

# Hypotheses and Methods

In this chapter, I introduce a series of questions resulting from the literature review of Chapter 1, which I will use to formulate hypotheses. I will then operationalize these hypotheses, and attempt to predict analytical outcomes based on the theories of Chapter 1. Following these hypotheses, I will outline key methods I will use to test them.

## Hypotheses

### Questions

Before moving in to outlining hypotheses, the first step necessary is to frame a series of questions, which the hypotheses will flow from. Based on relevant research, the most obvious first question to ask would be:

Q1: *What is the effect of mail voting on turnout?*

I went through this question substantially in the previous chapter; it should be clear that depending on which paradigm of participation choice is present, the answer here can be radically different. In order to best answer the previous question, it is necessary to establish some conditions on importance of effect. Therefore it is also necessary to ask the following question:

Q2: *Is this effect significant when compared to other metrics that affect turnout?*

The last question asked in this thesis is more specific to a particular formulation of Aldritch's hypothesis on voting "at the margins". I mentioned in the previous section that VBM could be theorized to have a more significant effect when discussing elections at the local level, or the regional level, rather than national general elections. Therefore a third question is:

Q3: *Is the effect of VBM more pronounced as significant, national determinants of turnout dull?*

### Hypotheses

Using the above questions I can now move on to formulate more clear hypotheses. Before diving right into that, I note that I intend this thesis to serve two purposes: first, to test voter choice theories between each other; second, to serve as an analytical tool for later evaluations of mail voting as policy. Based on the theoretical review of the previous chapter it should be apparent that of these two purposes, the former is primarily addressed, with the later tangentially arising from my conclusions. The hypotheses in this section spring mostly from a wish to test theories of voter choice, and in particular a wish to defend the theory of voting "at the margins" as introduced by Aldritch. Therefore all hypotheses in this section will be phrased from the perspective of this theory, with the competing alternate hypotheses being counter-claims potentially rooted in different theories of voter participation.

In response to Q1, Q2, a first hypothesis is:

H1: *Mail voting is another marginal effect on voting decisions, and therefore does not significantly affect turnout*

The alternative hypothesis would be:

H1': *Mail voting significantly affects turnout, even compared to other metrics*

Similarly, for the third question, a corresponding hypothesis derived from Aldritch's paradigm is:

H2: *The effect of VBM on turnout is more pronounced as national effects dull*

The alternative hypothesis is:

H2': *The effect of VBM on turnout is consistent and independent*

## Criteria

A first, glaring issue that needs to be clarified is the apparent contradictions between my two hypothesized results. This becomes clear, however, if I define “significant effect” in the context of my first hypothesis. Aldritch’s paradigm does state that “conveniences” like mail voting should not have significant effects, but those effects are defined in the context of huge, clashing forces that vastly outweigh them. This does not necessarily mean that they are literally non-existent, but that they are poor indicators of consistently increased turnout. Therefore, I will confirm my first hypothesis not only if the effect of mail voting on turnout is statistically insignificant, but also if it is relatively small in comparison to the effects of other variables I include. I will confirm the alternative hypothesis if, across multiple of the models I will parametrize and fit, VBM retains a consistent, significant effect on turnout. If the effect is negative, this may point to a habitual or structural voting paradigm being present. If the effect is positive, this may be a signifier that issues of convenience in voting—having a mail delivered ballot, voting from your kitchen table etc.—have a particularly strong effect in the examined elections.

Moving on to the second hypothesis. It is extremely hard to correctly operationalize and account for all variables going into turnout. Therefore, instead of trying to include all national effects into a model and try to see how they interact with VBM, I will test my hypothesis on more localized elections. At least in theory, I can assume that if mail voting significantly impacts people’s decision to vote, it will be in a context where the convenience of voting significantly outweighs information effects from national media, communal pressures, or national campaigns. This can be found to some extent in primary elections, but much more significantly in off year local state elections. A potential re-formulation of the second hypothesis, that makes it more specific to the criteria I have set, is:

H3: *The effect of VBM on turnout is more pronounced in local or off year elections*

I will confirm this hypothesis if mail voting has significantly larger positive effects on turnout in smaller, local elections.

## Importance of Hypotheses

The importance of these hypotheses is intrinsically tied to the importance of different theories of electoral participation. Confirming or rejecting each hypothesis—even when only applied to a single state—serves as an argument for or against one of the aforementioned theories. The theories in and of themselves are significant, since they form a part of a broader literature on elections, democracy, and electoral processes, that can be said to be foundational to political science as a whole. Elections are the root from which all democratic governing springs; understanding why people participate in them is understanding how they choose to be included or excluded from the process of policy-building, and how they interact with the state.

Additionally, from a public policy perspective, these hypotheses are significant since they serve as metrics for the effectiveness of mail voting as an electoral reform. Whether, in general, mail voting increases turnout is directly connected to whether it is successful in expanding the democratic franchise. If it is not, questions can be raised as to the effectiveness of expanding voter access through elections administration, rather than education, or even measures like voting-day-holidays or local transportation to polling places. In local elections in particular, significant effects of mail voting could be precursors to more general involvement of individuals in their local politics. This may open the way to numerous comparative studies on local politics between states that apply VBM and states that do not.

Lastly, from a narrower perspective specific to the study of early and mail voting, my first hypothesis can still be said to be significant, yet mundane. It does its job according to the particular state I chose to look at—in this case Colorado—to add to existing literature on mail voting effects in different parts of the country. However, my second and third hypotheses are much more unique in their scope. There have not been many studies that look at VBM at a more localized level, and any addition to the literature on this front—however limited—could be significant.

## Methodology

Before directly defining all parameters of the models I will later use in writing this thesis, I will go through each type of method to provide some background on the statistics behind the models. In the next chapter, I will introduce the data and fully outline my models. This section should serve as a general introduction to the methods. I will not extensively go through the statistics behind linear or multiple regression, but will assume that it is common knowledge. For an extensive introduction to such methods, James et al.(2017) or Chihara and Hesterberg (2011) are particularly useful.

### Logistic Regression

Let function  $f : [0, 1] \rightarrow \mathbb{R}$  be defined as:

$$f(p) = \text{logit}(p) = \log\left(\frac{p}{1-p}\right)$$

This is called the logit function or, when  $p$  refers to a probability, the log-odds function. When modelling a binary response  $Y$ , which follows a Bernoulli distribution:

$$Y \sim \text{Bernoulli}(p),$$

the logit function can be used as a link function to model  $Y$  in a generalized linear model. The generic form of a generalized linear model looks like:

$$f(Y) = XB,$$

where  $Y$  is a vector of response variable values,  $X$  is a matrix of predictors, and  $B$  is a matrix of coefficients to be estimated. The function  $f$  is called a link function, because it “links” the response variable with the set of predictors included in the model. This is typically done to ensure that the range of values outputted by the model are consistent with the range of the response variable. When wanting to compute a model on a binary response through its corresponding Bernoulli distribution probability parameter, the inverse logit function should be a perfect fit for a link function, since it maps values from all real numbers to a range between 0 and 1. Using the inverse logit function, we arrive at the final form of logistic regression, which is:

$$\mathbb{P}(Y_i = 1) = \text{logit}^{-1}(XB)$$

Conveniently, despite the use of a link function, there is an easy way to interpret the coefficients of such a regression. While obviously individual values from the  $B$  matrix will not be particularly helpful,  $e^B$  can be used as a matrix of multiplicative, one-unit shifts in the value of the probability that  $Y_i = 1$ . This means that a one unit increase in any predictor will cause an effect equal to multiplying  $p$  by the exponent of the corresponding coefficient. [James et al. 2017]

### Generalized Additive Models

In simple logistic or linear regression, there is an assumption made on the functional form of the relationship between predictors and response variable. These are called parametric models, where the data is exclusively used to estimate values for coefficients. Non-parametric models, on the other hand, use the data to estimate both coefficients and the function that serves to connect response to predictors. While on the surface this seems like a great idea (more reliance on your data and less assumptions!), such an exclusively non-parametric model would suffer greatly from the curse of dimensionality—where the addition of multiple predictors or overreliance on data leads to substantial over-fitting.

The solution, then, is a Generalized Additive Model, or GAM. This model lets us fit a different functional form to each observation, allowing for assumptions to be made on the data where it is safe to do so, and for non-parametric fitting when it is necessary. This model looks like:

$$y_i = \alpha + \sum_{j=1}^p \beta_j f_j(x_{ij}),$$

where  $y_i$  the  $i$ -th response variable,  $\alpha$  is the intercept term,  $f_j, \beta_j$  a series of  $p$  functions and coefficients, and  $x_{ij}$  the  $i$ -th observation for the  $j$ -th predictor. Note that for  $f_j(x_j) = x_j$ , this is a multilinear regression! [James\_introduction\_2017]

A type of most commonly fit functions—and the type I will make use of—are smoothing splines. These are cubic functions connected at specific points called “knots”, with the limitation that the full function must be continuous and smooth. These are particularly useful when modeling time variables, as they can be fitted to variables like years or months in order to distinguish a secular trend from a general trend over time. In terms of this thesis, this will help when responding to Q2 as it was framed earlier in this chapter. [Barr\_comprehensive\_2012]

## Multilevel Models

Multilevel models—otherwise known as hierarchical or “mixed effects” models—can be intuitively pictured in two ways: either as a set of models working on different “levels”, where one is calculated first, with its effects having implications for the second, or as a model where some of the parameters estimated act under a particular series of constraints. Multilevel models are, in essence, a compromise between levels of “pooling” data. If the dataset on which parameters are being estimate operates in different units of observation—say on the individual and county level—you could run a model that treats all individuals as coming from the same larger group; this would be a complete pooling model. You could also add indicator variables for each and every group, de facto estimating  $n$  different models for  $n$  groups; this would be a no pooling model. Multilevel modelling offers partial pooling [Gelman\_data\_2006].

To consider what this model looks like, let’s assume a dataset comprising of a vector of values for the response variable  $Y$ , a matrix of  $i$  individual level predictors  $X$ , a matrix of  $j$  group level predictors  $U$ , intercept terms  $\alpha$ , individual level coefficients  $B$ , and group level coefficients  $\Gamma$ . Based on this, a multilevel model with intercept terms varying by group looks like:

$$Y_i = \alpha_{[i],j} + X_i B, \quad \alpha_{[i],j} \sim N(U_{j[i]}\Gamma, \sigma_\alpha^2)$$