

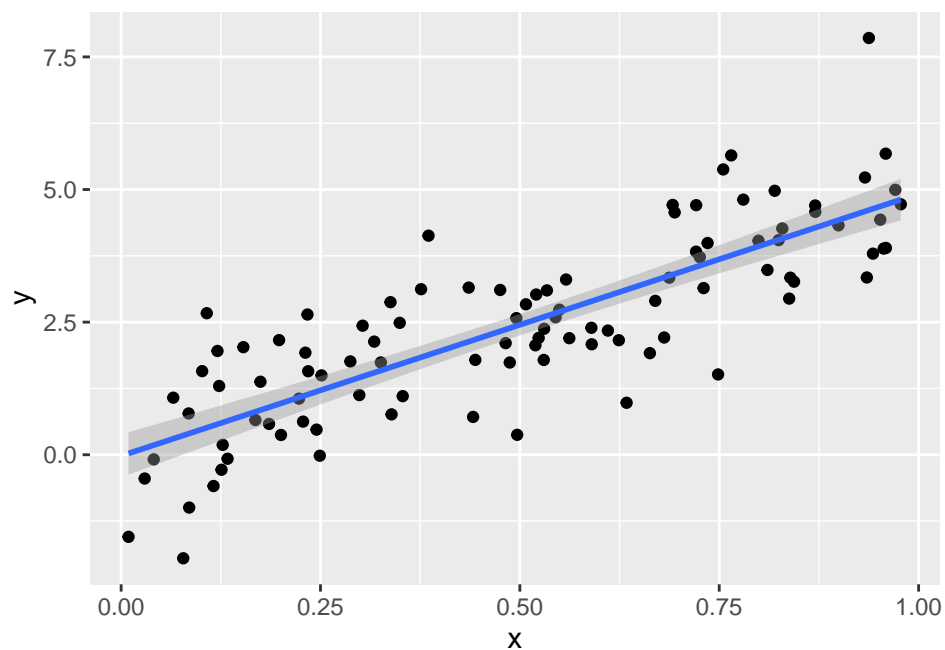
Stan demo

key

Part 1. Linear Regression First consider a simple linear regression model with a single continuous variable.

```
n <- 100
beta <- 5
sigma <- 1
x <- runif(n)
y <- rnorm(n, x * beta, sigma)

lm_dat <- tibble(x = x, y = y)
lm_dat %>% ggplot(aes(y = y, x = x)) + geom_point() + geom_smooth(formula = 'y~x', method = 'lm')
```



We have used both `lm` and `stan_glm` to fit these models.

```
lm_dat %>% lm(y ~ x, data = .) %>% display()
```

```
## lm(formula = y ~ x, data = .)
##           coef.est coef.se
## (Intercept) -0.02    0.20
## x           4.94    0.35
## ---
## n = 100, k = 2
## residual sd = 1.01, R-Squared = 0.67
```

```
lm_dat %>% stan_glm(y ~ x, data = ., refresh = 0) %>% print(digits = 2)
```

```
## stan_glm
## family:      gaussian [identity]
## formula:     y ~ x
## observations: 100
## predictors:  2
## -----
##              Median MAD_SD
## (Intercept) -0.02   0.21
## x           4.94   0.36
##
## Auxiliary parameter(s):
##              Median MAD_SD
## sigma 1.02   0.07
##
## -----
## * For help interpreting the printed output see ?print.stanreg
## * For info on the priors used see ?prior_summary.stanreg
```

Unsurprisingly, `stan_lm` uses `stan`. The code can be extracted, but is not particularly easy to follow. However, we can fairly easily write code Stan code for this model, or obtain it see https://mc-stan.org/docs/2_29/stan-users-guide/linear-regression.html.

```
data {
  int<lower=0> N;
  vector[N] x;
  vector[N] y;
}
parameters {
  real alpha;
  real beta;
  real<lower=0> sigma;
}
model {
  y ~ normal(alpha + beta * x, sigma);
}
```

```
Reg_params <- stan("lm.stan",
  data=list(N = n,
            y = y,
            x = x),
  iter = 2000)
```

```
print(Reg_params, pars = c('alpha', 'beta', 'sigma'))
```

```
## Inference for Stan model: anon_model.
## 4 chains, each with iter=2000; warmup=1000; thin=1;
## post-warmup draws per chain=1000, total post-warmup draws=4000.
##
##      mean se_mean  sd  2.5%  25%   50%  75%  97.5% n_eff Rhat
## alpha -0.03     0.00 0.20 -0.45 -0.16 -0.03  0.11  0.35 1709   1
## beta   4.95     0.01 0.35  4.31  4.72  4.95  5.18  5.64 1791   1
## sigma  1.02     0.00 0.07  0.89  0.97  1.02  1.07  1.18 2134   1
##
## Samples were drawn using NUTS(diag_e) at Fri Apr 22 12:22:09 2022.
```

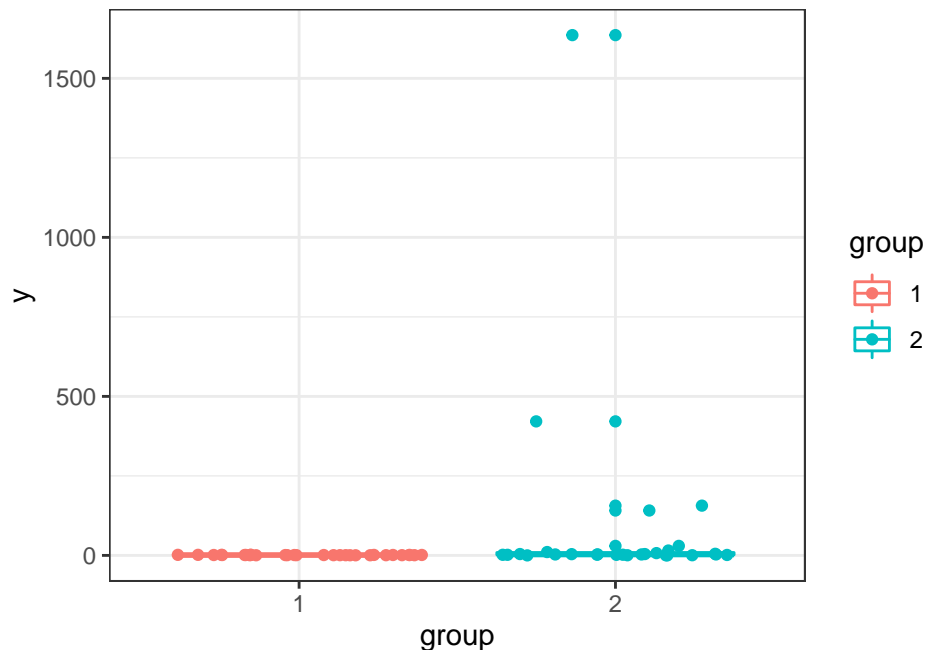
```
## For each parameter, n_eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor on split chains (at
## convergence, Rhat=1).
```

Part 2. Bayesian “t-test” The standard 2-sample t-test, typically has an assumption of constant variance between the groups. However, consider simulated data where the variance terms are different for each group.

```
n1 <- 30
mu1 <- 1
sigma1 <- .5
y1 <- rnorm(n1, mu1, sigma1)

n2 <- 27
mu2 <- .5
sigma2 <- 1.5
y2 <- rt(n2, mu2, sigma2)

tibble(y = c(y1, y2), group = factor(c(rep(1, n1), rep(2, n2)))) %>%
  ggplot(aes(y = y, x = group, color = group)) +
  geom_boxplot() + theme_bw() +
  geom_jitter()
```



The default settings for the `t.test()` function do not account for different variances. There is an option for non-equal variances, but isn't necessarily clear what the procedure does and it doesn't directly return estimated variances.

```
t.test(y1, y2)

##
## Welch Two Sample t-test
##
## data: y1 and y2
## t = -1.4587, df = 26, p-value = 0.1566
```

```
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -216.96522 36.84491
## sample estimates:
## mean of x mean of y
## 1.034236 91.094394
```

Thus, we can easily construct this procedure in Stan. Consider the following code for a two-sample t-test.

```
data {
  int<lower=1> n1; // number of observations in group 1
  vector[n1] y1; // observations from group 1
  int<lower=1> n2; // number of observations in group 1
  vector[n2] y2; // observations from group 1
}
parameters {
  real<lower=0> sigma; // variance parameter
  real<lower=0> mu1; // group 1 mean
  real<lower=0> mu2; // group 2 mean
}
transformed parameters{
  real diff;
  diff = mu1 - mu2;
}
model {
  y1 ~ normal(mu1, sigma);
  y2 ~ normal(mu2, sigma);
}
```

```
ttest_params <- stan("t_test_key.stan",
  data=list(n1 = n1,
            y1 = y1,
            n2 = n2,
            y2 = y2),
  iter = 2000)
```

```
## Running /Library/Frameworks/R.framework/Resources/bin/R CMD SHLIB foo.c
## clang -mmacosx-version-min=10.13 -I"/Library/Frameworks/R.framework/Resources/include" -DNDEBUG -I
## In file included from <built-in>:1:
## In file included from /Library/Frameworks/R.framework/Versions/4.1/Resources/library/StanHeaders/inc
## In file included from /Library/Frameworks/R.framework/Versions/4.1/Resources/library/RcppEigen/inclu
## In file included from /Library/Frameworks/R.framework/Versions/4.1/Resources/library/RcppEigen/inclu
## /Library/Frameworks/R.framework/Versions/4.1/Resources/library/RcppEigen/include/Eigen/src/Core/util
## namespace Eigen {
## ~
## /Library/Frameworks/R.framework/Versions/4.1/Resources/library/RcppEigen/include/Eigen/src/Core/util
## namespace Eigen {
## ~
## ;
## In file included from <built-in>:1:
## In file included from /Library/Frameworks/R.framework/Versions/4.1/Resources/library/StanHeaders/inc
## In file included from /Library/Frameworks/R.framework/Versions/4.1/Resources/library/RcppEigen/inclu
## /Library/Frameworks/R.framework/Versions/4.1/Resources/library/RcppEigen/include/Eigen/Core:96:10: f
## #include <complex>
## ~~~~~
```

```
## 3 errors generated.
## make: *** [foo.o] Error 1

print(ttest_params)

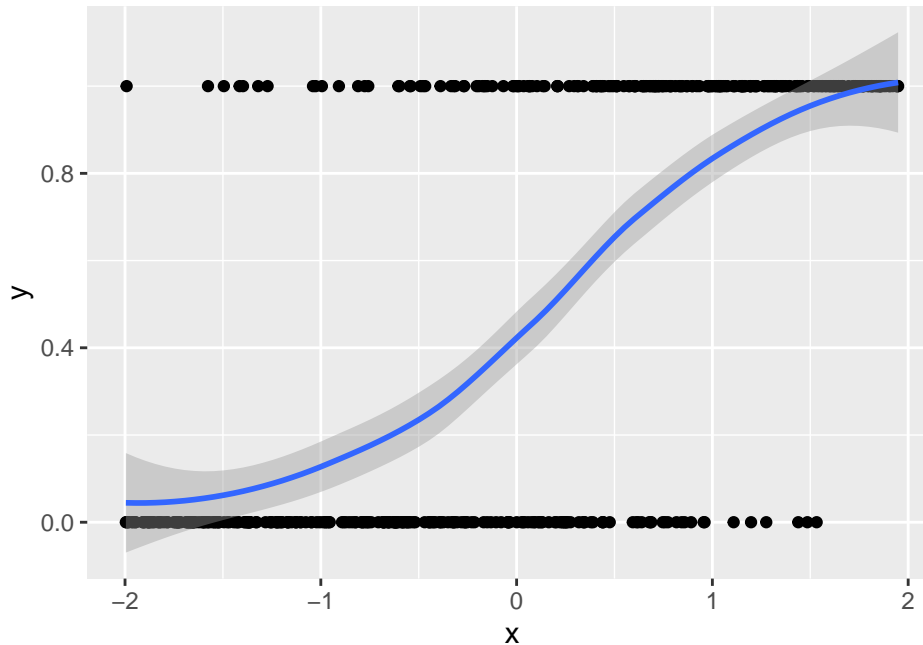
## Inference for Stan model: anon_model.
## 4 chains, each with iter=2000; warmup=1000; thin=1;
## post-warmup draws per chain=1000, total post-warmup draws=4000.
##
##          mean se_mean    sd    2.5%    25%    50%    75%    97.5% n_eff Rhat
## sigma1    0.47     0.00  0.07     0.36     0.42     0.46     0.51     0.62  3103    1
## sigma2  334.64     0.80 48.50   253.27   301.03   329.11   362.96   445.87  3711    1
## mu1        1.04     0.00  0.08     0.87     0.98     1.04     1.09     1.21  3462    1
## mu2       101.33     1.15 55.34    10.04    59.66    96.24   137.77   221.19  2305    1
## diff     -100.29     1.15 55.34  -220.27  -136.71  -95.24   -58.59    -8.98  2305    1
## lp__     -151.98     0.04  1.53  -155.99  -152.74  -151.64  -150.85  -150.07  1258    1
##
## Samples were drawn using NUTS(diag_e) at Fri Apr 22 12:23:20 2022.
## For each parameter, n_eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor on split chains (at
## convergence, Rhat=1).
```

Q1: Update the code to allow for different variances between the two groups.

Part 3. Logistic Regression Similarly for a logistic regression model with a single continuous variable.

```
n <- 500
beta <- 2
x <- runif(n, -2, 2)
p <- invlogit(x * beta)
y <- rbinom(n, 1, p)

logistic_dat <- tibble(x = x, y = y)
logistic_dat %>%
  ggplot(aes(y = y, x = x)) +
  geom_point() +
  geom_smooth(formula = 'y~x', method = 'loess')
```



We have used both `glm` and `stan_glm` to fit these models.

```
logistic_dat %>% glm(y ~ x, data = ., family = binomial) %>% display()
```

```
## glm(formula = y ~ x, family = binomial, data = .)
##      coef.est coef.se
## (Intercept) -0.17   0.13
## x           1.88   0.16
## ---
##  n = 500, k = 2
##  residual deviance = 402.4, null deviance = 692.9 (difference = 290.6)
```

```
logistic_dat %>% stan_glm(y ~ x, data = ., family = binomial, refresh = 0) %>% print(digits = 2)
```

```
## stan_glm
## family:      binomial [logit]
## formula:     y ~ x
## observations: 500
## predictors:  2
## -----
##              Median MAD_SD
## (Intercept) -0.18   0.13
## x           1.88   0.16
##
## -----
## * For help interpreting the printed output see ?print.stanreg
## * For info on the priors used see ?prior_summary.stanreg
```

Similarly Stan code can be used to fit a logistic regression model.

```
data {
  int <lower = 0> N;
  int <lower = 0, upper = 1> y [N];
  vector [N] x;
}
```

```

parameters {
  real alpha;
  real beta;
}

model {
  y ~ bernoulli_logit(alpha + beta * x);

  // alpha ~ normal(0, 1);
  // beta ~ normal(1, 1);
}

log_params <- stan("logistic.stan",
  data=list(N = n,
            y = y,
            x = x),
  iter = 2000)

print(log_params, pars = c('alpha', 'beta'))

## Inference for Stan model: anon_model.
## 4 chains, each with iter=2000; warmup=1000; thin=1;
## post-warmup draws per chain=1000, total post-warmup draws=4000.
##
##      mean se_mean   sd  2.5%  25%   50%   75% 97.5% n_eff Rhat
## alpha -0.17      0 0.13 -0.43 -0.26 -0.17 -0.09  0.07  2472    1
## beta   1.90      0 0.16  1.60  1.79  1.89  2.00  2.22  2536    1
##
## Samples were drawn using NUTS(diag_e) at Fri Apr 22 12:24:37 2022.
## For each parameter, n_eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor on split chains (at
## convergence, Rhat=1).

```

Q2: Update the STAN code to place priors on alpha and beta.

Part 4. Hierarchical Logistic Regression The stan reference book contains code for a hierarchical logistic regression model https://mc-stan.org/docs/2_29/stan-users-guide/hierarchical-logistic-regression.html.

```

data {
  int<lower = 0> K;
  int<lower = 0> N;
  int<lower = 1, upper = K> kk[N];
  vector[N] x;
  int<lower = 0, upper = 1> y[N];
}

parameters {
  matrix[K,2] beta;
  vector[2] mu;
  vector<lower=0>[2] sigma;
}

model {
  mu ~ normal(0, 2);
  sigma ~ normal(0, 2);
  for (i in 1:2)
    beta[, i] ~ normal(mu[i], sigma[i]);
}

```

```

    for (n in 1:N){
      y[n] ~ bernoulli_logit(beta[kk[n], 1] + beta[kk[n], 2] * x[n]);
    }
  }
}

```

Q3: Simulate hierarchical logistic regression data

```

N <- 500
K <- 10
kk <- rep(1:K, each = N/K)
beta0 <- rnorm(K, sd = .5)
beta1 <- rnorm(K, sd = .5)

x <- runif(N)
p <- invlogit(rep(beta0, each = N/K) + rep(beta1, each = N/K) * x)
y <- rbinom(N, 1, p)

```

```

hl_params <- stan("hier_logistic.stan",
  data=list(N = N,
    y = y,
    x = x,
    kk = kk,
    K = K),
  iter = 10000)

```

```
print(hl_params)
```

```

## Inference for Stan model: anon_model.
## 4 chains, each with iter=10000; warmup=5000; thin=1;
## post-warmup draws per chain=5000, total post-warmup draws=20000.
##
##               mean se_mean  sd   2.5%   25%   50%   75%  97.5% n_eff
## beta[1,1]    -0.09   0.00 0.37  -0.85  -0.32  -0.07   0.15   0.62  8349
## beta[1,2]    -0.78   0.01 0.71  -2.22  -1.24  -0.75  -0.28   0.53  6333
## beta[2,1]    -0.26   0.00 0.36  -1.03  -0.48  -0.23  -0.01   0.39  7627
## beta[2,2]     0.24   0.01 0.61  -0.96  -0.16   0.23   0.64   1.49  8674
## beta[3,1]    -0.26   0.00 0.35  -1.00  -0.49  -0.24  -0.02   0.38  6673
## beta[3,2]    -0.23   0.01 0.64  -1.53  -0.65  -0.21   0.20   0.99  8004
## beta[4,1]    -0.12   0.00 0.33  -0.81  -0.32  -0.10   0.10   0.52  9947
## beta[4,2]     0.52   0.01 0.62  -0.64   0.11   0.49   0.92   1.81  9053
## beta[5,1]    -0.18   0.00 0.37  -0.99  -0.41  -0.16   0.07   0.49  8852
## beta[5,2]     0.30   0.01 0.58  -0.81  -0.08   0.28   0.67   1.51  9153
## beta[6,1]     0.38   0.01 0.39  -0.30   0.10   0.35   0.64   1.21  4545
## beta[6,2]     0.59   0.01 0.71  -0.77   0.11   0.57   1.05   2.06  6082
## beta[7,1]     0.46   0.01 0.47  -0.33   0.11   0.41   0.78   1.46  3786
## beta[7,2]     2.02   0.02 1.22  -0.01   1.11   1.96   2.82   4.64  3846
## beta[8,1]     0.10   0.00 0.35  -0.57  -0.12   0.09   0.33   0.82  8528
## beta[8,2]     0.66   0.01 0.64  -0.55   0.23   0.64   1.07   2.00  8132
## beta[9,1]    -0.30   0.00 0.37  -1.08  -0.54  -0.27  -0.03   0.36  5936
## beta[9,2]    -0.31   0.01 0.65  -1.61  -0.75  -0.30   0.14   0.94  6825
## beta[10,1]     0.09   0.00 0.33  -0.54  -0.12   0.08   0.30   0.79  9979
## beta[10,2]   -0.15   0.01 0.62  -1.44  -0.55  -0.13   0.26   1.01  9888
## mu[1]        -0.02   0.00 0.24  -0.50  -0.17  -0.02   0.14   0.47  7562
## mu[2]         0.28   0.01 0.48  -0.65  -0.03   0.27   0.58   1.27  5822
## sigma[1]      0.46   0.00 0.26   0.08   0.27   0.43   0.61   1.04  2795
## sigma[2]      1.07   0.01 0.51   0.23   0.71   1.01   1.36   2.26  4112

```



```

## lp__      -327.28    0.14 5.55 -338.17 -330.89 -327.37 -323.81 -315.65 1648
##          Rhat
## beta[1,1] 1.00
## beta[1,2] 1.00
## beta[2,1] 1.00
## beta[2,2] 1.00
## beta[3,1] 1.00
## beta[3,2] 1.00
## beta[4,1] 1.00
## beta[4,2] 1.00
## beta[5,1] 1.00
## beta[5,2] 1.00
## beta[6,1] 1.00
## beta[6,2] 1.00
## beta[7,1] 1.00
## beta[7,2] 1.00
## beta[8,1] 1.00
## beta[8,2] 1.00
## beta[9,1] 1.00
## beta[9,2] 1.00
## beta[10,1] 1.00
## beta[10,2] 1.00
## mu[1]     1.00
## mu[2]     1.00
## sigma[1]  1.00
## sigma[2]  1.00
## lp__      1.01
##
## Samples were drawn using NUTS(diag_e) at Fri Apr 22 12:27:56 2022.
## For each parameter, n_eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor on split chains (at
## convergence, Rhat=1).

```