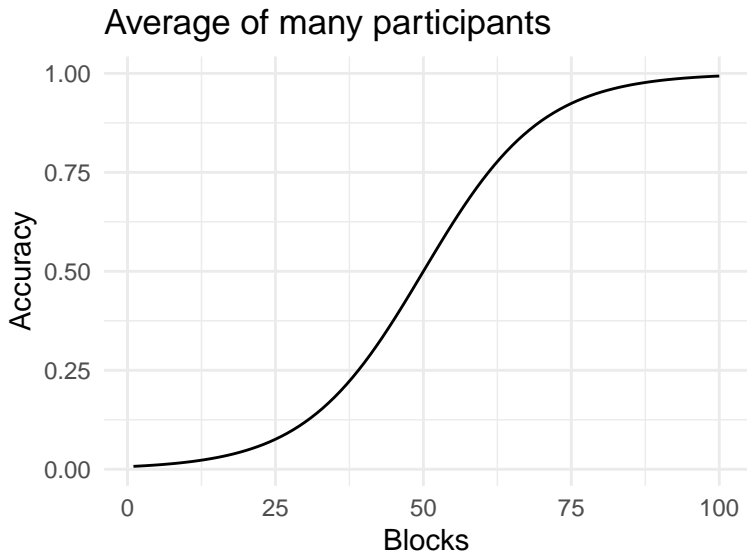
Do we want the **average model parameters of the data** or the **model parameters of the average data**?

$$\overline{f(x)} \neq f(\bar{x})$$

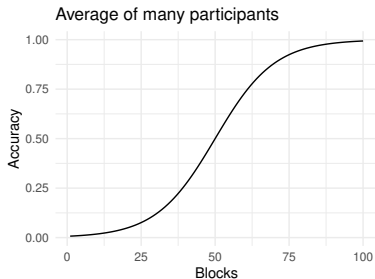
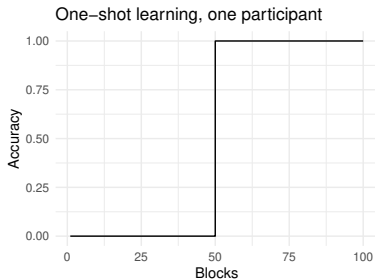
One-shot learning, one participant







The “average” learning curve looks nothing like the person-specific learning curve!



Instead, what we want here is to acknowledge that each person has their own trajectory, and then say something about the (average) properties of the trajectories.

Instead, what we want here is to acknowledge that each person has their own trajectory, and then say something about the (average) properties of the trajectories.

We want to make an abstraction of the data, which are something complicated that is generated by a process with parameters θ_p (for person p), and instead focus on parameters.

$$\mathcal{M}_h : \begin{cases} x_p \sim \text{one-shot}(\theta_p) \\ \theta_p \sim N(\mu, \tau) \end{cases}$$

The hierarchical model contains one set of assumptions about the data given the model (the likelihood level), and another set of assumptions about structure among parameters.

$$\mathcal{M}_h : \begin{cases} x_p \sim \text{one-shot}(\theta_p) \\ \theta_p \sim N(\mu, \tau) \end{cases}$$

The hierarchical model contains one set of assumptions about the data given the model (the likelihood level), and another set of assumptions about structure among parameters.

Here, the hierarchical parameters μ and τ tell us something about the population of participants, each with their own change point θ_p .

Population assumptions

An assumption made in hierarchical models is that it is possible to know things about participants merely because they are members of some sampled population.

Population assumptions

An assumption made in hierarchical models is that it is possible to know things about participants merely because they are members of some sampled population.

For example, just by knowing someone is a human, we can make reasonable estimates about their number of thumbs.

Population assumptions

An assumption made in hierarchical models is that it is possible to know things about participants merely because they are members of some sampled population.

For example, just by knowing someone is a human, we can make reasonable estimates about their number of thumbs.

If $n_t > 2$, a measurement error may have occurred.

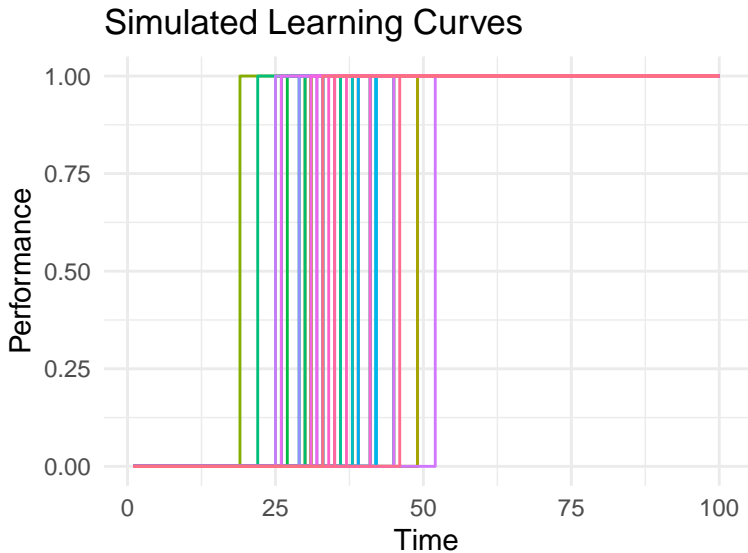
Back to the one-shot learning, we could try to characterize the population of participants in terms of their change point distribution $\theta_p \sim N(\mu, \tau)$.

Back to the one-shot learning, we could try to characterize the population of participants in terms of their change point distribution $\theta_p \sim N(\mu, \tau)$.

- Let's generate learning curves for a population of $n = 50$ participants, each with a random change point.

Back to the one-shot learning, we could try to characterize the population of participants in terms of their change point distribution $\theta_p \sim N(\mu, \tau)$.

- Let's generate learning curves for a population of $n = 50$ participants, each with a random change point.
- Let's use a step function to represent the learning curves.



- We define a hierarchical model in JAGS to estimate each participant's change point.

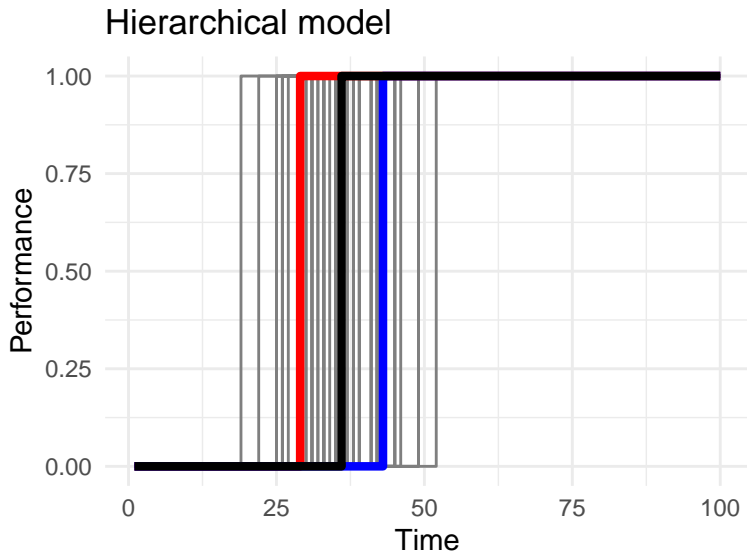
Hierarchical model of insight learning

- We define a hierarchical model in JAGS to estimate each participant's change point.
- The change point for each participant is modeled as a uniform distribution between 1 and 100.

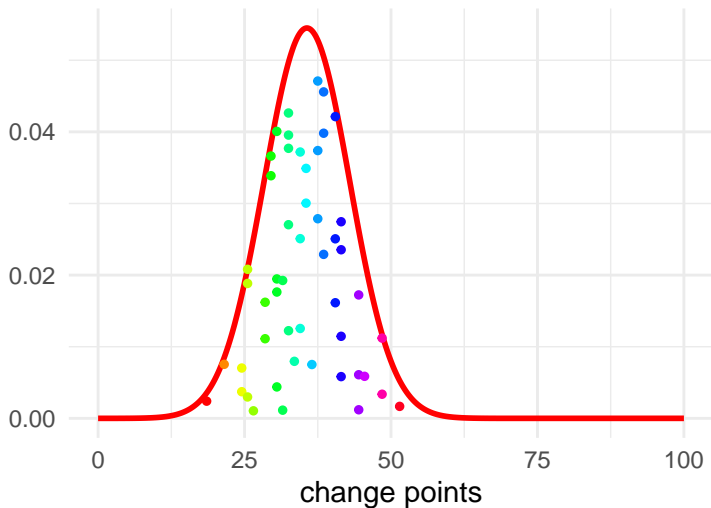
Hierarchical model of insight learning

```
##  
## model {  
##   for (j in 1:P) {  
##     for (i in 1:N) {  
##       y[i, j] ~ dbern(p[i, j])  
##       p[i, j] <- ifelse(time[i] < theta[j], 0, 1)  
##     }  
##     theta[j] ~ dnorm(mu, tau)T(0,100)  
##   }  
##   mu ~ dnorm(50, 0.05)T(0,100)  
##   tau ~ dnorm(20, 0.10)T(0,)  
##   sigma <- pow(tau, -0.5)  
## }
```

Hierarchical model of insight learning



Population distribution of change points



Parameter estimates

```
summary(samples_hrcl)[1]$  
  statistics[c("mu","sigma"),  
             c("Mean", "SD", "Time-series SE")] %>%  
  print()
```

##		Mean	SD	Time-series SE
## mu	35.642928	1.0245142	0.011413396	
## sigma	7.320065	0.7517616	0.008193004	

Parameter estimates

```
summary(samples_hrc1)[1]$  
  statistics[c("mu","sigma"),  
             c("Mean", "SD", "Time-series SE")] %>%  
  print()
```

##		Mean	SD	Time-series SE
## mu	35.642928	1.0245142	0.011413396	
## sigma	7.320065	0.7517616	0.008193004	

Compare to:

- Simulated $\mu = 35$

Parameter estimates

```
summary(samples_hrc1)[1]$  
  statistics[c("mu","sigma"),  
             c("Mean", "SD", "Time-series SE")] %>%  
  print()
```

##		Mean	SD	Time-series SE
## mu	35.642928	1.0245142	0.011413396	
## sigma	7.320065	0.7517616	0.008193004	

Compare to:

- Simulated $\mu = 35$
- Simulated $\sigma = 8$

Unorthodox things you can do if you're a Bayesian

Unorthodox things you can do if you're a Bayesian

"If you're a Bayesian you can do everything that God forbids."
(- Willem Heiser)

Unorthodox things you can do if you're a Bayesian

"If you're a Bayesian you can do everything that God forbids."

(- Willem Heiser)

- Accumulate evidence for the absence of an effect

Unorthodox things you can do if you're a Bayesian

"If you're a Bayesian you can do everything that God forbids."

(- Willem Heiser)

- Accumulate evidence for the absence of an effect
- Draw valid conclusions with almost no data

Unorthodox things you can do if you're a Bayesian

"If you're a Bayesian you can do everything that God forbids."

(- Willem Heiser)

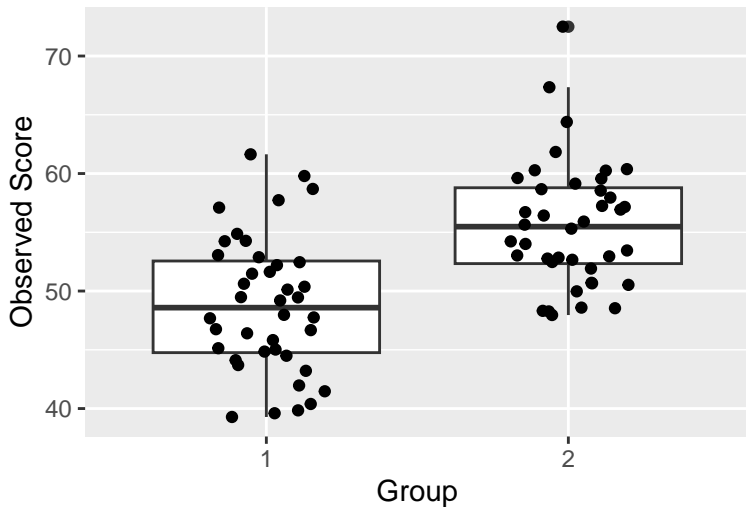
- Accumulate evidence for the absence of an effect
- Draw valid conclusions with almost no data
- Say things about means without knowing much about the basic units

Means without knowing much about the basic units

```
N <- 80
group <- rep(1:2, each = N/2) # Make two groups
# True score for each participant depends on group
true_score <-
  rnorm(N,
        mean = rep(c(50, 55), each = N/2),
        sd    = 2)
# Observations are noisy
observed_score <- true_score +
  rnorm(N, mean = 0, sd = 5)

data <- data.frame(participant = 1:N,
                   group = factor(group),
                   true_score, observed_score)
```

Observed Scores by Group



```
model_string <- "  
model {  
  for (i in 1:N) {  
    # Likelihood  
    observed_score[i] ~ dnorm(true_score[i], tau)  
    # Hierarchical model  
    true_score[i]      ~ dnorm(mu[group[i]], tau_true)  
  }  
  diff_mu ~ dnorm(0, 0.001)  
  mu[1]    ~ dunif(0, 100)  
  mu[2] <- mu[1] + diff_mu  
  tau      ~ dgamma(0.1, 0.1)  
  tau_true ~ dgamma(0.1, 0.1)  
}  
"
```

```
data_list <- list(  
  N = N,  
  observed_score = data$observed_score,  
  group = as.numeric(data$group)  
)  
  
# Initial values  
inits <- function() {  
  list(  
    tau = rgamma(1, 0.1, 0.1),  
    tau_true = rgamma(1, 0.1, 0.1),  
    true_score = runif(N, 0, 100)  
  )  
}
```



```
jags_model <- jags.model(textConnection(model_string),  
                          data      = data_list,  
                          inits     = inits,  
                          n.chains = 3 ,  
                          n.adapt  = 1000 )
```

```
## Compiling model graph  
##   Resolving undeclared variables  
##   Allocating nodes  
## Graph information:  
##   Observed stochastic nodes: 80  
##   Unobserved stochastic nodes: 84  
##   Total graph size: 250  
##  
## Initializing model
```

```

update(jags_model, n.iter = 1000)
samples <- coda.samples(jags_model,
                        variable.names = c("diff_mu",
                                           "true_score"),
                        n.iter = 5000)

```

```

# Summarize the results
summary_diff_mu <- summary(samples)$
  statistics["diff_mu", ]
print(summary_diff_mu)

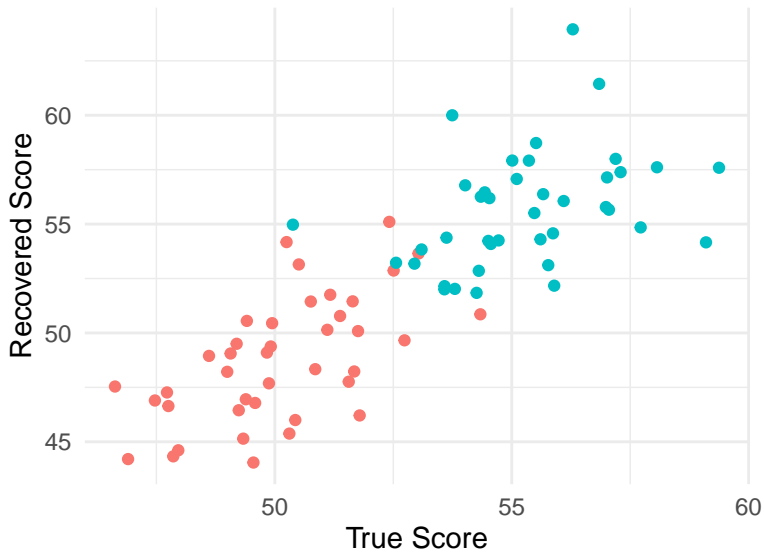
```

##	Mean	SD	Naive SE	Time-series SE
##	6.875807265	1.214229771	0.009914145	0.045074384

```
# Extract true scores from the posterior samples
true_scores_posterior <- as.data.frame(as.matrix(samples)) %>%
  select(starts_with("true_score")) %>%
  apply(2, mean)

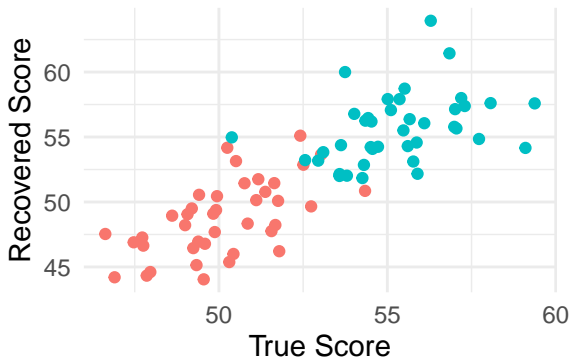
# Add the recovered true scores to the data frame
data$recovered_true_score <- true_scores_posterior

# Scatter plot of true scores vs recovered true scores
p2 <- ggplot(data, aes(x = true_score, y = recovered_true_score,
                       color = group)) +
  geom_point() +
  labs(x = "True Score", y = "Recovered Score") +
  theme_minimal() +
  theme(legend.position = "none")
```



Hierarchical recovery beats individual recovery

The difference between groups is well recovered as 6.8758073 ± 1.2142298 , even though the true scores within each group are poorly recovered.



Perform inference with no actual data

Suppose these were scores from a test (and suppose the two groups are an “on-track” group and an “advanced” group). One student from the advanced group missed class. What do we know about student $N + 1$?

```

model_string <- "
model {
  for (i in 1:N) {
    # Likelihood
    observed_score[i] ~ dnorm(true_score[i], tau)
    # Hierarchical model
    true_score[i]      ~ dnorm(mu[group[i]], tau_true)
  }
  diff_mu ~ dnorm(0, 0.001)
  mu[1]    ~ dunif(0, 100)
  mu[2] <- mu[1] + diff_mu
  tau      ~ dgamma(0.1, 0.1)
  tau_true ~ dgamma(0.1, 0.1)

  true_score[N+1] ~ dnorm(mu[2], tau_true)
  observed_score[N+1] ~ dnorm(true_score[N+1], tau)
}

```

```
data_list <- list(  
  N = N,  
  observed_score = c(data$observed_score, NA),  
  group = as.numeric(data$group)  
)  
  
# Initial values  
inits <- function() {  
  list(  
    tau = rgamma(1, 0.1, 0.1),  
    tau_true = rgamma(1, 0.1, 0.1),  
    true_score = runif(N+1, 0, 100)  
  )  
}
```



```
jags_model <- jags.model(textConnection(model_string),  
                          data      = data_list,  
                          inits     = inits,  
                          n.chains = 3 ,  
                          n.adapt  = 1000 )
```

```
## Compiling model graph  
##   Resolving undeclared variables  
##   Allocating nodes  
## Graph information:  
##   Observed stochastic nodes: 80  
##   Unobserved stochastic nodes: 86  
##   Total graph size: 252  
##  
## Initializing model
```

```
update(jags_model, n.iter = 1000)
samples <- coda.samples(jags_model,
                        variable.names = c("diff_mu",
                                           "true_score",
                                           "observed_score"),
                        n.iter = 5000)
```

Summarize the results

```
summary_new <- summary(samples)$
  statistics[c("true_score[81]", "observed_score[81]"),
            c("Mean", "SD")]
print(summary_new)
```

##	Mean	SD
## true_score[81]	55.62454	4.451111
## observed_score[81]	55.66355	5.828375

True score and observed score

