

# PrivacyUnbiased

## Introduction to PrivacyUnbiased

This package implements methods developed in:

- Evans, Georgina, and Gary King (2020): “*Statistically Valid Inferences from Differentially Private Data Releases*”. In: URL: <https://gking.harvard.edu/dpd>.

In a major development for research data sharing, data providers are beginning to supplement insecure privacy protection strategies, such as “de-identification” with a formal approach called “differential privacy”. One version of differential privacy adds specially calibrated random noise to a dataset, which is then released to researchers. This offers mathematical guarantees for the privacy of research subjects while still making it possible to learn about aggregate patterns of interest. Unfortunately, adding random noise creates measurement error, which, if ignored, induces statistical bias — including in different situations attenuation, exaggeration, switched signs, and incorrect uncertainty estimates. The procedures implemented in **PrivacyUnbiased** account for these biases, producing statistically consistent point estimates from differentially private data.

**PrivacyUnbiased**, which corrects statistical problems with privacy protective procedures added to data, is designed to complement **UnbiasedPrivacy**, which corrects statistical problems with privacy protective procedures added to the results of statistical analyses [Evans et al., Working paper].

## Installing PrivacyUnbiased

To install **PrivacyUnbiased**, run:

```
devtools::install_github("georgieevans/PrivacyUnbiased")
library(PrivacyUnbiased)
```

## Example

We demonstrate the capabilities of **PrivacyUnbiased** by simulating the scenario described above. We start with a hypothetical private data set (**private\_data**). We then add random error to every cell of the data by drawing errors,  $\epsilon_{ik}$ , from a mean 0 normal distribution,  $\epsilon_{ik} \sim \mathcal{N}(0, \sigma_k^2)$ . We set  $\sigma_k$  for each of the  $k$  columns of the data. This produces a differentially private data set (**dp\_data**). In practice, the data analyst would not have access to **private\_data** and would only be able to observe **dp\_data**.

This example data can be loaded into the R environment (after loading the package) by running the following code:

```
# Load the private data
data("private_data")

# Load the DP data
data('dp_data')
```

## lmdp()

**lmdp()** is the primary function of the package. It returns estimates of bias corrected coefficients from differentially private data, alongside several other quantities. Users can interact with it in a similar way to **lm()**. There are only two required inputs, the **formula** and **data**. For instance:

```
lmdp_test <- lmdp(Y ~ X1 + X2 + X3, data = dp_data)
```

You can read the documentation for `lmdp()` by running the code:

```
?lmdp
```

An important distinction between `lmdp()` and `lm()` is that the first row of `data` must indicate the standard deviations of the DP error added to the rest of the data matrix. For instance, if we look at `dp_data`, we see by looking at row 1 that no noise was added to `Y`, the standard error of noise added to `X1` was 0.7, and so on.

```
head(dp_data)
```

```
##           Y           X1           X2           X3
## 1  0.00000 0.700000  1.200000  1.000000
## 2 75.91382 5.783422 13.789933 3.519179
## 3 86.24916 7.654997 14.050685 1.177196
## 4 73.44615 5.953728 8.070672 3.759313
## 5 42.39201 4.804003 16.633445 1.129892
## 6 39.04964 4.083823 12.605909 1.751404
```

An exception to this rule is if the argument `noise` is set to something other than its default (`= NULL`). If `noise = x` (where `x` is any real number), then `lmdp()` will automatically set the error for every column to `x`. In this situation, the first row of the data matrix will be ignored.

The output from `lmdp()` can be summarized using `summary()`, just like a standard `lm` object.

```
summary(lmdp_test)
```

```
##           Estimate Std. Error   t value Pr(>|t|)
## (Intercept)  10.1021     0.1941   52.0427     0
## X1           11.9862     0.0238  503.3449     0
## X2           -2.9960     0.0154 -194.7213     0
## X3            9.0030     0.0292  308.7954     0
```

The additional output from `lmdp()` is stored in a list that can be accessed as follows:

```
# This summarizes the output of an lmdp object
```

```
str(lmdp_test)
```

```
## List of 13
## $ b           : Named num [1:4] 13.9 10.37 -2.1 6.75
##   ..- attr(*, "names")= chr [1:4] "(Intercept)" "X1" "X2" "X3"
## $ b_vcov       : num [1:4, 1:4] 0.019853 -0.000006 -0.000969 -0.000881 -0.000006 ...
## $ beta_tilde   : Named num [1:4] 10.1 12 -3 9
##   ..- attr(*, "names")= chr [1:4] "(Intercept)" "X1" "X2" "X3"
## $ beta_tilde_vcov: num [1:4, 1:4] 0.037679 0.000265 -0.001872 -0.002569 0.000265 ...
## $ var_sims     : num [1:500, 1:8] 14.1 13.8 13.7 13.8 14.1 ...
## $ Sigma_sq_hat : num [1, 1] 3.72
## $ vc_pos_def   : logi TRUE
## $ boot         : logi FALSE
## $ est_vc       : num [1:14, 1:14] 0 0 0 0 0 0 0 0 0 0 ...
## $ Y            : num [1:100001] 0 75.9 86.2 73.4 42.4 ...
## $ X            : num [1:100001, 1:4] 0 1 1 1 1 1 1 1 1 1 ...
##   ..- attr(*, "dimnames")=List of 2
```

```
## .. ..$ : chr [1:100001] "" "2" "3" "4" ...
## .. ..$ : chr [1:4] "(Intercept)" "X1" "X2" "X3"
## $ S : num [1:4, 1:4] 0 0 0 0 0 0.49 0 0 0 0 ...
## $ formula :Class 'formula' language Y ~ X1 + X2 + X3
## .. ..- attr(*, ".Environment")=<environment: R_GlobalEnv>
## - attr(*, "class")= chr "lmdp"

# For instance we can access the variance covariance matrix as follows
lmdp_test$beta_tilde_vcov

##           [,1]           [,2]           [,3]           [,4]
## [1,]  0.0376793753  2.649809e-04 -1.872160e-03 -2.569057e-03
## [2,]  0.0002649809  5.670616e-04 -2.709794e-04  2.044941e-05
## [3,] -0.0018721602 -2.709794e-04  2.367360e-04 -6.382436e-06
## [4,] -0.0025690567  2.044941e-05 -6.382436e-06  8.500320e-04
```

## The impact of bias correction

It is informative to compare `lmdp()` estimates to the estimates produced from `lm()` that do not adjust for the random error in `dp_data`:

```
lm_test <- lm(Y ~ X1 + X2 + X3, data = dp_data)

# Biased OLS estimates
round(summary(lm_test)$coef, 4)
```

```
##           Estimate Std. Error   t value Pr(>|t|)
## (Intercept)  13.8938     0.1552   89.4991     0
## X1           10.3655     0.0170  611.4476     0
## X2           -2.1032     0.0111 -189.5703     0
## X3             6.7456     0.0184  366.8663     0
```

*# Notice that if we set noise = 0, lmdp gives the same point estimates as lm()  
# Standard errors differ since we use a different estimate procedure*

```
lmdp_test_0 <- lmdp(Y ~ X1 + X2 + X3, data = dp_data, noise = 0)

summary(lmdp_test_0)
```

```
##           Estimate Std. Error   t value Pr(>|t|)
## (Intercept)  13.8938     0.1495   92.9225     0
## X1           10.3655     0.0162  639.6195     0
## X2           -2.1032     0.0108 -195.6122     0
## X3             6.7456     0.0179  376.5787     0
```

We can compare the `lmdp()` estimates and `lm()` estimates to the unbiased estimates on private data.

```
lm_true <- lm(Y ~ Z1 + Z2 + Z3, data = private_data)

# We see that the lmdp estimates are very close to the lm estimates on private data
```

```
# In contrast, the lm estimates appear biased
round(summary(lm_true)$coef, 4)
```

```
##           Estimate Std. Error    t value Pr(>|t|)
## (Intercept)  10.0071     0.0283   353.6061     0
## Z1          11.9973     0.0032  3757.4418     0
## Z2          -2.9997     0.0021 -1421.4332     0
## Z3           9.0020     0.0037  2454.6365     0
```

## Variance estimation

The default setting of `lmdp()` is to estimate the standard errors using the simulation method developed in Evans and King [Working paper]. We also offer the option to bootstrap the standard errors by setting the argument `bootstrap_var` to `TRUE`. In general the two methods will produce very similar estimates. The advantage of the simulation method is computational. For large datasets, the bootstrap is essentially infeasible without access to large amounts of computing power. In contrast, the computational time of our simulation procedure scales only slowly in dataset size.

```
# Timing simulation variance estimation
system.time(simulation <- lmdp(Y ~ X1 + X2 + X3, data = dp_data))
```

```
##      user  system elapsed
##   1.503   0.128   2.162
```

```
# Timing bootstrap variance estimation
system.time(bootstrap <- lmdp(Y ~ X1 + X2 + X3, data = dp_data, bootstrap_var = TRUE))
```

```
##      user  system elapsed
##  31.296   7.090  40.181
```

```
# Bootstrap takes ~30 times longer than simulation for dataset of size N = 100000
```

```
# The standard error estimates are similar between the two methods:
```

```
# Bootstrap Std. Error
summary(bootstrap)[, "Std. Error"]
```

```
## (Intercept)          X1          X2          X3
##      0.1936      0.0214      0.0147      0.0292
```

```
# Simulation Std. Error
summary(simulation)[, "Std. Error"]
```

```
## (Intercept)          X1          X2          X3
##      0.2019      0.0219      0.0153      0.0295
```

On small datasets with a relatively large amount of DP error, the variance-covariance matrix we estimate as a parameter to draw random variables may not be positive definite. If this happens, then we use the function `nearPD()` from the package `Matrix`, which finds a close positive definite matrix [Bates and Maechler, 2019]. `lmdp()` will produce the warning message `VC matrix not positive definite` to alert users to this. The function also returns an indicator variable which records whether the matrix was positive definite which can be accessed as follows:

```
lmdp_test$vc_pos_def
```

```
## [1] TRUE
```

## Variable transformation

As discussed in Evans and King [Working paper], transforming variables with random error poses additional complications for estimation. `PrivacyUnbiased` can currently accomodate two types of variable transformation: interaction variables, and squared variables. For example:

```
# Interaction variable
lmdp_interaction <- lmdp(Y ~ X1 + X2 + X3 + X1*X2, data = dp_data)
summary(lmdp_interaction)
```

	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	9.2477	0.4265	21.6805	0.0000
## X1	12.1158	0.0625	193.9338	0.0000
## X2	-2.9425	0.0275	-107.0530	0.0000
## X3	9.0034	0.0296	304.4886	0.0000
## X1:X2	-0.0076	0.0034	-2.2372	0.0253

```
# lmdp with interactions produces similar estimates to lm on private data
# Standard errors are lower since Z's do not contain random noise
lm_interaction <- lm(Y ~ Z1 + Z2 + Z3 + Z1*Z2, data = private_data)
round(summary(lm_interaction)$coefficients, 4)
```

	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	9.9577	0.0622	159.9740	0.0000
## Z1	12.0047	0.0090	1334.0903	0.0000
## Z2	-2.9966	0.0041	-736.7444	0.0000
## Z3	9.0020	0.0037	2454.6331	0.0000
## Z1:Z2	-0.0004	0.0005	-0.8903	0.3733

Other variable transformations, or multiple variable transformations, are not allowed in this version of the package and their inclusion will induce an error message. Note also that `lmdp()` currently only supports bootstrap estimation when the model includes transformed variables. For future release, we are working on expanding the set of admissible variable transformations and introducing the simulation approach to variance estimation for these cases.

## Descriptive statistics

`PrivacyUnbiased` also includes a function to calculate estimates of descriptive statistics of the *private data*, using the moment estimation methods described in Evans and King [Working paper]. The primary function for descriptive statistics is `descriptiveDP()`, which takes two arguments, the variable name and the data frame.

```
# Estimate Descriptive statistics
descriptiveDP(X3, dp_data)
```

```
## Loading required package: polynom
```

	Mean	Std.Dev	Skewness	Kurtosis
## Before DP noise (estimated)	2.9867	1.7250	0.5924	3.4429
## After DP noise (observed)	2.9867	1.9939	0.3836	3.2481

```

# We can compare these estimates to the true values from private_data:
true_descriptive <- round(c(mean(private_data$Z3),
  sd(private_data$Z3),
  moments::skewness(private_data$Z3),
  moments::kurtosis(private_data$Z3)), 4)

names(true_descriptive) <- c('Mean', 'Std.Dev', 'Skewness', 'Kurtosis')

true_descriptive

##      Mean  Std.Dev Skewness Kurtosis
##  2.9933   1.7256   0.5798   3.3534

```

We will look at simulated heteroskedastic data with zero-inflated covariates to demonstrate methods described in Evans and King [Working paper] for data of this nature. The simulated data can be loaded as follows.

```

# Load the private data
data('private_data2')

# Load the differentially private data
data('dp_data2')

```

## Estimating private histograms

The function `distributionDP()` estimates private histograms from the differentially private data. There are three required arguments. `variable` should be the name of the column in the data, `data` is that data name, and `distributions` is a vector of distributions to parameterize.

```

set.seed(02138)

dist_test <- distributionDP(variable = X1, data = dp_data2, distributions = c('Normal', 'Poisson', 'ZIP'))

## Warning in paramsZIP(X = X, S = S, R = moments_fit, moments_df =
## moments_df, : ZIP estimates outside logical bounds

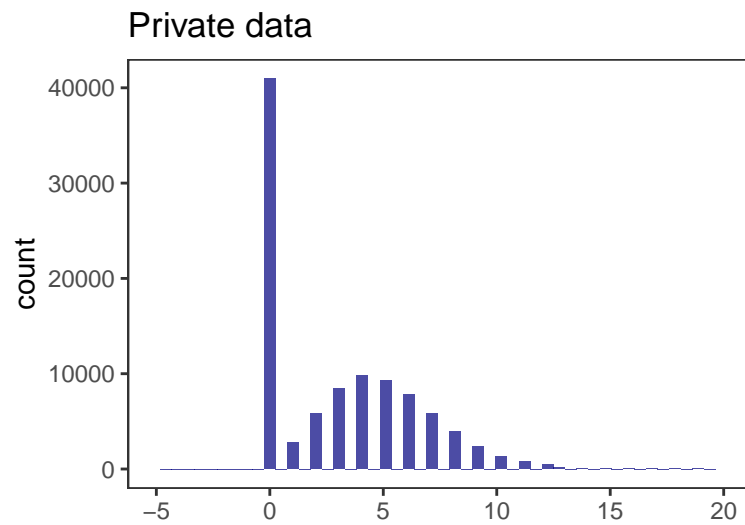
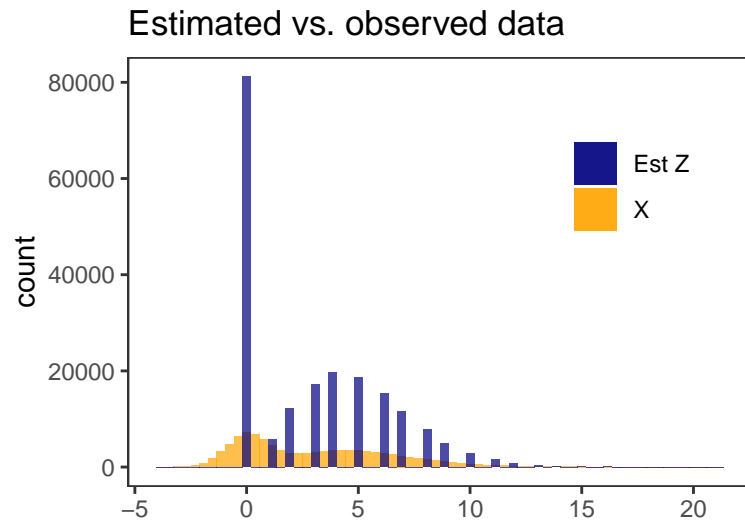
##
## Estimated Raw Moments/Implied Distribution Moments:
##
##      mu1  mu2  mu3  mu4  mu5  mu6
## Normal   1  1.00  1.21  1.29  1.49  1.67
## Poisson   1  1.57  2.44  3.76  5.79  8.93
## NB        1  1.00  0.82  0.57  0.36  0.19
## ZINB       1  1.00  1.01  1.01  1.02  1.03

# Check that moments were estimated with sufficient precision
dist_test$Normal$moment_precision

## [1] 332049.11   356.25   221.65   140.64   89.74   57.28

# So it looks like ZINB fits the data well
# Here we can check this is right by comparing the parameterized distribution to the private data
# Normally we would not have access to the private data and would rely on the moment fit
dist_test$ZINB$plot

```



## Regression diagnostics

The function `diagnosticsDP()` runs regression diagnostics of the kind that would be run after running OLS on private data. It takes a single argument which must be an `lmdp()` object.

```
# First we run a bias corrected regression
lmdp_obj <- lmdp(Y ~ X1, data = dp_data2)

# Now we can run regression diagnostics
diagnostics <- diagnosticsDP(lmdp_obj)
```

```
##
## Heteroskedasticity test via Variance Regression (Bias Corrected Residuals):
##
##      Estimate Std. Error t value Pr(>|t|)
## (Intercept)   9.5007    0.4780  19.8762     0
## X1            2.1201    0.1054  20.1189     0
##
## Error Normality Test:
##
```

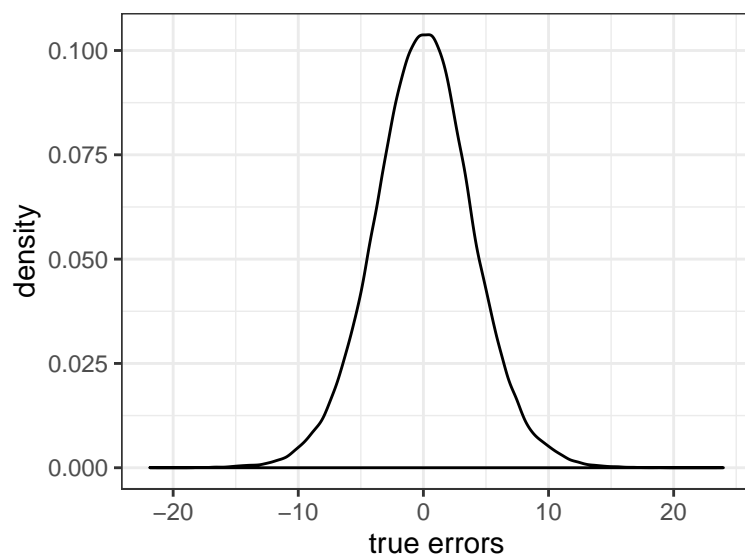
```
## Warning: Error moment estimates are imprecise - unable to accurately test normality of errors
## Skewness 0.07748076      0
## Kurtosis 3.79233169     3

# We see that there is no evidence of non-normality (which is right, true errors were drawn from normal.
# We also detect heteroskedasticity since coefficient on X1 is positive

# To validate, here can check this against the private data

# Calculate the true errors
true_coef <- c(2, 8)
true_errors <- private_data2$Y - cbind(1, private_data2$Z1)%*%true_coef

# Validate they are ormallly distributed
ggplot() + geom_density(aes(x = true_errors)) + theme_bw() + labs(x = 'true errors')
```



```
# Test for heteroskedasticity (we obtain similar coefficients, showing heteroskedasticity)

summary(lm(true_errors^2 ~ private_data2$Z1))$coefficients

##              Estimate Std. Error  t value Pr(>|t|)
## (Intercept)    9.922274  0.10570836  93.86462     0
## private_data2$Z1 2.015053  0.02445956  82.38307     0
```

## References

- Douglas Bates and Martin Maechler. *Matrix: Sparse and Dense Matrix Classes and Methods*, 2019. URL <https://CRAN.R-project.org/package=Matrix>. R package version 1.2-18.
- Georgina Evans and Gary King. Statistically valid inferences from differentially private data releases, Working paper. URL <https://gking.harvard.edu/dpd>.
- Georgina Evans, Gary King, Margaret Schwenzfeier, and Abhradeep Thakurta. Statistically valid inferences from privacy protected data, Working paper. URL <https://gking.harvard.edu/dp>.