

Bayesian Nonparametric and Meta Analyses of COVID-19 Studies

Haixin Yu

12/20/2020

Introduction

Background

Coronavirus disease 2019 (COVID-19) is a contagious disease caused by severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2). The first case was reported in Wuhan, China, in December 2019. Then the disease rapidly spread across the world. On March 11, 2020, the World Health Organization (WHO) declared the outbreak a pandemic. As of today, a total of 76,381,409 and 17,702,516 confirmed cases have been recorded globally and in the United States. The rapid spread of the virus is largely attributed to the fact that infected patients may be asymptomatic. An asymptomatic patient is defined as a person with laboratory-confirmed COVID-19 infection that has no symptoms at the time of first clinical assessment and at the end of follow-up. Although asymptomatic COVID-19 infections have been reported in various studies, the proportion of asymptomatic infections varied widely between 5% and 95%. Thus, we conducted the Bayesian nonparametric and meta analyses to answer three questions:

- 1) What is the proportion of asymptomatic infections among people who infected with COVID-19?
- 2) What are the factors of asymptomatic infections in COVID-19?
- 3) Most importantly, is it reasonable to synthesize all the studies to produce a single estimate?

Data Source

The file listing all the studies used for my analysis is available from the Harvard Dataverse Database. Four online databases including PubMed, Embase, bioRxiv and medRxiv were searched through 31 July 2020. All included studies reported a specific number of asymptomatic COVID-19 infections. There was no language restriction. Risk of bias may be included in studies.

Analysis and Results

Data Preprocessing

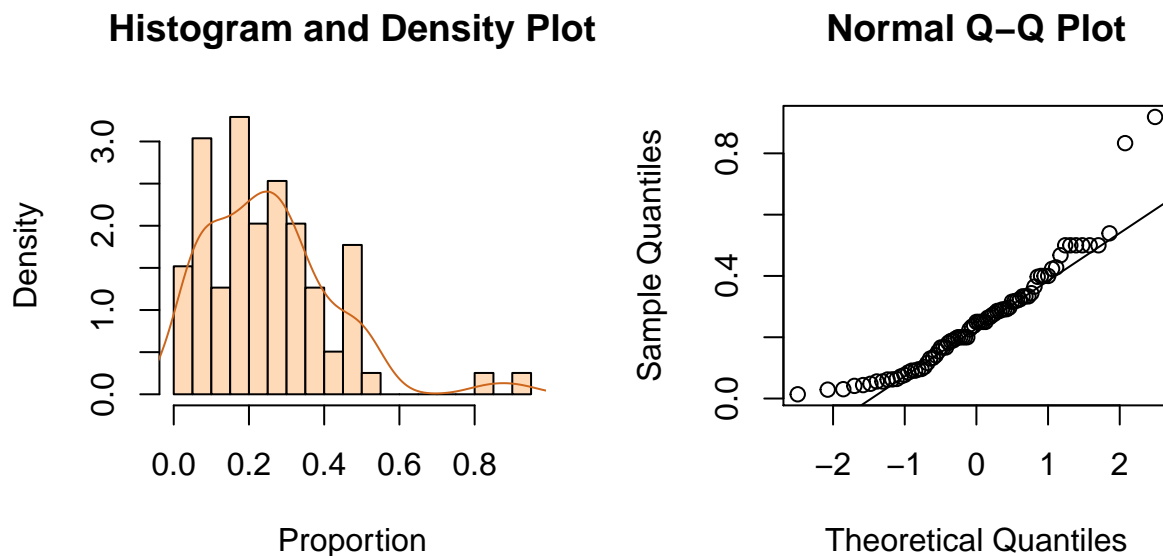
All analyses were conducted using R version 4.0.2. A total of 6,602 patients from 79 studies were included. The original dataset had 7 variables. I separated the variable '**setting**' into two variables: '**code**' and '**setting**'. I added a new variable '**proportion**' that calculated by dividing the number of asymptomatic infections by the number of total COVID-19 infections. Thus, there are 79 observations of 9 variables in our dataset. The nine variables are:

- **id**: record id
- **author**: first author
- **collection**: setting of data collection
- **code**: setting code
- **setting** (seven study settings):
 - contact investigation
 - contact investigation, aggregated

- outbreak investigation
- screening
- hospitalized adults
- hospitalized children
- hospitalized children and adults
- **cases:** the number of asymptomatic infections
- **total:** the number of total infections
- **proportion:** cases/total
- **source** (four databases):
 - PubMed
 - Embase
 - bioRxiv
 - medRxiv

For the seven categories of the variable '**setting**': contact investigation included 16 studies, contact investigation, aggregated included 9 studies, outbreak investigation included 12 studies, screening included 7 studies, hospitalized adults included 15 studies, hospitalized children included 10 studies, and hospitalized children and adults included 10 studies.

The histogram and density plot below gave us a general idea of how the data was distributed. The normal Q-Q plot below examined the normality of the data. I performed two types of analysis on the dataset. The first one was meta-analysis using **meta** and **metafor** packages. The second one was Bayesian nonparametric analysis using **DPpackage**.

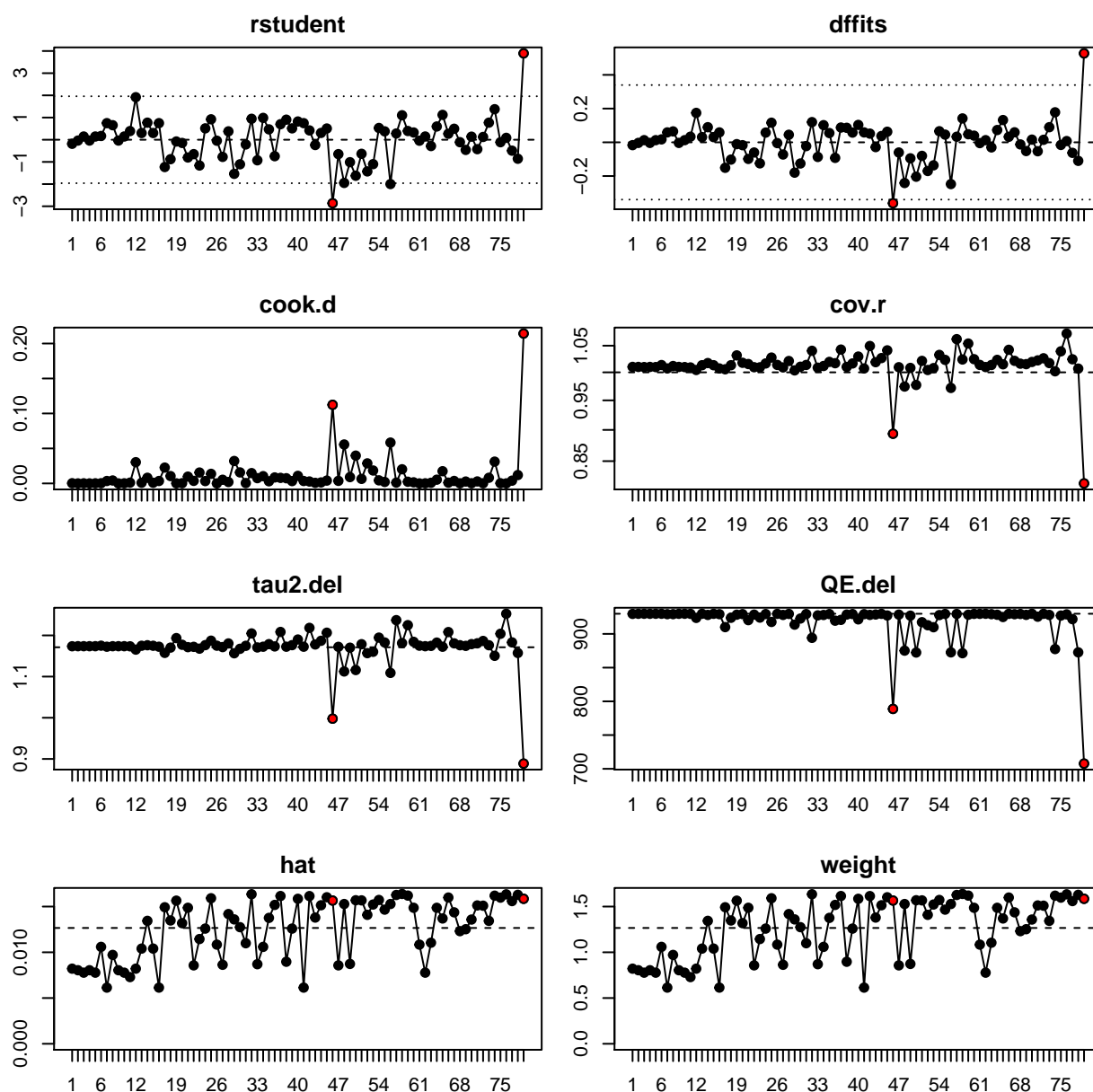


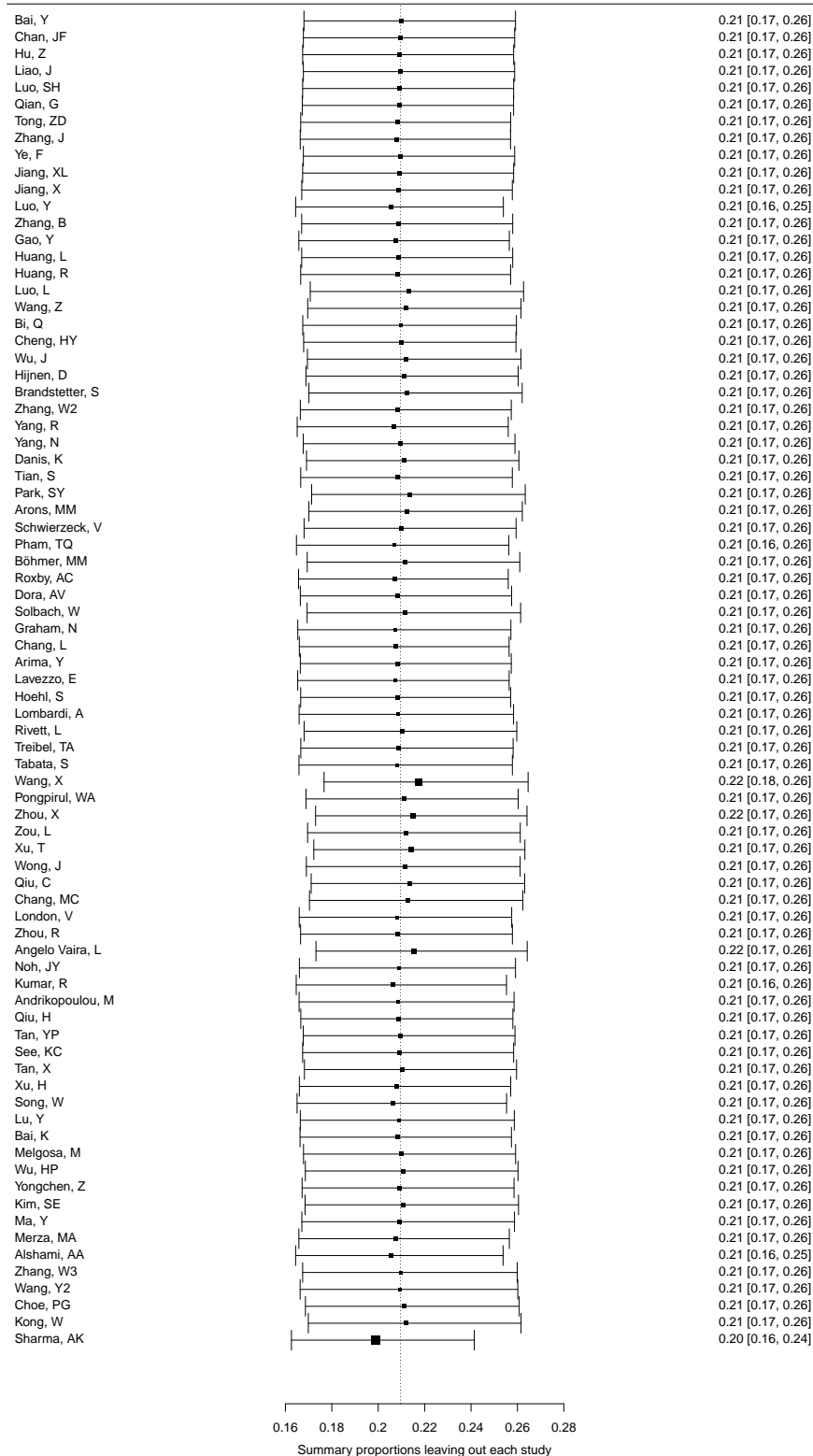
Meta Analysis Using meta and metafor Packages

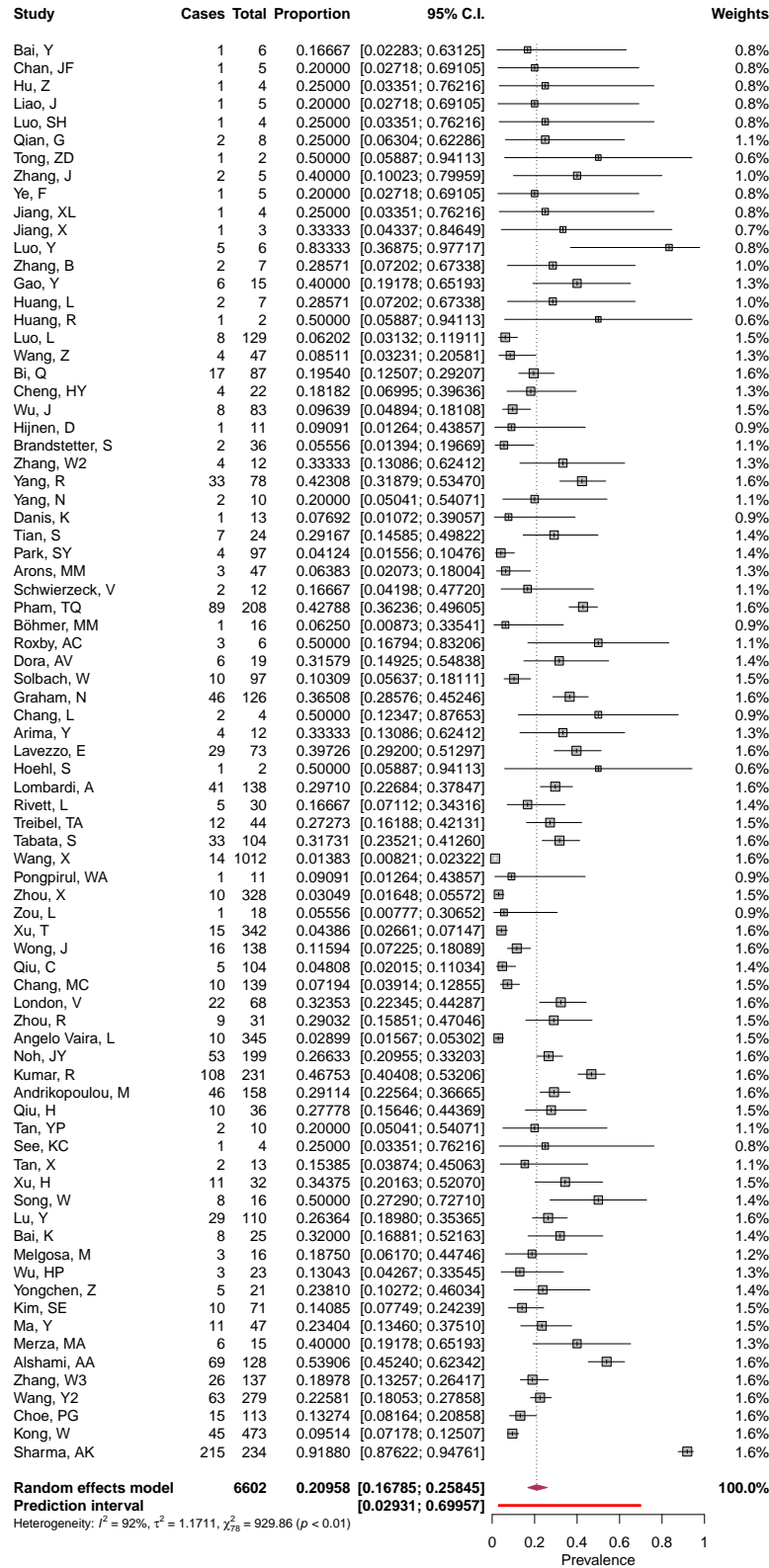
Based on the QQ plot above, the proportion of asymptomatic infections was transformed using the logit transformation to make it conform to the normal distribution. The **escalc()** function in the **metafor** package was used to calculate individual effect sizes (y_i) and its sampling variances (v_i). The overall effect size and its 95% confidence interval was based on the random-effects-model using DerSimonian-Laird estimator and normal approximation interval respectively. Unlike the fixed-effects-model, which assumes all studies along with their effect sizes stem from a single homogeneous population, the random-effects-model assumes a distribution of true effect sizes from a “universe” of populations, which takes both within- and between-study variances into account.

In the forest plot below, we displayed the overall estimate of asymptomatic COVID-19 infections through the `metaprop()` and `forest()` function in the `meta` package. Of 6,602 confirmed cases, 1,273 were defined as asymptomatic infections. The overall estimate of the proportion of people who become infected with COVID-19 and remain asymptomatic throughout the course of infection was 20.96%, with a 95% confidence interval of 16.8% to 25.8% and a prediction interval of 3.0% to 70.0%. Apparently, the estimates of asymptomatic infections were quite different from study to study in the forest plot. The high heterogeneity was also confirmed by the τ^2 and I^2 statistic ($\tau^2 = 1.17$, $I^2 = 92\%$, $p < 0.01$).

There are three methods to quantify heterogeneity: 1) the test for heterogeneity (Q), 2) the estimate of between-study variance (τ^2), 3) the estimate for the proportion of the observed variability that reflects the between-study variance (I^2). Together, they can inform us if the effects are consistent. According to our results, it is not appropriate to combine all the studies to make a summary estimate of the proportion of asymptomatic COVID-19 infections. To identify and visually inspect the influential studies, leave-one-out analysis and diagnostic test were used.



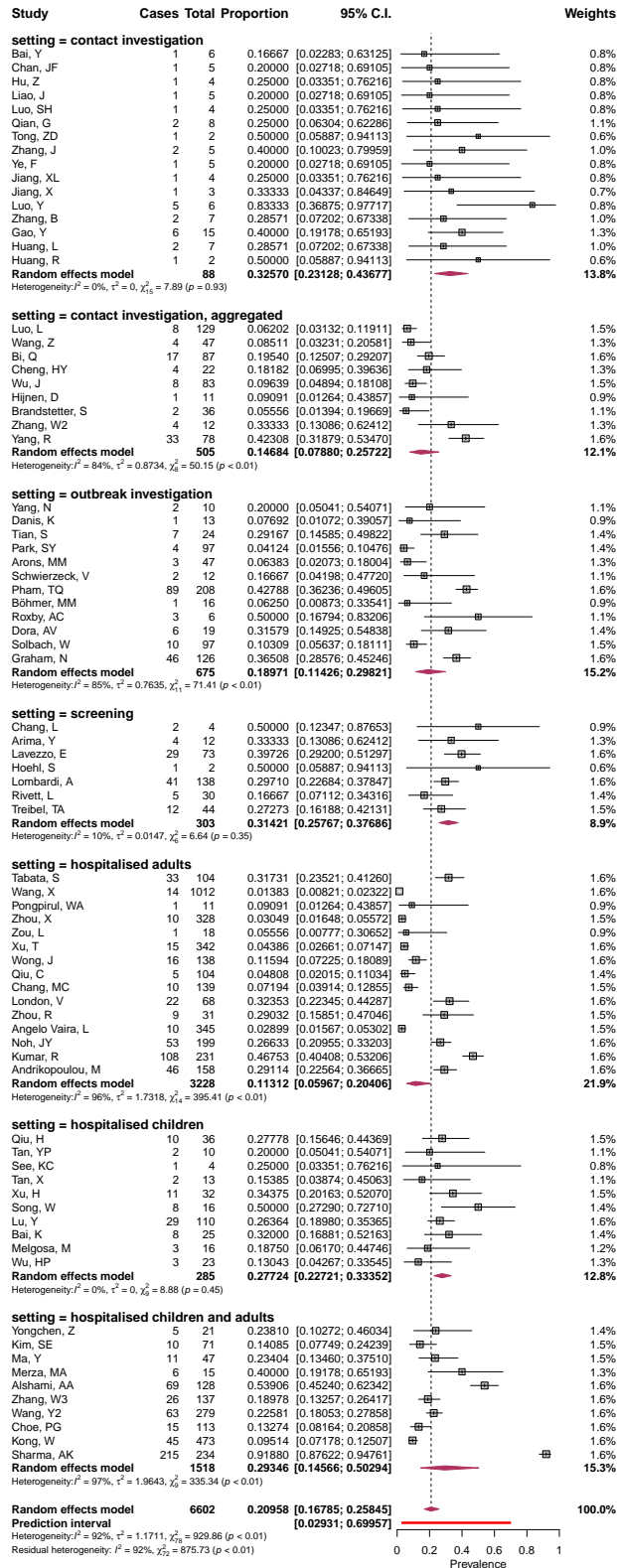




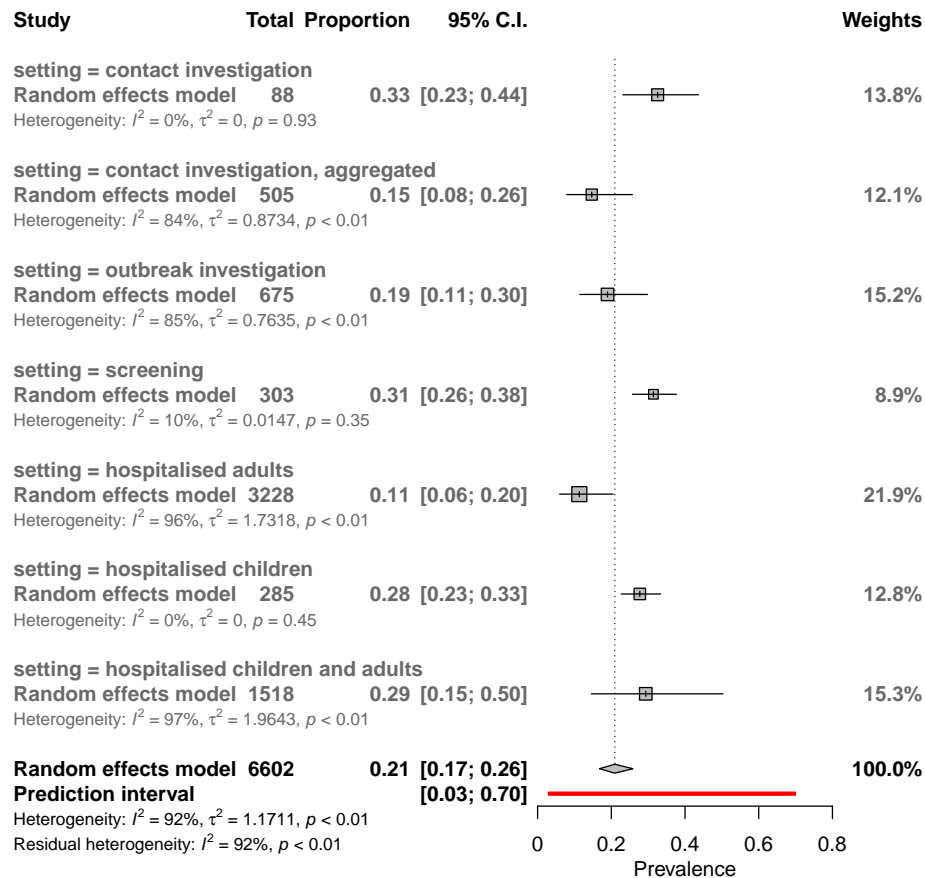
Subgroup Analysis

Subgroup analysis were conducted to explain potential factors contributing to heterogeneity. There are seven defined subgroups based on the method of selection of asymptomatic status. As shown in the forest plot in presence of subgroups, the summary estimate of the proportion of people with asymptomatic COVID-19 infections varied from setting to setting. The first three subgroups involved studies that reported on different types of contact investigation. In 16 studies of single-family contact investigation, the summary estimate was 33% (95% CI 23%-44%). In 9 studies of aggregated contact investigation from both asymptomatic and symptomatic people, the summary estimate was 15% (95% CI 8%-26%). There were 12 studies that reported on outbreak investigation arising from a single symptomatic case in nursing homes or occupational settings. The summary estimate of these studies was 19% (95% CI 11%-30%). People with asymptomatic COVID-19 infections that detected through screening in defined populations were reported in 7 studies. The screened populations included healthcare workers, evacuees from transmission sites, and the entire population of an Italian village. The summary estimate of these studies was 31% (95% CI 26%-38%). The remaining studies were done in hospital settings, included adult only (15 studies), children only (10 studies), and children and adults (10 studies). The summary estimates of these subgroups were 11% (95% CI 6%-20%), 28% (95% CI 23%-33%), and 29% (95% CI 15%-50%) respectively.

```
forest(pes.summary,
      xlim=c(0,1),
      leftcols=c("studlab", "event", "n", "effect", "ci"),
      leftlabs=c("Study", "Cases", "Total", "Proportion", "95% C.I."),
      rightcols=c("w.random"),
      rightlabs=c("Weights"),
      xlab="Prevalence",
      prediction=TRUE,
      col.predict.lines="red",
      fs.xlab=12,
      fs.study=12,
      fs.study.lables=12,
      fs.heading=14,
      squaresize=0.5,
      col.square="grey",
      col.square.lines="black",
      col.diamond="maroon",
      col.diamond.lines="maroon",
      col.by="black",
      pooled.totals=TRUE,
      comb.fixed=FALSE,
      lty.random=2,
      type.study="square",
      type.random="diamond",
      ff.random="bold",
      hetlab="Heterogeneity:",
      fs.hetstat=10,
      smlab="",
      print.tau2=TRUE,
      print.Q=TRUE,
      print.pval.Q=TRUE,
      print.I2=TRUE,
      digits.Q=2,
      digits=5)
```



Based on the subgroup results showing in the forest plot, it is more reasonable to make summary estimates of the proportion of asymptomatic COVID-19 infections in different settings. Contact investigation and screening were both at a relatively high asymptomatic infection rate of 33% and 31% respectively. The summary estimate of the proportion of asymptomatic infections in children (28%) is much higher than in adults (11%). However, we can still observe some overlap in the confidence intervals of these subgroups, which was confirmed by the significant heterogeneity among studies conducted in aggregated contact investigation ($I^2 = 84\%$, $p < 0.01$), outbreak investigation ($I^2 = 85\%$, $p < 0.01$), hospitalized adults ($I^2 = 96\%$, $p < 0.01$) and hospitalized children and adults ($I^2 = 97\%$, $p < 0.01$). That is, when we grouped the included studies according to different study settings, there was no significant difference in effect sizes between the seven subgroups. This was also supported by the significant unexplained heterogeneity left in the data ($I^2 = 92\%$, $p < 0.01$), which implied that there might be one or more important missing factors could better explain the heterogeneity.



Bayesian Nonparametric Analysis Using DPmeta {DPpackage}

Bayesian nonparametric analysis was conducted using **DPpackage**. This package contains a series of Bayesian nonparametric models under Dirichlet process. A Bayesian nonparametric model is a Bayesian model on an infinite dimensional parameter space. Its high flexibility and robustness allow for better data modeling than parametric models. A Dirichlet process is a probability distribution whose range is itself a set of probability distributions. It means that instead of generating a single parameter, a single draw from the Dirichlet process outputs another distribution. Therefore, the **DPmeta()** function in the **DPpackage** will generate a posterior density sample for a semiparametric linear mixed effects meta-analysis model using a Dirichlet process or a mixture of Dirichlet process prior for the distribution of the random effects. The highest posterior density (HPD) interval will be used to represent the Bayesian credible intervals. It is the shortest interval among all of the Bayesian credible intervals.

The **DPmeta()** function includes five arguments:

- **formula**: a two-sided linear formula object.
- **prior**: a list giving the prior information.
 - alpha (the value of the precision parameter)
 - mu (the value of the mean of the centering distribution)
 - sigma2 (the value of the variance of the centering distribution)
- **mcmc**: a list giving the MCMC parameters.
 - nburn (the number of burn-in scans)
 - nskip (the thinning interval)
 - nsave (the total number of scans to be saved)
 - ndisplay (the number of saved scans to be displayed on screen)
- **state**: a list giving the current value of the parameters.
- **status**: a logical variable.
 - the run is new (TRUE)
 - the continuation of a previous analysis (FALSE)

formula describes how to predict the proportion of asymptomatic COVID-19 infections (y). **prior** gives the prior information (alpha, mu, tau1 and tau2). **mcmc** gives the MCMC parameters (nburn, nsave, nskip, ndisplay). **state** is set to NULL. **status** indicates the run is new. Since we know the distribution of the included studies, we can fairly estimate the value of the mean of the centering distribution (mu) and the hyperparameters for the prior distribution of the variance of the centering distribution (tau1 and tau2). Then by tuning the precision parameter (alpha), the minimum number of clusters was obtained when alpha was 0.01 (See alpha Tuning Plot). The model summary indicated that the overall estimate of the proportion of people who become infected with COVID-19 and remain asymptomatic throughout the course of infection was 22.15%, with a 95% HPD interval of 18.1% to 26.3%. Five clusters were generated. To extract the predictive information of random effects, we used the **DPrandom()** function in the **DPpackage**. It showed that the prediction of the proportion of asymptomatic COVID-19 infection was 22.19%, along with a 95% HPD interval of 1.5% to 46.0%. The trace and density of model parameters are shown in the model summary plots below.

```
set.seed(1234)
# Prior information
prior = list(alpha=0.01, mu=mean(effects), tau1=0.01, tau2=0.01)

# Initial state
state = NULL

# MCMC parameters
nburn = 20000    # the number of burn-in scans
nskip = 20       # the thinning interval
nsave = 10000    # the total number of scans to be saved
ndisplay = 400   # the number of saved scans to be displayed on screen
```

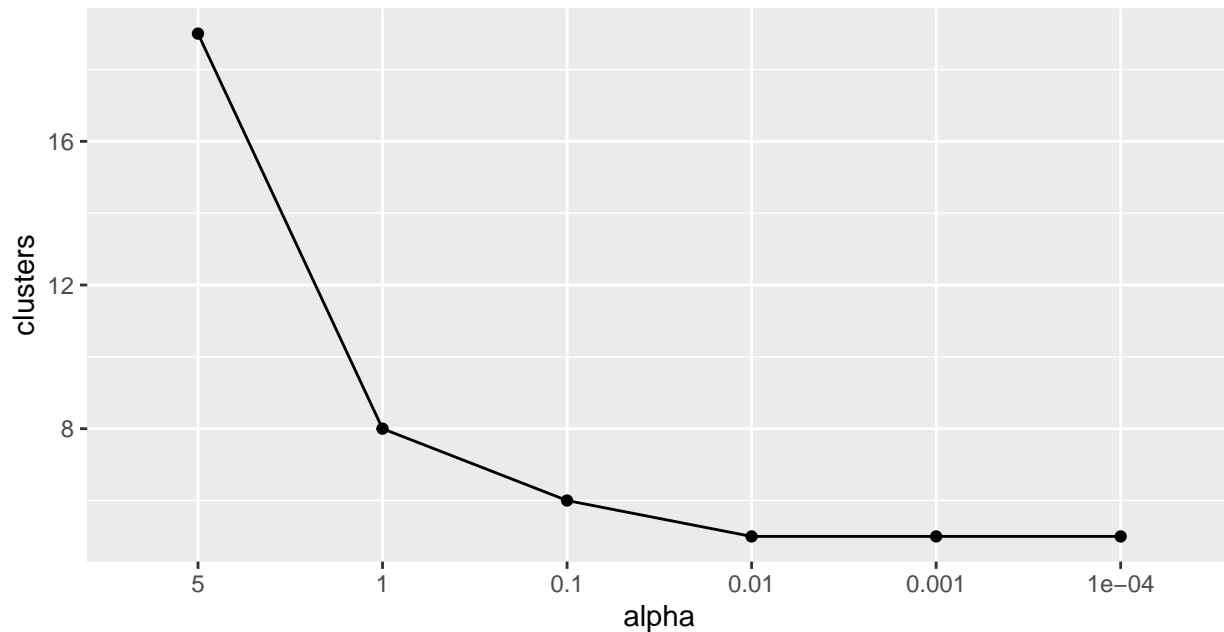
```
mcmc = list(nburn=nburn, nsave=nsave, nskip=nskip, ndisplay=ndisplay)

# Fit the model: First run
fit = DPmeta(formula=y~1, prior=prior, mcmc=mcmc, state=state, status=TRUE)
```

```
##
## MCMC scan 400 of 10000 (CPU time: 1.578 s)
## MCMC scan 800 of 10000 (CPU time: 2.062 s)
## MCMC scan 1200 of 10000 (CPU time: 2.547 s)
## MCMC scan 1600 of 10000 (CPU time: 3.031 s)
## MCMC scan 2000 of 10000 (CPU time: 3.516 s)
## MCMC scan 2400 of 10000 (CPU time: 3.984 s)
## MCMC scan 2800 of 10000 (CPU time: 4.500 s)
## MCMC scan 3200 of 10000 (CPU time: 4.984 s)
## MCMC scan 3600 of 10000 (CPU time: 5.469 s)
## MCMC scan 4000 of 10000 (CPU time: 5.953 s)
## MCMC scan 4400 of 10000 (CPU time: 6.438 s)
## MCMC scan 4800 of 10000 (CPU time: 6.922 s)
## MCMC scan 5200 of 10000 (CPU time: 7.406 s)
## MCMC scan 5600 of 10000 (CPU time: 7.906 s)
## MCMC scan 6000 of 10000 (CPU time: 8.406 s)
## MCMC scan 6400 of 10000 (CPU time: 8.875 s)
## MCMC scan 6800 of 10000 (CPU time: 9.375 s)
## MCMC scan 7200 of 10000 (CPU time: 9.844 s)
## MCMC scan 7600 of 10000 (CPU time: 10.328 s)
## MCMC scan 8000 of 10000 (CPU time: 10.812 s)
## MCMC scan 8400 of 10000 (CPU time: 11.312 s)
## MCMC scan 8800 of 10000 (CPU time: 11.797 s)
## MCMC scan 9200 of 10000 (CPU time: 12.281 s)
## MCMC scan 9600 of 10000 (CPU time: 12.781 s)
## MCMC scan 10000 of 10000 (CPU time: 13.250 s)
```

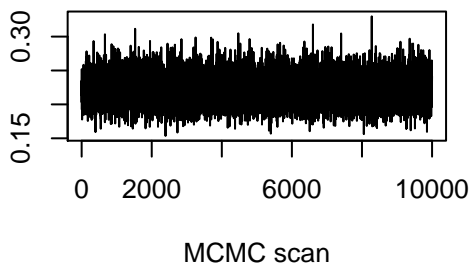
```
fit
```

```
##
## Bayesian semiparametric linear mixed effects meta-analysis
##
## Call:
## DPmeta.default(formula = y ~ 1, prior = prior, mcmc = mcmc, state = state,
##      status = TRUE)
##
## Posterior Inference of Parameters:
##      effects      mu    sigma2  ncluster
##    0.2215    0.2545    0.1882    5.0671
##
## Number of Studies: 79
```

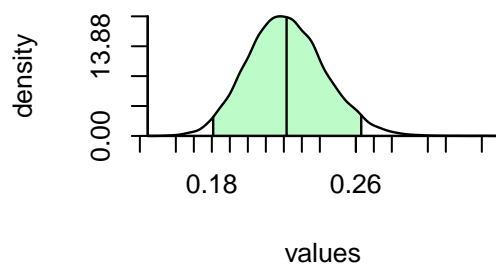


```
##
## Bayesian semiparametric linear mixed effects meta-analysis
##
## Call:
## DPmeta.default(formula = y ~ 1, prior = prior, mcmc = mcmc, state = state,
##   status = TRUE)
##
## Posterior Predictive Distributions (log):
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -2.8917  0.3107  0.7064  0.7165  1.1084  3.2312
##
## Model's performance:
##   Dbar   Dhat    pD     DIC    LPML
##  -189.15 -245.36  56.21 -132.94   56.60
##
## Regression coefficients:
##           Mean      Median   Std. Dev. Naive Std.Error 95%HPD-Low
## effects 0.2215233 0.2205796 0.0213156 0.0002132      0.1807334
##           95%HPD-Upp
## effects 0.2629920
##
## Baseline distribution:
##           Mean      Median   Std. Dev. Naive Std.Error 95%HPD-Low 95%HPD-Upp
## mu      0.254521 0.254521 0.000000 0.000000      0.254521 0.254521
## sigma2 0.188225 0.128675 0.229249 0.002292      0.025506 0.513459
##
## Precision parameter:
##           Mean      Median   Std. Dev. Naive Std.Error 95%HPD-Low
## ncluster 5.067100 5.000000 0.256913 0.002569      5.000000
##           95%HPD-Upp
## ncluster 6.000000
##
## Number of Studies: 79
```

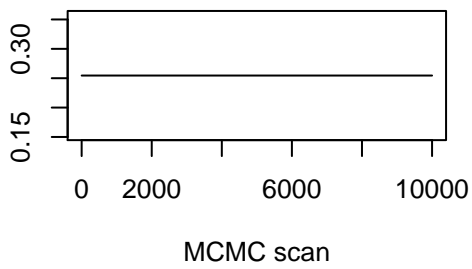
Trace of effects



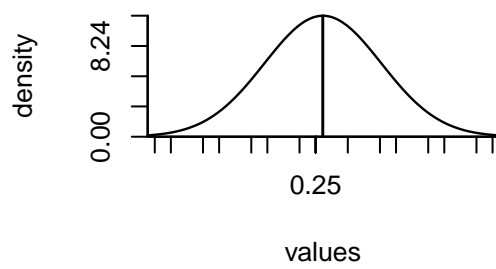
Density of effects



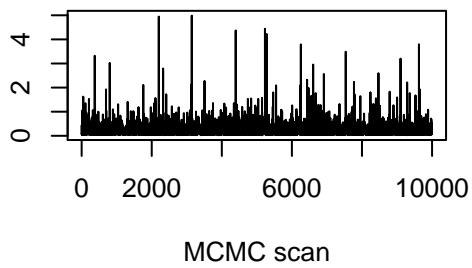
Trace of mu



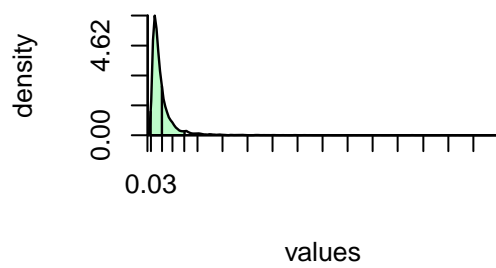
Density of mu



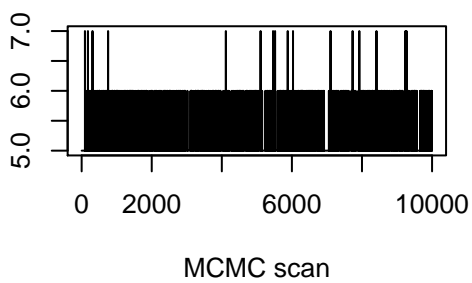
Trace of sigma2



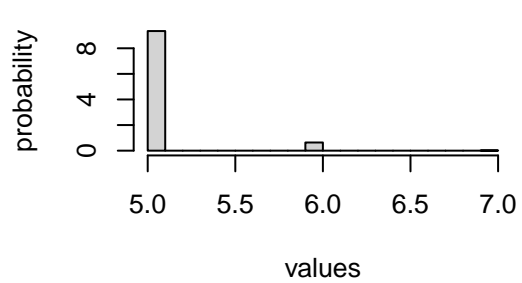
Density of sigma2

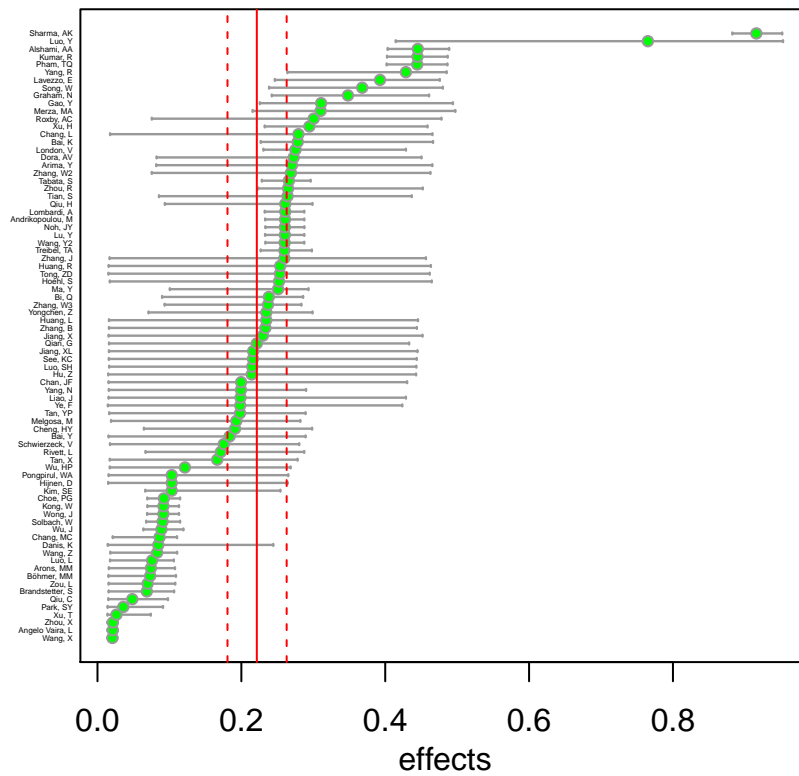


Trace of ncluster

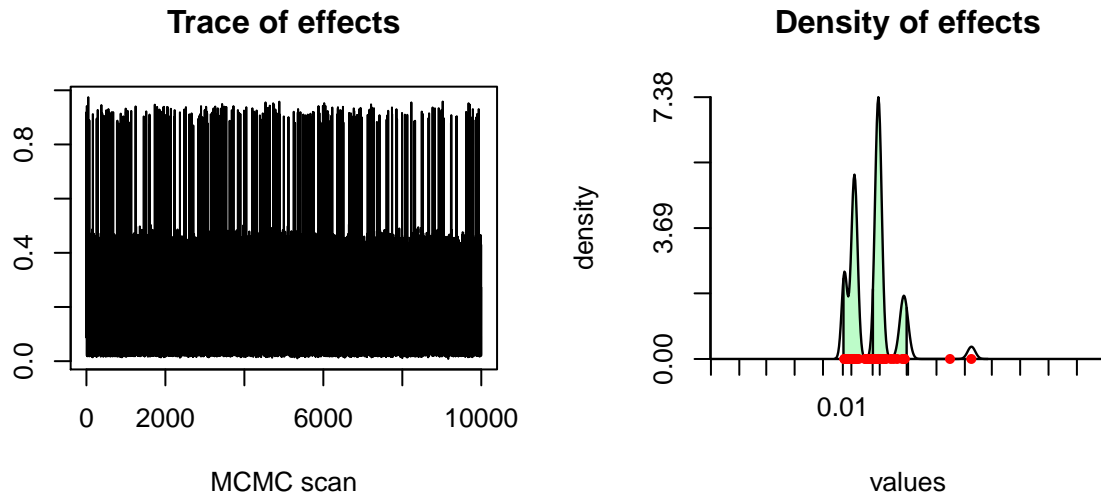


Density of ncluster



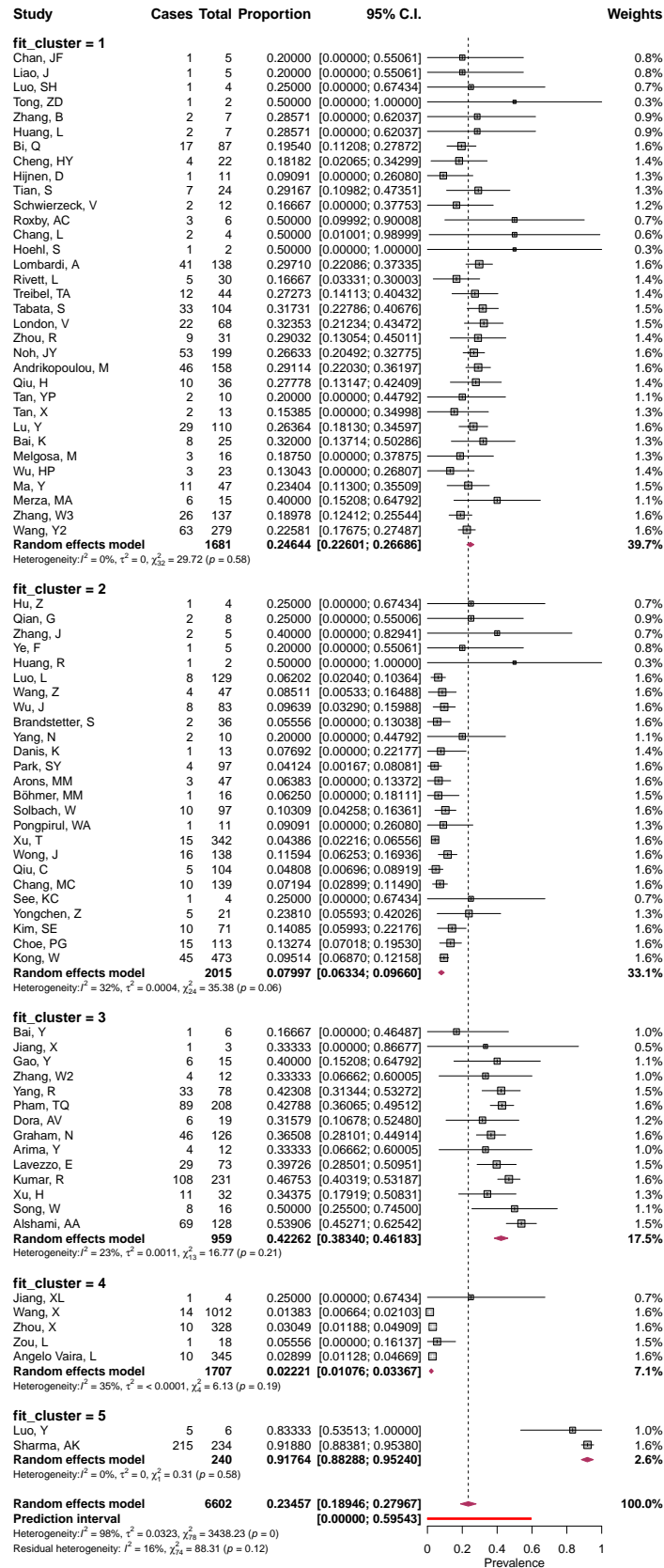


```
##
## Random effect information for the DP object:
##
## Call:
## DPmeta.default(formula = y ~ 1, prior = prior, mcmc = mcmc, state = state,
##   status = TRUE)
##
##
## Predictive distribution:
##      Mean      Median   Std. Dev.  Naive Std.Error  95%HPD-Low  95%HPD-Upp
## theta  0.221928  0.249146  0.167730   0.001677      0.014302   0.459616
```

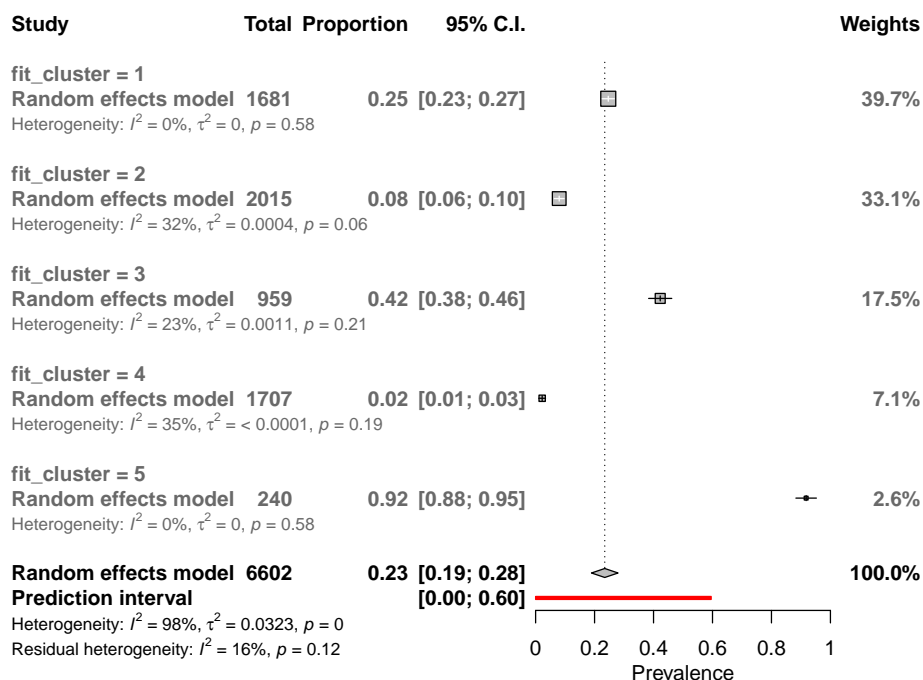


Cluster Analysis

The **DPmeta** model divided the included studies into five clusters. For further study and comparison, I also conducted a subgroup analysis on the five clusters. As shown in the forest plot in presence of clusters, the summary estimates of the proportion of people with asymptomatic COVID-19 infections varied significantly from cluster to cluster. Cluster 2 and cluster 4 were to the left of the overall summary estimate. Cluster 3 and cluster 5 were to the right of the overall summary estimate. Cluster 1 was on the overall summary estimate. There were 33 studies in cluster 1. The summary estimate of these studies was 25% (95% CI 23%-27%). It included studies from all seven settings. In cluster 2, there were 25 studies. The summary estimate of them was 8% (95% CI 6%-10%). It included studies under six settings other than screening. In 14 studies of cluster 3, the summary estimate of the proportion with asymptomatic infections was 42% (95% CI 38%-46%). It also included studies from all seven settings. Cluster 4 had 5 studies with a summary estimate of 2% (95% CI 1%-3%). Two settings were included: contact investigation and hospitalized adults. The remaining two studies were in cluster 5. They came from contact investigation and hospitalized children and adults. The summary estimate of these two studies was 92% (95% CI 88%-95%).



As shown in the forest plot of cluster results below, there were significant differences in the proportion of asymptomatic COVID-19 infections among the five clusters generated by the **DPmeta** model. Cluster 5 had an extremely high asymptomatic infection rate of 92% (95% CI 88%-95%). Cluster 4 had the lowest asymptomatic infection rate of 2% (95% CI 1%-3%). It is reasonable to say that the summary estimates of the proportion of asymptomatic COVID-19 infections based on the five clusters are more appropriate than based on the seven study settings. There was only a small amount of heterogeneity in cluster 2 (32%), cluster3 (23%), and cluster 4 (35%). The residual heterogeneity in the data was also only 16%. However, due to the lack of individual level data on COVID-19, we cannot explain the underlying factors behind the clusters at this stage.



Comparison

The overall summary estimate of the proportion of asymptomatic COVID-19 infections by Bayesian nonparametric analysis using the **DPpackage** (22.15%, 95% HPD 18.1%-26.3%) is slightly higher than by meta-analysis using the **meta** and **metafor** packages (20.96%, 95% CI 16.8%-25.8%). However, the prediction interval of Bayesian nonparametric analysis (3.0%-70%) is much narrower than that of meta-analysis (1.5%-46.0%), which enables a more precise proportion estimate of the asymptomatic COVID-19 infections. Also, cluster analysis based on the **DPmeta** model can better explain the heterogeneity among included studies than subgroup analysis based on study settings. However, the factors behind the clusters cannot be summarized since individual level data on COVID-19 is not available.

Conclusion

We conducted the Bayesian nonparametric and meta analyses to address three questions: 1) What is the proportion of asymptomatic infections among people who infected with COVID-19? 2) What are the factors of asymptomatic infection in COVID-19? 3) And most importantly, is it reasonable to synthesize all the studies to produce a single estimate? The overall summary estimate of the proportion of asymptomatic COVID-19 infections in both methods were around 20%. However, it is not appropriate to synthesize the proportions of all the studies into a single estimate. Subgroup analysis can be conducted to explain the heterogeneity between studies, but not all factors can be summarized in the available data. The variable 'setting' can only be part of the contributing factors. Further analysis is needed as more studies and features become available. Therefore, preventive measures such as masks and social distancing will continue to be needed to reduce transmission.

References

- WHO Coronavirus Disease (COVID-19) Dashboard. <https://covid19.who.int/>
- Müller P, Quintana FA, Jara A, Hanson (2015) Bayesian Nonparametric Data Analysis. Springer, New York
- Harrer, M., Cuijpers, P., Furukawa, T.A., & Ebert, D. D. (2019). Doing Meta-Analysis in R: A Hands-on Guide. DOI: 10.5281/zenodo.2551803. https://bookdown.org/MathiasHarrer/Doing_Meta_Analysis_in_R/
- Riley RD, Higgins JP, Deeks JJ. Interpretation of random effects meta-analyses. *BMJ*. 2011 Feb 10;342:d549. doi: 10.1136/bmj.d549. PMID: 21310794.
- He W, Yi GY, Zhu Y. Estimation of the basic reproduction number, average incubation time, asymptomatic infection rate, and case fatality rate for COVID-19: Meta-analysis and sensitivity analysis. *J Med Virol*. 2020 May 29;10.1002/jmv.26041. doi: 10.1002/jmv.26041. Epub ahead of print. PMID: 32470164; PMCID: PMC7283745.
- Zheng B, Wang H, Yu C. An increasing public health burden arising from children infected with SARS-CoV2: A systematic review and meta-analysis. *Pediatr Pulmonol*. 2020 Dec;55(12):3487-3496. doi: 10.1002/ppul.25008. Epub 2020 Sep 25. PMID: 32757374; PMCID: PMC7436588.
- He J, Guo Y, Mao R, Zhang J. Proportion of asymptomatic coronavirus disease 2019: A systematic review and meta-analysis. *J Med Virol*. 2020;1–11. <https://doi.org/10.1002/jmv.26326>
- Newcombe RG. Two-sided confidence intervals for the single proportion: comparison of seven methods. *Stat Med*. 1998 Apr 30;17(8):857-72. doi: 10.1002/(sici)1097-0258(19980430)17:8<857::aid-sim777>3.0.co;2-e. PMID: 9595616.
- Heneghan C, Brassey J, Jefferson T. Covid-19: What Proportion Are Asymptomatic? Oxford: Centre for Evidence Based Medicine; 2020 [cited 2020 Jul 27]. Available from: <https://www.cebm.net/covid-19/covid-19-what-proportion-are-asymptomatic/>
- Buitrago-Garcia D, Egli-Gany D, Counotte MJ, Hossmann S, Imeri H, Ipekci AM, et al. (2020) Occurrence and transmission potential of asymptomatic and presymptomatic SARS-CoV-2 infections: A living systematic review and metaanalysis. *PLoS Med* 17(9): e1003346. <https://doi.org/10.1371/journal.pmed.1003346>
- Jara, A., Hanson, T. E., Quintana, F. A., Müller, P., & Rosner, G. L. (2011). DPpackage: Bayesian semi- and nonparametric modeling in R. *Journal of Statistical Software*, 40(5), 1-30. <https://doi.org/10.18637/jss.v040.i05>