# 431 Class 11

Thomas E. Love, Ph.D.

2022-10-04

431

# Today's Agenda

- Building models for `sbp_2` using `sbp_1` and `insur_1`

  - without an interaction term

  - with an interaction

- Comparing our four models (two from Class 10 and two today)

Version 2022-09-25 23:07:19

431 CASE WESTERN RESERVE UNIVERSITY

# Today's R packages

```
 1  library(broom)
 2  library(equatiomatic)
 3  library(haven)           ## import SPSS .sav file
 4  library(rstanarm)        ## fit stan_glm() model
 5  library(janitor)
 6  library(kableExtra)
 7  library(naniar)
 8  library(patchwork)
 9  library(tidyverse)
10
11  theme_set(theme_bw())
```

431

# Today's Data

Today's data describe 1,500 adults with hypertension living in Cuyahoga County, whose (systolic) blood pressure was measured at baseline, and then again one year later. We also have information on (baseline) primary insurance, and other things.

- We created and partitioned the data back in Class 10

```
1  bp_full <- read_rds("c11/data/bp_full.Rds")
2  bp_train <- read_rds("c11/data/bp_train.Rds")
3  bp_test <- read_rds("c11/data/bp_test.Rds")
```

431 CASE WESTERN RESERVE UNIVERSITY

# Research Questions

1. Can we build an effective model to predict sbp_2 (SBP after a year) using sbp_1 (SBP at baseline)? (addressed in class 10)

2. Is the effectiveness of such a model for prediction of sbp_2 materially affected by whether we also include information about ins_1 (Primary insurance at baseline)? (today)

# Modeling Goals

## Class 10

- Model `sbp_2` on the basis of `sbp_1`
    - using a linear regression model
    - using a (naive) Bayesian model

## Today

- Model `sbp_2` using `sbp_1` and `ins_1`
    - without an interaction term
    - including an `sbp_1*ins_1` interaction term

Build models with **training** sample, evaluate performance in **testing** sample.

431 CASE WESTERN RESERVE UNIVERSITY

# Previous models (m1 and m2)

Fit in training sample, then evaluate in testing sample.

```
1  m1_train <- lm(sbp_2 ~ sbp_1, data = bp_train)
2  m1_test_aug <- augment(m1_train, newdata = bp_test)
3
4  m2_train <- stan_glm(sbp_2 ~ sbp_1, data = bp_train, refresh = 0)
5  m2_test_aug <- bp_test |> select(record, sbp_2, sbp_1) |>
6    mutate(.fitted = predict(m2_train, newdata = bp_test),
7           .resid = sbp_2 - .fitted)
```

431 CASE WESTERN RESERVE UNIVERSITY

# Which priors did we use in `m2_train`?

For more, visit https://mc-stan.org/rstanarm/articles/priors.html.

```
1  prior_summary(m2_train)
```

```
Priors for model 'm2_train'
------
Intercept (after predictors centered)
  Specified prior:
    ~ normal(location = 132, scale = 2.5)
  Adjusted prior:
    ~ normal(location = 132, scale = 41)

Coefficients
  Specified prior:
    ~ normal(location = 0, scale = 2.5)
  Adjusted prior:
    ~ normal(location = 0, scale = 2.4)

Auxiliary (sigma)
  Specified prior:
```
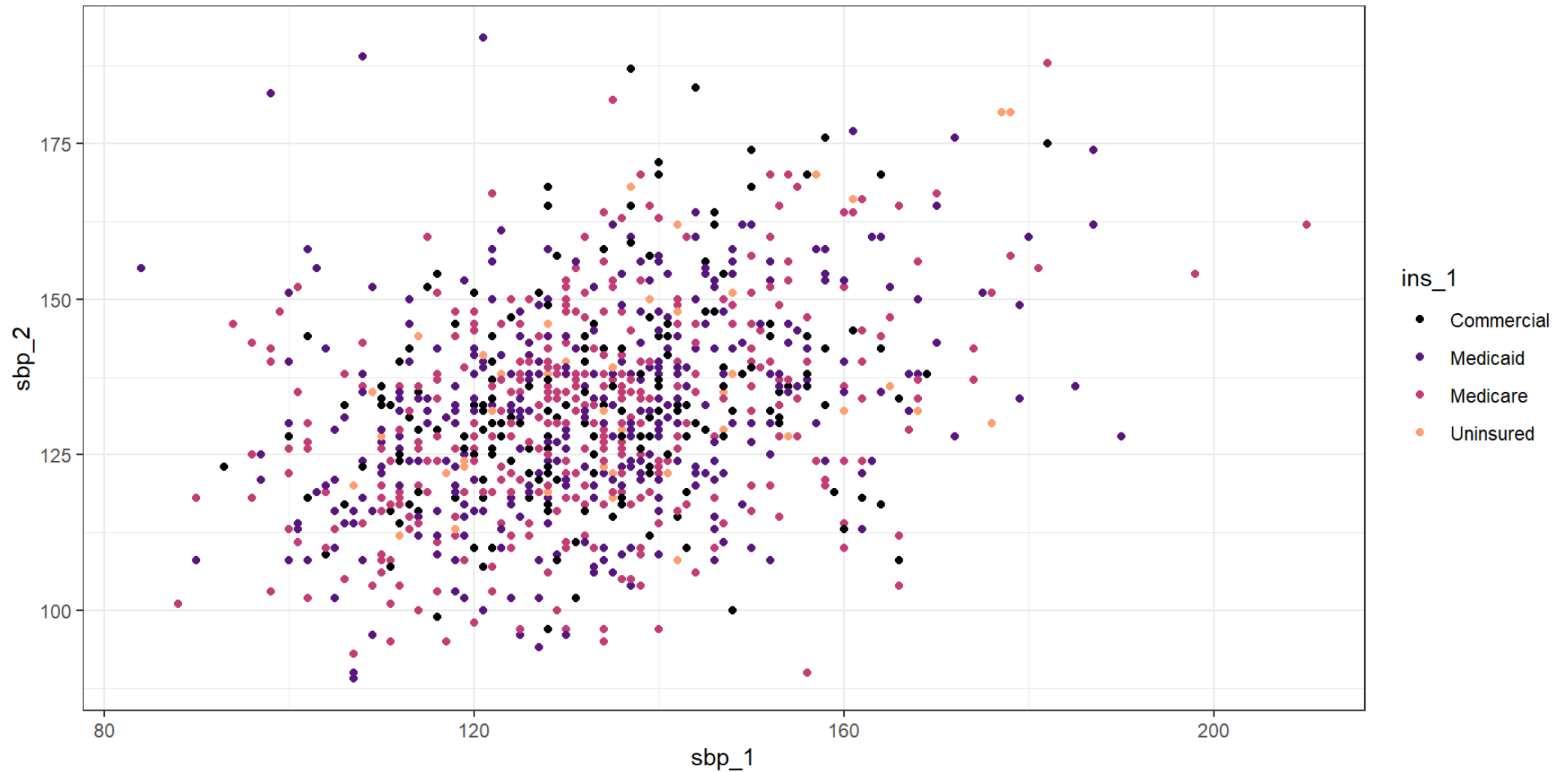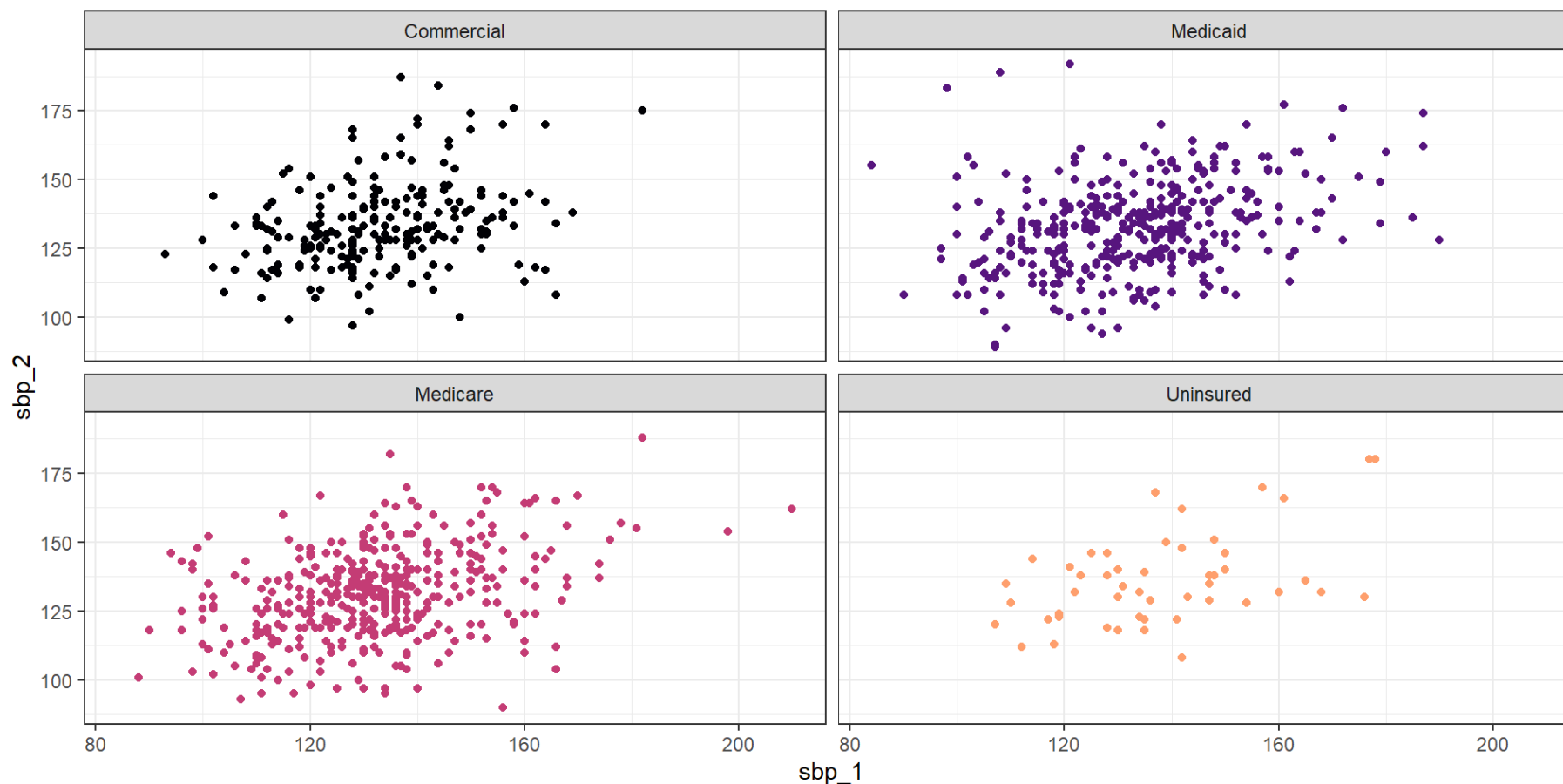
431 CASE WESTERN RESERVE UNIVERSITY

# Add in `ins_1` information

# Faceting by `ins_1` group

```
1  ggplot(data = bp_train, aes(x = sbp_1, y = sbp_2, col = ins_1)) +
2    geom_point() + scale_color_viridis_d(option = "A", end = 0.8) +
3    facet_wrap(~ ins_1) + guides(col = "none")
```

# Two possible models

```
1  m3_train <- lm(sbp_2 ~ sbp_1 + ins_1, data = bp_train)
2  m4_train <- lm(sbp_2 ~ sbp_1 * ins_1, data = bp_train)
```

- What is the difference between m3 and m4?

- Model m3 does not include an interaction term, while m4 does.

- How does this work in practice?

431 CASE WESTERN RESERVE UNIVERSITY

# Equation for m3

```
m3_train <- lm(sbp_2 ~ sbp_1 + ins_1, data = bp_train)
```

```
1  extract_eq(m3_train, use_coefs = TRUE, operator_location = "start",
2             wrap = TRUE, terms_per_line = 2, coef_digits = 2)
```

$$\widehat{\text{sbp\_2}} = 89.36 + 0.33(\text{sbp\_1})$$
$$- 0.83(\text{ins\_1}_{\text{Medicaid}}) - 2.41(\text{ins\_1}_{\text{Medicare}})$$
$$+ 1.38(\text{ins\_1}_{\text{Uninsured}})$$

In model m3, the intercept term of the sbp_1-sbp_2 relationship varies depending on insurance.

# Model **m3** by Insurance Type

$$\widehat{\text{sbp\_2}} = 89.36 + 0.33(\text{sbp\_1}) - 0.83(\text{ins\_1}_{\text{Medicaid}}) - 2.41(\text{ins\_1}_{\text{Medicare}})$$
$$+ 1.38(\text{ins\_1}_{\text{Uninsured}})$$

| Insurance | Estimated **sbp_2** |
|:---:|:---:|
| Commmercial | 89.36 + 0.33 sbp_1 |
| Medicaid | ?? |
| Medicare | ?? |
| Uninsured | ?? |

# Model m3 by Insurance Type

$$\widehat{\text{sbp\_2}} = 89.36 + 0.33(\text{sbp\_1}) - 0.83(\text{ins\_1}_{\text{Medicaid}}) - 2.41(\text{ins\_1}_{\text{Medicare}})$$
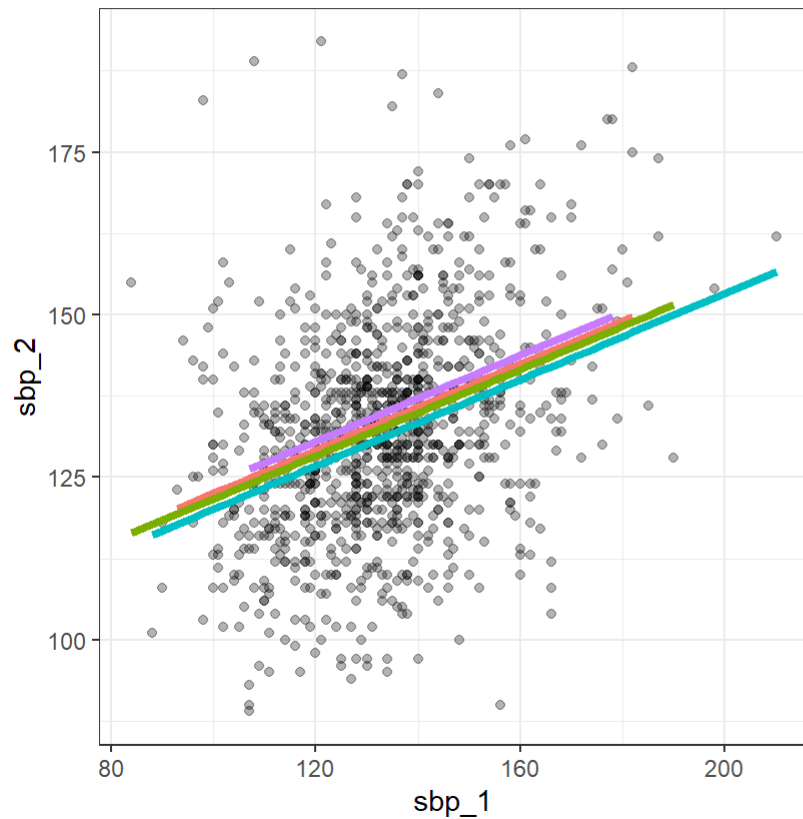$$+ 1.38(\text{ins\_1}_{\text{Uninsured}})$$

| Insurance | Estimated sbp_2 |
|---|---|
| Commmercial | 89.36 + 0.33 sbp_1 |
| Medicaid | (89.36 - 0.83) + 0.33 sbp_1 <br> = **88.53** + 0.33 sbp_1 |
| Medicare | (89.36 - 2.41) + 0.33 sbp_1 <br> = **86.95** + 0.33 sbp_1 |
| Uninsured | (89.36 + 1.38) + 0.33 sbp_1 <br> = **90.74** + 0.33 sbp_1 |

431 CASE WESTERN RESERVE UNIVERSITY

# The m3 model (pictured)

```
1  m3_train_aug <- augment(m3_train, data = bp_train)
2
3  p1 <- ggplot(m3_train_aug, aes(x = sbp_1, y = sbp_2, group = ins_1)) +
4    geom_point(alpha = 0.3) +
5    geom_line(aes(x = sbp_1, y = .fitted, col = ins_1), lwd = 1.5) +
6    labs(title = "m3: Same Slope, Intercepts vary by insurance")
7
8  p2 <- ggplot(m3_train_aug, aes(x = sbp_1, y = sbp_2,
9                                 col = ins_1, group = ins_1)) +
10   geom_point() + geom_line(aes(x = sbp_1, y = .fitted), col = "black") +
11   facet_wrap( ~ ins_1) + guides(col = "none")
12
13 p1 + p2
```
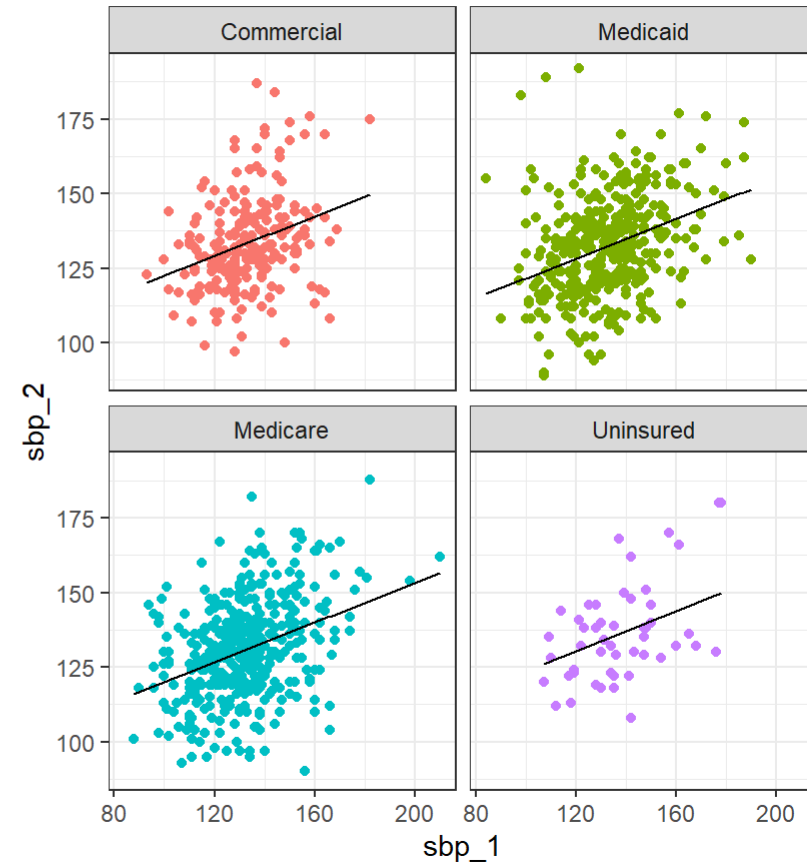
431 CASE WESTERN RESERVE UNIVERSITY

# The m3 model (pictured)

m3: Same Slope, Intercepts vary by insurance

# Tidied Model m3 coefficients

Again, in model m3, only the intercept of the sbp_1 to sbp_2 model varies depending on the ins_1 category.

```
1  tidy(m3_train, conf.int = TRUE, conf.level = 0.90) |>
2    select(term, estimate, std.error, conf.low, conf.high) |>
3    kbl(digits = c(0, 2, 2, 2, 2)) |> kable_styling(font_size = 28)
```

| term | estimate | std.error | conf.low | conf.high |
|------|----------|-----------|----------|-----------|
| (Intercept) | 89.36 | 3.82 | 83.07 | 95.64 |
| sbp_1 | 0.33 | 0.03 | 0.29 | 0.38 |
| ins_1Medicaid | -0.83 | 1.28 | -2.95 | 1.28 |
| ins_1Medicare | -2.41 | 1.28 | -4.51 | -0.30 |
| ins_1Uninsured | 1.38 | 2.41 | -2.59 | 5.34 |

431 CASE WESTERN RESERVE UNIVERSITY

# Equation for m4

```
1  extract_eq(m4_train, use_coefs = TRUE, operator_location = "start", wrap =
2           terms_per_line = 1, coef_digits = 2, font_size = "small")
```

$$\widehat{\text{sbp\_2}} = 90.33$$
$$+ 0.32(\text{sbp\_1})$$
$$+ 1.56(\text{ins\_1}_{\text{Medicaid}})$$
$$- 4.86(\text{ins\_1}_{\text{Medicare}})$$
$$- 16.75(\text{ins\_1}_{\text{Uninsured}})$$
$$- 0.02(\text{sbp\_1} \times \text{ins\_1}_{\text{Medicaid}})$$
$$+ 0.02(\text{sbp\_1} \times \text{ins\_1}_{\text{Medicare}})$$
$$+ 0.13(\text{sbp\_1} \times \text{ins\_1}_{\text{Uninsured}})$$

431

# Model m4 by Insurance Type

$$\widehat{\text{sbp\_2}} = 90.33 + 0.32(\text{sbp\_1}) + 1.56(\text{ins\_1}_{\text{Medicaid}})$$
$$- 4.86(\text{ins\_1}_{\text{Medicare}}) - 16.75(\text{ins\_1}_{\text{Uninsured}}) - 0.02(\text{sbp\_1} \times \text{ins\_1}_{\text{Medicaid}})$$
$$+ 0.02(\text{sbp\_1} \times \text{ins\_1}_{\text{Medicare}}) + 0.13(\text{sbp\_1} \times \text{ins\_1}_{\text{Uninsured}})$$

| Insurance | Estimated sbp_2 |
|---|---|
| Commmercial | 90.33 + 0.32 sbp_1 |
| Medicaid | ?? |
| Medicare | ?? |
| Uninsured | ?? |

# Model m4 by Insurance Type

$$\widehat{\text{sbp\_2}} = 90.33 + 0.32(\text{sbp\_1}) + 1.56(\text{ins\_1}_{\text{Medicaid}})$$
$$- 4.86(\text{ins\_1}_{\text{Medicare}}) - 16.75(\text{ins\_1}_{\text{Uninsured}}) - 0.02(\text{sbp\_1} \times \text{ins\_1}_{\text{Medicaid}})$$
$$+ 0.02(\text{sbp\_1} \times \text{ins\_1}_{\text{Medicare}}) + 0.13(\text{sbp\_1} \times \text{ins\_1}_{\text{Uninsured}})$$

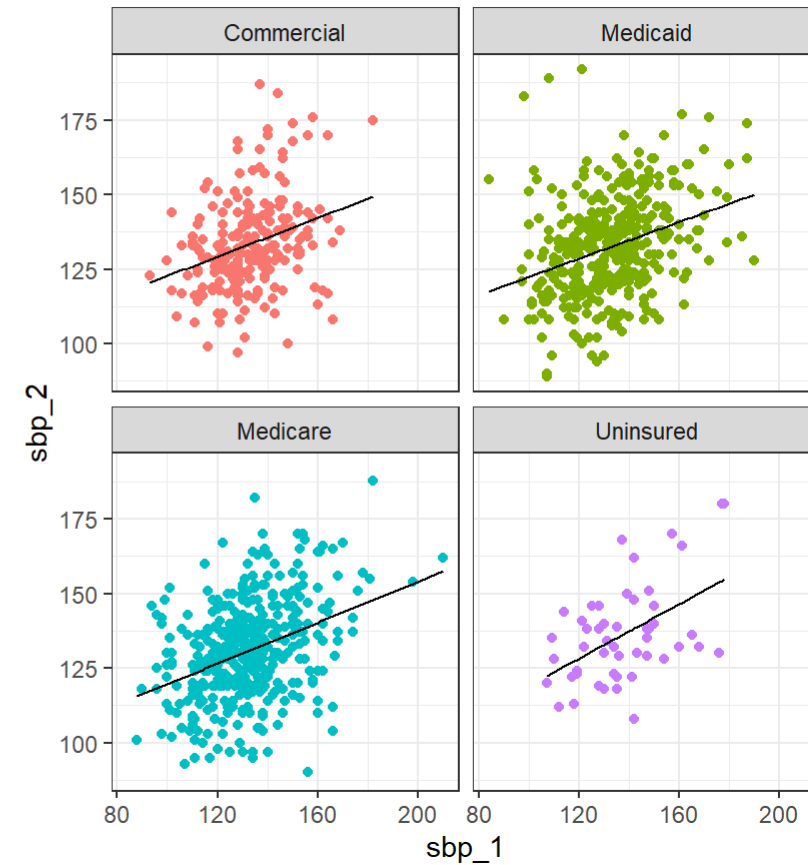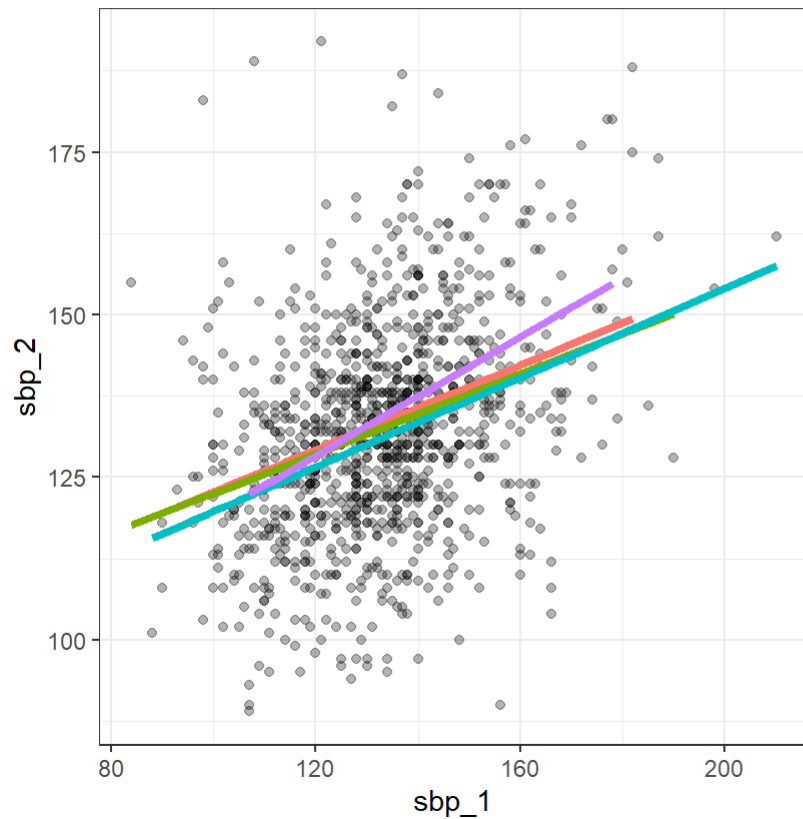| Insurance | Estimated sbp_2 |
|---|---|
| Commmercial | 90.33 + 0.32 sbp_1 |
| Medicaid | (90.33 + 1.56) + (0.32 - 0.02) sbp_1 <br> = **91.89** + **0.30** sbp_1 |
| Medicare | (90.33 - 4.86) + (0.32 + 0.02) sbp_1 <br> = **85.47** + **0.34** sbp_1 |
| Uninsured | (90.33 - 16.75) + (0.32 + 0.13) sbp_1 <br> = **73.58** + **0.45** sbp_1 |

# The **m4** model (pictured)

```
 1  m4_train_aug <- augment(m4_train, data = bp_train)
 2
 3  p1 <- ggplot(m4_train_aug, aes(x = sbp_1, y = sbp_2, group = ins_1)) +
 4    geom_point(alpha = 0.3) +
 5    geom_line(aes(x = sbp_1, y = .fitted, col = ins_1), lwd = 1.5) +
 6    labs(title = "m4: Slopes and Intercepts vary by insurance")
 7
 8  p2 <- ggplot(m4_train_aug, aes(x = sbp_1, y = sbp_2,
 9                          col = ins_1, group = ins_1)) +
10    geom_point() + geom_line(aes(x = sbp_1, y = .fitted), col = "black") +
11    facet_wrap( ~ ins_1) + guides(col = "none")
12
13  p1 + p2
```

# The m4 model (pictured)

m4: Slopes and Intercepts vary by insurance

# Models m3 and m4

```
1  m3_train <- lm(sbp_2 ~ sbp_1 + ins_1, data = bp_train)
2  m4_train <- lm(sbp_2 ~ sbp_1 * ins_1, data = bp_train)
```

- What is the difference between m3 and m4?

  - Model m3 will allow **only the intercept** term of the sbp_1-sbp_2 relationship to vary depending on insurance.

  - Model m4 will allow **both the slope and intercept** of the sbp_1-sbp_2 relationship to vary depending on insurance.

# Tidied Model m4 coefficients

```
1  tidy(m4_train, conf.int = TRUE, conf.level = 0.90) |>
2    select(term, estimate, std.error, conf.low, conf.high) |>
3    kbl(digits = c(0, 2, 2, 2, 2)) |> kable_styling(font_size = 24)
```

| term | estimate | std.error | conf.low | conf.high |
|---|---|---|---|---|
| (Intercept) | 90.33 | 9.30 | 75.02 | 105.64 |
| sbp_1 | 0.32 | 0.07 | 0.21 | 0.44 |
| ins_1Medicaid | 1.56 | 11.05 | -16.63 | 19.74 |
| ins_1Medicare | -4.86 | 10.96 | -22.90 | 13.18 |
| ins_1Uninsured | -16.75 | 19.35 | -48.61 | 15.10 |
| sbp_1:ins_1Medicaid | -0.02 | 0.08 | -0.15 | 0.12 |
| sbp_1:ins_1Medicare | 0.02 | 0.08 | -0.12 | 0.15 |
| sbp_1:ins_1Uninsured | 0.13 | 0.14 | -0.10 | 0.36 |

# Fit within the Training Sample

## Model m3 (no interaction)

```
1  glance(m3_train) |> select(r.squared, sigma, AIC, df, df.residual, nobs) |>
2    kbl(digits = c(3, 1, 1, 0, 0, 0)) |> kable_styling(font_size = 32)
```

| r.squared | sigma | AIC | df | df.residual | nobs |
|-----------|-------|------|-----|-------------|------|
| 0.128 | 15.2 | 8707 | 4 | 1045 | 1050 |

## Model m4 (with sbp_1-insurance interaction)

```
1  glance(m4_train) |> select(r.squared, sigma, AIC, df, df.residual, nobs) |>
2    kbl(digits = c(3, 1, 1, 0, 0, 0)) |> kable_styling(font_size = 32)
```

| r.squared | sigma | AIC | df | df.residual | nobs |
|-----------|-------|--------|-----|-------------|------|
| 0.129 | 15.3 | 8711.5 | 7 | 1042 | 1050 |

# Augmenting and Testing Models m3 and m4

```
1   ## in the training sample (for residual plots)
2
3   m3_train_aug <- augment(m3_train, data = bp_train)
4   m4_train_aug <- augment(m4_train, data = bp_train)
5
6   # in the test sample (calculating prediction errors)
7
8   m3_test_aug <- augment(m3_train, newdata = bp_test)
9   m4_test_aug <- augment(m4_train, newdata = bp_test)
```
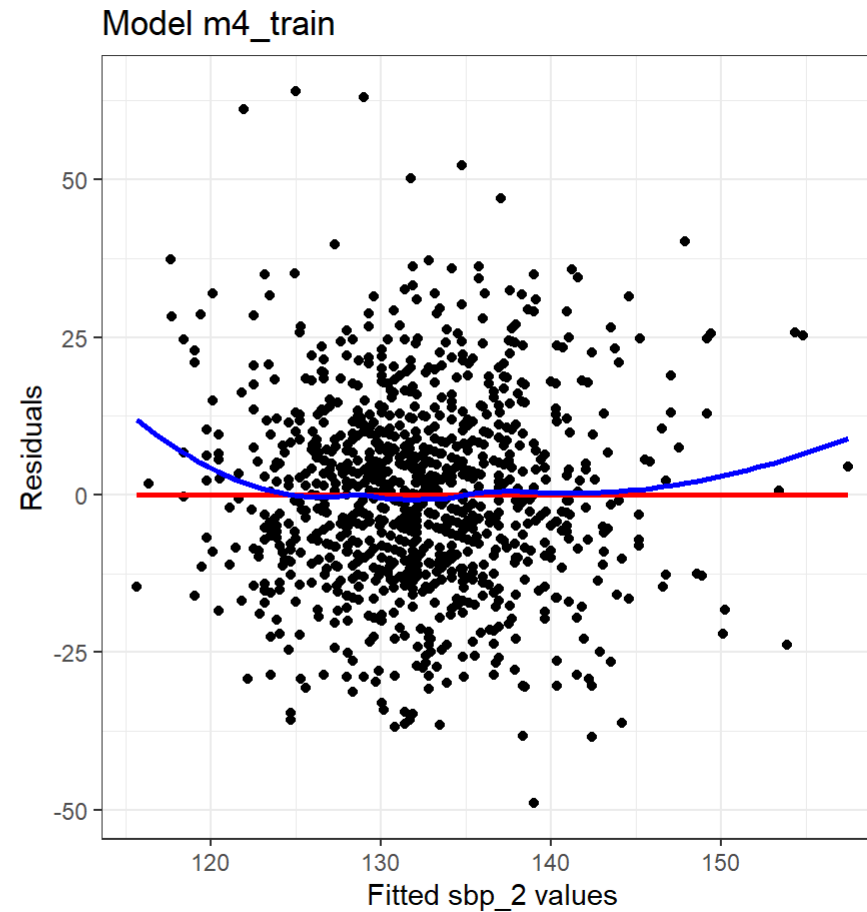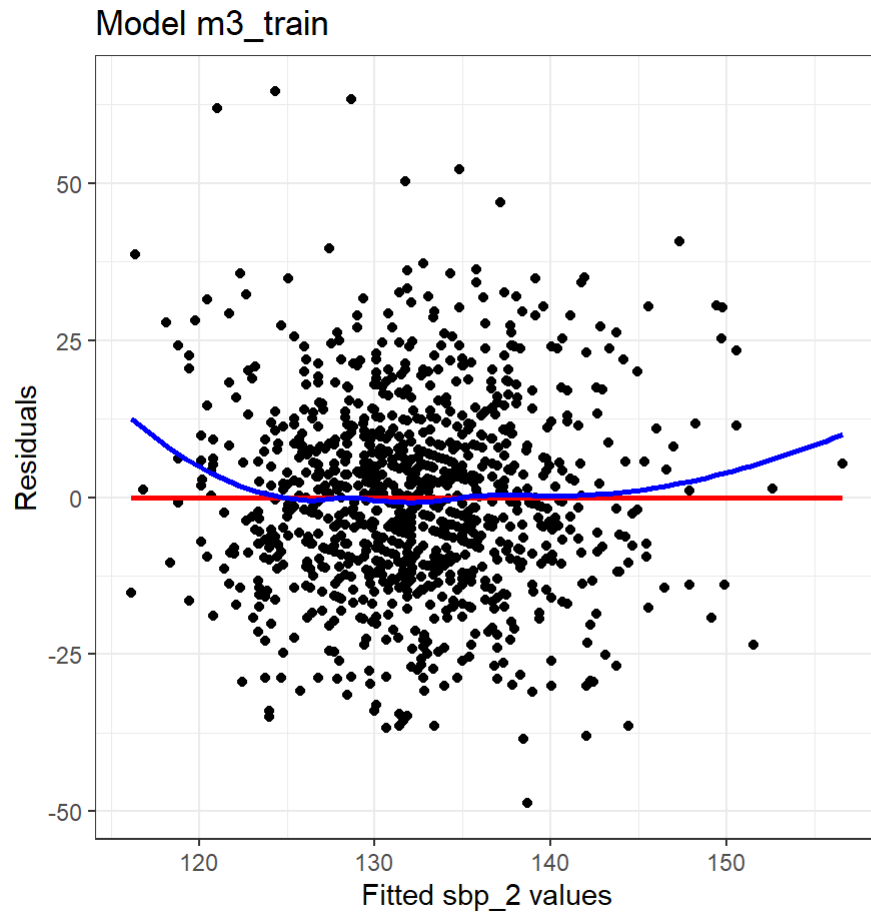
431 CASE WESTERN RESERVE UNIVERSITY

# Residuals vs. Fitted Values Plots

```
1   p1 <- ggplot(m3_train_aug, aes(x = .fitted, y = .resid)) +
2     geom_point() +
3     geom_smooth(method = "lm", col = "red",
4                 formula = y ~ x, se = FALSE) +
5     geom_smooth(method = "loess", col = "blue",
6                 formula = y ~ x, se = FALSE) +
7     theme(aspect.ratio = 1) +
8     labs(title = "Model m3_train",
9          x = "Fitted sbp_2 values", y = "Residuals")
10
11  p2 <- ggplot(m4_train_aug, aes(x = .fitted, y = .resid)) +
12    geom_point() +
13    geom_smooth(method = "lm", col = "red",
14                formula = y ~ x, se = FALSE) +
15    geom_smooth(method = "loess", col = "blue",
16                formula = y ~ x, se = FALSE) +
17    theme(aspect.ratio = 1) +
18    labs(title = "Model m4_train",
19         x = "Fitted sbp_2 values", y = "Residuals")
```

# Residuals vs. Fitted Values Plots

# m3 and m4: Same predictions?

```r
1   t1 <- bind_cols(m3_train_aug$record, m3_train_aug$ins_1, m3_train_aug$.fitt
2                   m4_train_aug$.fitted)
3
4   names(t1) <- c("record", "ins_1", "m3_fit", "m4_fit")
5
6   p1 <- ggplot(data = t1, aes(x = m3_fit, y = m4_fit)) +
7     geom_abline(aes(col = "black"), intercept = 0, slope = 1) +
8     geom_point(size = 2) +
9     theme(aspect.ratio = 1) +
10    labs(title = "Figure 1. Predicted sbp_2 from m3, m4")
11
12  p2 <- ggplot(data = t1, aes(x = m3_fit, y = m4_fit, col = ins_1)) +
13    geom_abline(aes(col = "black"), intercept = 0, slope = 1) +
14    geom_point(size = 2) +
15    theme(aspect.ratio = 1) +
16    facet_wrap( ~ ins_1) +
17    guides(col = "none") +
18    labs(title = "Figure 2. Predicted sbp_2 by ins_1")
19
```

# m3 and m4: Same predictions?



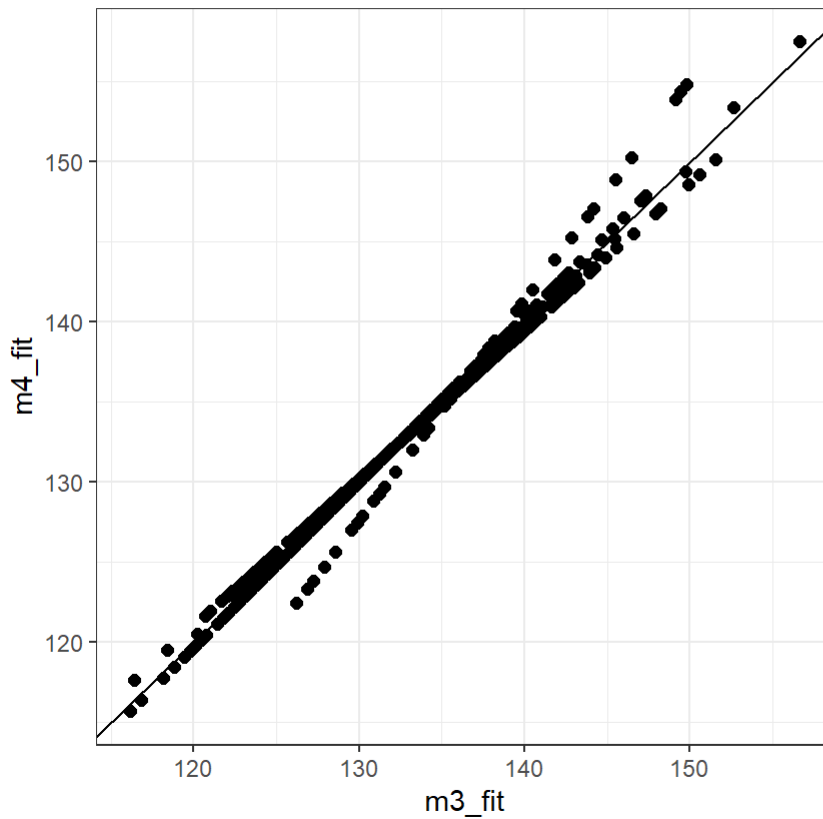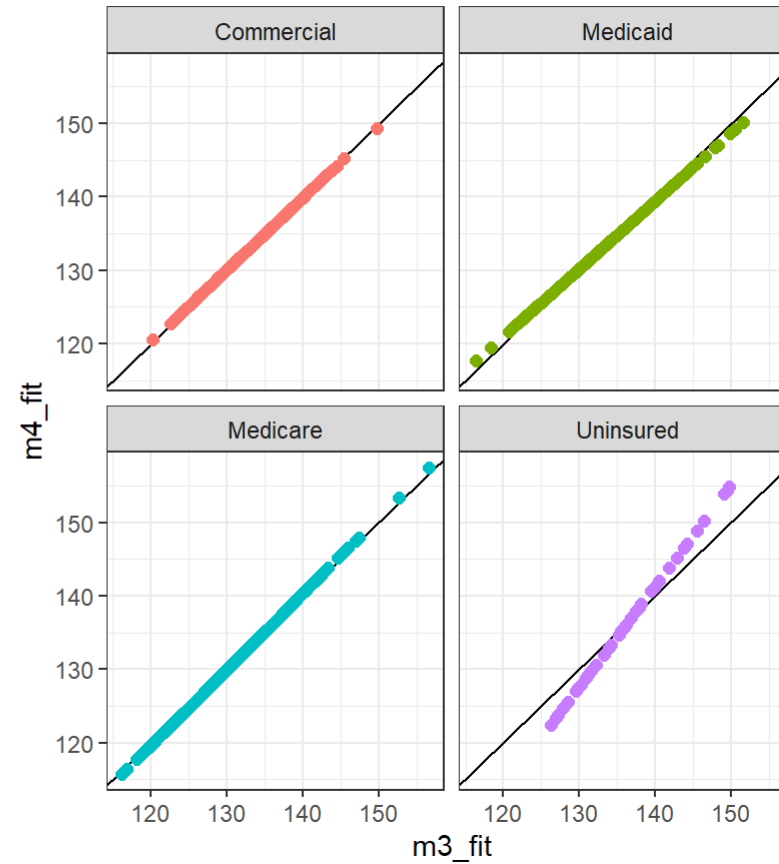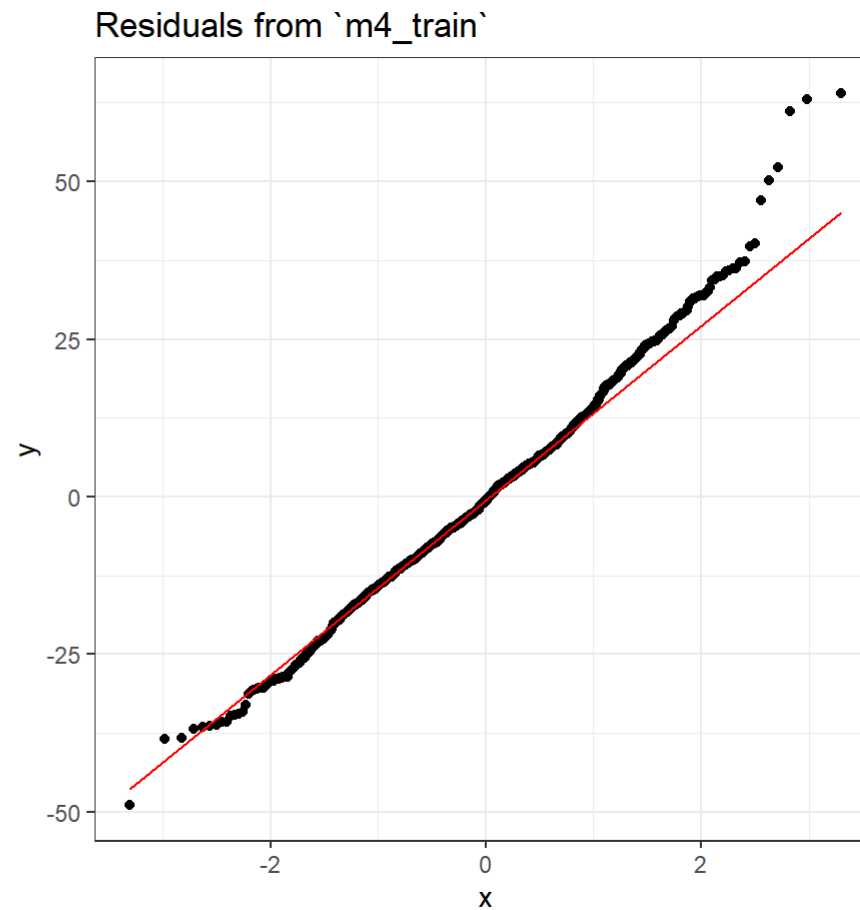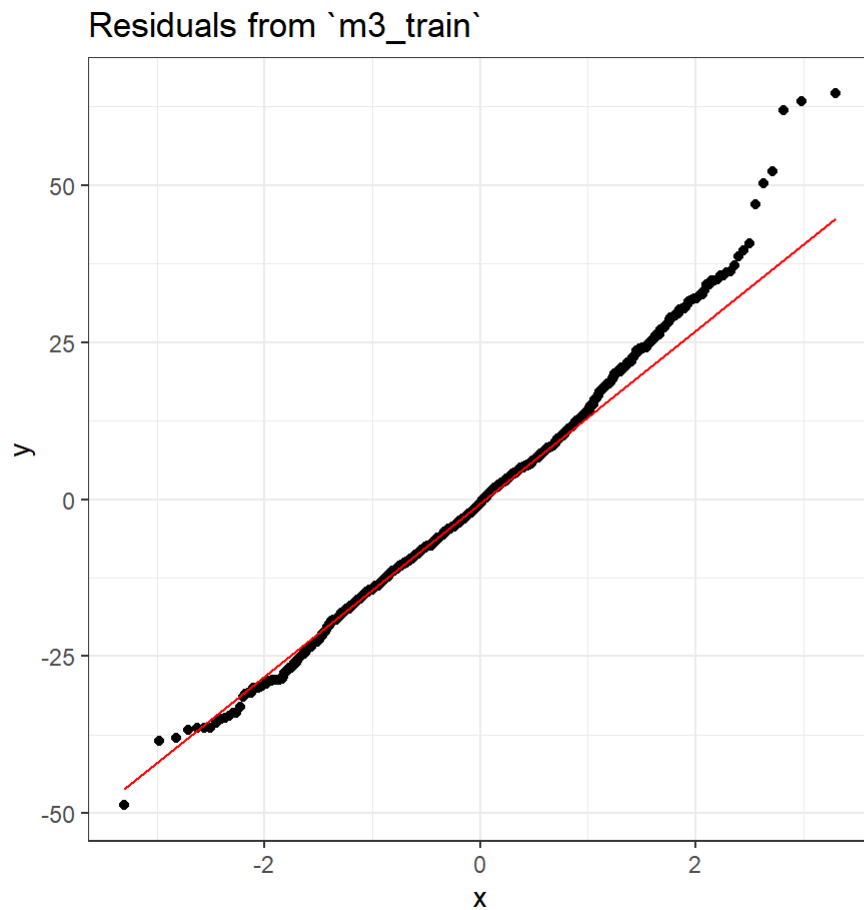Figure 1. Predicted sbp_2 from m3, m4



Figure 2. Predicted sbp_2 by ins_1

# Normality of Residuals?

```
1  p1 <- ggplot(m3_train, aes(sample = .resid)) +
2    geom_qq() + geom_qq_line(col = "red") + theme(aspect.ratio = 1) +
3    labs(title = "Residuals from `m3_train`")
4
5  p2 <- ggplot(m4_train, aes(sample = .resid)) +
6    geom_qq() + geom_qq_line(col = "red") + theme(aspect.ratio = 1) +
7    labs(title = "Residuals from `m4_train`")
8
9  p1 + p2
```

# Normality of Residuals?

# Training Set Performance

```
1  bind_rows(glance(m1_train), broom.mixed::glance(m2_train), glance(m3_train)
2            glance(m4_train)) |>
3    mutate(model = c("m1", "m2", "m3", "m4")) |>
4    select(model, r2 = r.squared, sigma, AIC) |>
5    kbl(digits = c(0, 3, 2, 1)) |> kable_styling(font_size = 28)
```

| model | r2 | sigma | AIC |
|-------|-------|-------|--------|
| m1 | 0.123 | 15.26 | 8706.4 |
| m2 | NA | 15.26 | NA |
| m3 | 0.128 | 15.24 | 8707.0 |
| m4 | 0.129 | 15.25 | 8711.5 |

- `glance()` produces different summaries for a Bayesian `stan_glm()` model like `m2`.

# Test Sample Results for Model m3

```
1  m3_test_aug <- augment(m3_train, newdata = bp_test)
2
3  ## Summarize absolute prediction errors
4  mosaic::favstats(~ abs(.resid), data = m3_test_aug) |>
5    kbl(digits = 2) |> kable_styling(font_size = 28)
```

| min | Q1 | median | Q3 | max | mean | sd | n | missing |
|---|---|---|---|---|---|---|---|---|
| 0.02 | 3.87 | 8.54 | 16.03 | 59.24 | 11.15 | 9.73 | 450 | 0 |

```
1  ## Summarize squared prediction errors
2  mosaic::favstats(~ .resid^2, data = m3_test_aug) |>
3    kbl(digits = 2) |> kable_styling(font_size = 28)
```

| min | Q1 | median | Q3 | max | mean | sd | n | missing |
|---|---|---|---|---|---|---|---|---|
| 0 | 14.99 | 72.94 | 256.9 | 3509.12 | 218.71 | 392.28 | 450 | 0 |

- MAPE = 11.15, max APE = 59.24

- RMSPE = $\sqrt{218.71}$ = 14.79

431 CASE WESTERN RESERVE UNIVERSITY

# Test Sample Results for Model m4

```
1  m4_test_aug <- augment(m4_train, newdata = bp_test)
2
3  ## Obtain mean, maximum absolute error and root mean squared error
4  m4_test_aug |> select(.resid) |>
5    summarize(MAPE = mean(abs(.resid)), maxAPE = max(abs(.resid)),
6              RMSPE = sqrt(mean(.resid^2))) |>
7    kbl(digits = 2) |> kable_styling(font_size = 32)
```

| MAPE | maxAPE | RMSPE |
| --- | --- | --- |
| 11.14 | 59.38 | 14.77 |

# Test Sample Correlation(fitted, actual)

Pearson correlation between fitted predictions and actual
sbp_2 within the test sample.

- We could also square this to get an $R^2$ result.

```
1  round_half_up(cor(m1_test_aug$.fitted, m1_test_aug$sbp_2),4)
```
[1] 0.3875
```
1  round_half_up(cor(m2_test_aug$.fitted, m2_test_aug$sbp_2),4)
```
[1] 0.3875
```
1  round_half_up(cor(m3_test_aug$.fitted, m3_test_aug$sbp_2),4)
```
[1] 0.391
```
1  round_half_up(cor(m4_test_aug$.fitted, m4_test_aug$sbp_2),4)
```
[1] 0.3945

# Comparing performance on the test data

- Which model performs best in our test sample?

| Summary | MAPE | Max APE | RMSPE | Cor(Fit,Obs) |
|---|---|---|---|---|
| m1: lm sbp_1 | 11.17 | 58.02 | 14.81 | 0.3875 |
| m2: stan_glm | 11.17 | 58.02 | 14.81 | 0.3875 |
| m3: sbp_1+ins | 11.15 | 59.24 | 14.79 | 0.391 |
| m4: sbp_1*ins | 11.14 | 59.38 | 14.77 | 0.3945 |

# Session Information

```
1  sessionInfo()
```

```
R version 4.2.1 (2022-06-23 ucrt)
Platform: x86_64-w64-mingw32/x64 (64-bit)
Running under: Windows 10 x64 (build 22000)

Matrix products: default

locale:
[1] LC_COLLATE=English_United States.utf8
[2] LC_CTYPE=English_United States.utf8
[3] LC_MONETARY=English_United States.utf8
[4] LC_NUMERIC=C
[5] LC_TIME=English_United States.utf8

attached base packages:
[1] stats     graphics  grDevices utils     datasets  methods   base
```