

1 Chapter 9: Classifiers

1.1 Playing with fire

On an August afternoon in 2007, a fire broke out on the 17th floor of the then-vacant Deutsche Bank Building, a skyscraper situated across from the former World Trade Center in New York City. The building, seriously damaged after the 9/11 attacks, had been undergoing hazard abatement and controlled demolition, leading to changes to the building floor plans and safety apparatus. When the New York City Fire Department (FDNY) responded to the scene, it was clear the fire was a serious one, quickly escalating to a seven-alarm fire incident requiring 87 units and 475 firefighters.¹ As standpipes had been disabled and floor plans altered, FDNY units found it difficult to navigate the skyscraper and put water on the fire, resorting to unconventional methods of supplying water to crews. Eventually, the fire was put out seven hours after it started, not before two firefighters lost their lives, succumbing to cardiac arrest from heavy smoke inhalation.² In response to the tragedy, a mayoral investigation found that the deaths could have been prevented had city agencies established information sharing protocols and leveraged a risk-based strategy to mitigate and avoid hazards.³ While an ideal end state would be to end all structural fires, the recommendations focused on reducing death and injury by ensuring that FDNY had the most up-to-date ground intelligence.

Risk mitigation strategy was indeed due for improvements. Since the 1950's, FDNY building inspections were managed using a manual index card system where inspection schedules were based on tiers of perceived risk, where the riskiest buildings needed to be inspected once a year and the least risky buildings were inspected once every seven years. While it was a longstanding process, it left a shortfall in inspection coverage. Of the one million buildings in New York City, only one-third are inspection-qualified. Of those 300-thousand buildings, FDNY had historically been only able to inspect at most 10% of the buildings in a given year due to other operational priorities. This meant that even on a seven-year schedule, not all buildings would be reached and the fixed timeline meant that both perfectly safe and guaranteed fire traps had equal chance of being inspected. This could be easily changed. By incorporating the latest information about where fires did and did not occur and associating it with building characteristics, a new data-driven strategy could direct how buildings are prioritized for inspection.

In 2013, the New York City Fire Department (FDNY) set out to address the risk management problem by melding data and technology with their field operations. On the surface, the idea of using data and technology to reduce the risk of fire is quite alluring. However, under the hood, there were notable obstacles. On the operational side, buy-in was required. Anyone who has observed firefighters on scene will notice that it is a well-choreographed operation – every person knows their part and abides by the established protocols as directed by leadership. For data to drive value, it needed to be integrated and accepted into the culture of a 10,000+ person fire fighting organization. On the technical side, decades worth of index cards needed to be digitized and a scheduling platform needed to be developed. Perhaps most importantly, the system had to work. Scheduling just any inspection is simple. But scheduling inspections to buildings with observable fire risks is far more challenging as such a system would need to be able to distinguish between fire-prone and fire-proof buildings. Without effective targeting, the entire effort would be for naught.

The Commissioner and First Deputy Commissioner at the time both believed that technology had a role to play at FDNY. Aligned with Mayor Michael Bloomberg's vision of smart, data-driven government, they saw an opportunity to set an example for the nation's fire services. They relied on the Assistant Commissioner for Management Initiatives to lead a change management process with fire chiefs and fire officers, information technology (IT) managers, among others to change the flow of operations so that data served as a pillar on which FDNY could rely. Alliances were forged with leading fire personnel such as the Deputy Chief of Fire Operations and Battalion Chiefs to formalize the role of data in the culture of the fire house, amending standard operating procedures (SOPs) to use a digital inspection system. On the IT front, a lead software engineer and project manager meticulously gathered specifications that were then used to construct a scheduling platform. Recognizing that the proof was in the pudding, a Director of Analytics was hired to lead

¹<https://cityroom.blogs.nytimes.com/2007/08/18/2-firefighters-are-dead-in-deutsche-bank-fire/>

²<http://www.nydailynews.com/news/firefighters-dead-7-alarm-deutsche-bank-blaze-article-1.238838>

³http://www1.nyc.gov/assets/doi/downloads/pdf/pr_db_61909_final.pdf

the overhaul of a prediction algorithm to rank buildings based on their risk and convincing stakeholders that a statistical representation of fire ignition was indeed trustworthy. The result was the Risk-Based Inspection System (RBIS), a firefighter-facing data platform that scheduled inspections at buildings with the greatest risk of fire. Three times a week for three hours per session, fire officers logged onto RBIS to obtain a list of buildings for scheduled inspection. Buildings were selected using FireCast, a statistical algorithm developed in-house to predict fires at the building level. Through FireCast, buildings no longer used assumed a static risk classification as in the index card system, but rather a dynamic risk score took into account the latest information.

Prediction often relies on accuracy measures to determine how well algorithms perform in the field; FireCast was no different. The algorithm was able to identify buildings with fires over 80% of the time – a degree of accuracy that superceded prior attempts at the problem. Upon implementing the new system, impacts were observed in leading operational indicators. In the first month, the number of safety violations issued grew by +19% relative to the trend under the index card system, but fell to +10% in the second month. This indicated that the riskiest buildings did indeed have more observable risks than less risky buildings, but the amount of observable risk fell as building inspection teams progressed down the risk list.

From a statistical perspective, the prediction should have yielded far more violations, but efficacy of the prediction program was limited by (1) a fire unit's time budget to conduct inspections; (2) a policy requiring that time had to be set aside for weekly inspections, which at times led to inspecting buildings that were not observably risky after all truly risky buildings were exhausted; (3) the rule of law giving residents the right to refuse inspection. To measure efficacy, FDNY developed an indicator known as the Pre-Arrival Coverage Rate (PACR), which measures the proportion of buildings that experienced a fire that were inspected within some period (90 days) before the fire occurred – essentially measuring if fire companies had the opportunity to evaluate risks of priority buildings. Under FireCast, FDNY had achieved a PACR of 16.5%, which was an eightfold improvement over the old strategy that yielded 1.5%.⁴ ⁵

Since Firecast was launched in 2014, other fire prediction efforts have emerged around the United States such as the Firebird open source system for Atlanta in 2016⁶ and a spatio-temporal fire prediction approach for Pittsburgh in 2018⁷.

1.2 What's a classifier?

The RBIS/FireCast is an example of a *classification* problem – a task in which a model determines which group or *class* does an observation belongs based on its attributes, doing so based by learning from known examples. Examples must include the factual (what happened) and counter-factual (what did not happen) in order to distinguish between potential fires from non-fires.

Classification is nothing new in everyday life as we use our own mental classification models to contextualize the world around us. For example, marketers and advertisers are always looking to get product offerings in front of prospective customers and will often purchase lists of people and apply models based on past customer behavior in order to identify those who are most likely to be interested.⁸ The criminal justice system has incorporated risk classification models to determine if those involved in alleged crime pose a flight risk if bail is posted.⁹ On a more futuristic front, the technology behind self-driving cars uses a complex array of sensors and cameras that are processed by classifiers in order to distinguish between cars, people, motorbikers and cyclists.

The same is true with fires.

⁴<http://www.nfpa.org/news-and-research/publications/nfpa-journal/2014/november-december-2014/features/in-pursuit-of-smart>

⁵<https://www.nist.gov/publications/research-roadmap-smart-fire-fighting>

⁶<http://firebird.gatech.edu/>

⁷<http://www.kdd.org/kdd2018/accepted-papers/view/a-dynamic-pipeline-for-spatio-temporal-fire-risk-prediction>

⁸ref needed

⁹ref needed from Vera Institute

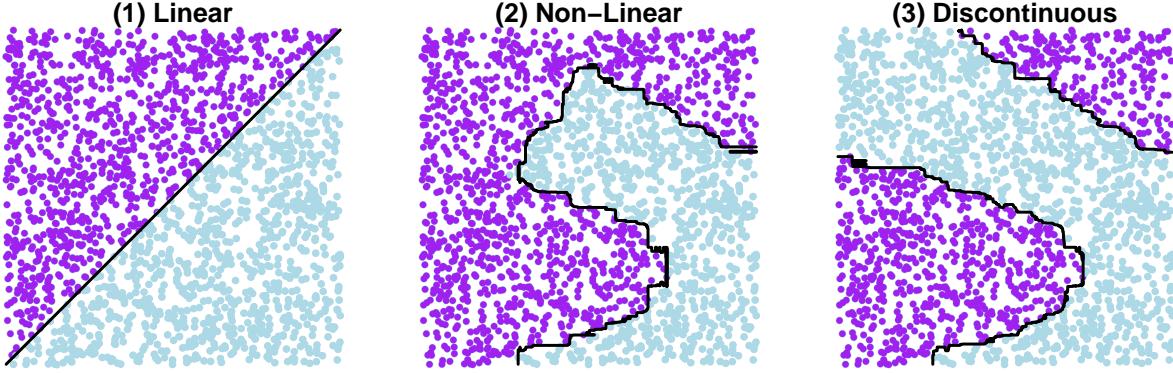


Figure 1: Types of classification problems

By examining buildings that have and have not caught fire in the past, we are able to learn what pattern characteristics are associated with greater risk of fire. We can *score* or apply the learned pattern to new records to obtaining the probability of fire so that fire fighters could have a sneak peak of what may happen to other buildings as their conditions change. There may be thousands of variables that can play a role in predicting fires.

If fires are truly predictable, we can employ supervised learning to map how input variables can distinguish buildings that had fires from those that did not. Otherwise stated, which fire status class does a building likely belong? Given a binary outcome $Y(\text{Fire})$ we can determine class membership as a function of the building's characteristics:

$$Y(\text{Fire}) = f(\text{Building characteristics, Location, Complaints, ...})$$

It sounds simple enough, but as it turns out, there are many different ways that algorithms can associate building characteristics to fires. We can illustrate the complexities of this task under three hypothetical scenarios, plotting two input variables on the X and Y axes and color-coding two classes in purple and light blue. The solid black lines represent the true *decision boundary*, or the threshold at which an observation is classified in a specified class.

At first, the idea of classifying records may seem straight forward, but can become a rabbit hole. In policy, we tend to start from a normative theory of how a phenomenon functions – perhaps a simple explanation backed by linear logic. For example, a building that is bigger and older may be a reasonable working hypothesis for identifying high risk buildings. But as we investigate more and deeper, we may add exceptions to the rule as we discover cases that do not conform. Those exceptions may improve the number of correctly classified records, but detract from the simple narrative.

The function that determines the decision boundary can take on many forms and each balances interpretability and accuracy. A simple linear boundary can be represented using a *logistic regression* – the workhorse of the natural and social sciences. While the method provides an analytically convenient answer, the resulting decision boundary may miss finer, more disjointed patterns. Non-parametric methods offer a far more flexible solution to prediction, but at the cost of interpretability. A simple technique known as *k-nearest neighbors* is useful when k-number of similar, comparable observations can serve as a reference for informing the prediction of a given point. *Decision tree learning* along with its many variants such as *Random Forest* learn patterns by partitioning a sample into finer more homogeneous sub-samples. When faced with infinite distributions of data, each method responds to the circumstances differently.

Each technique has its own set of assumptions that in turn make it better suited for certain policy problems and certain types of data. But there are basic considerations that underlie our use of any classifier, namely (1) the idea of separability, (2) balancing interpretability with prediction, and (3) the different definitions of accuracy.

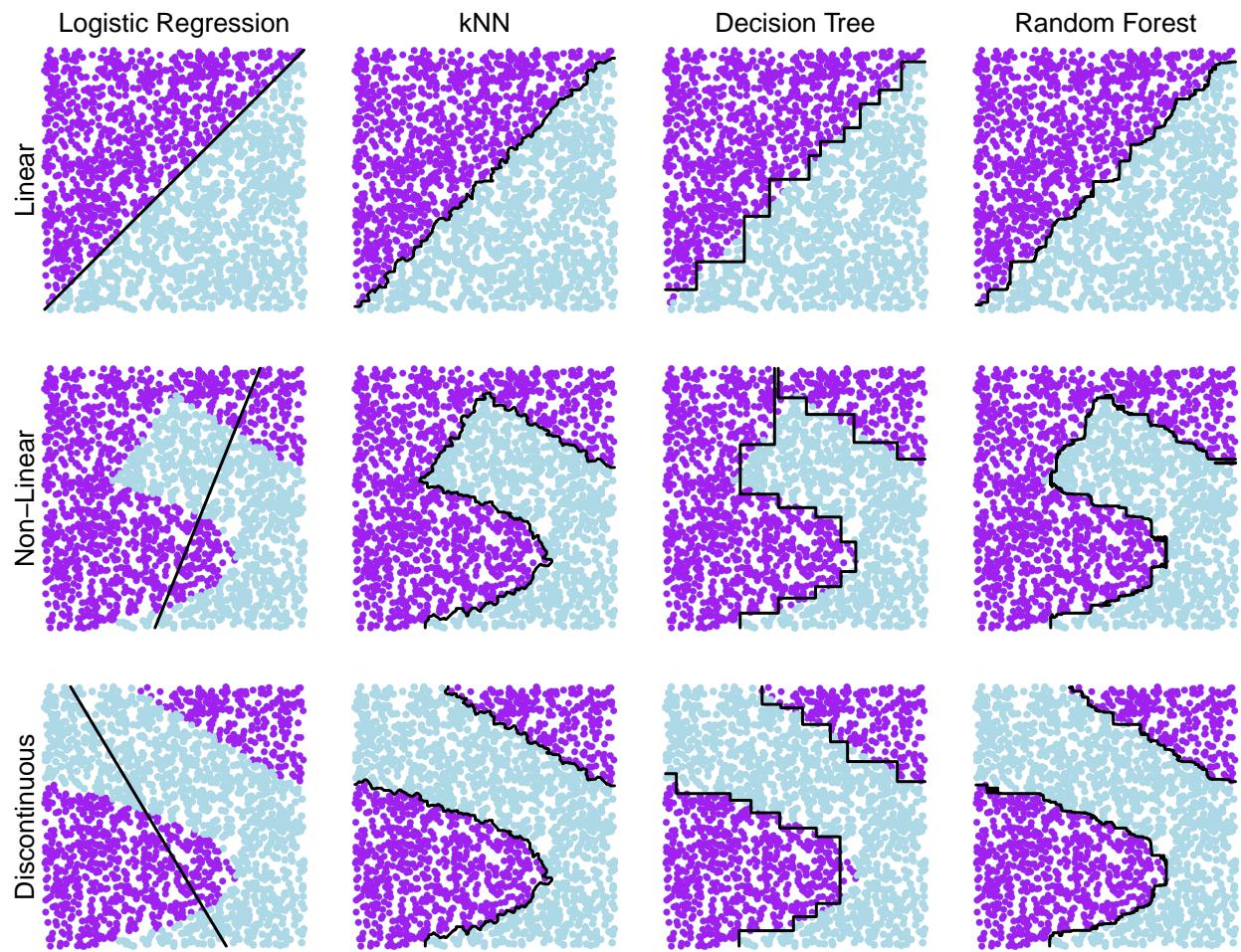


Figure 2: Linear, Non-Linear and Discontinuous Classification Problems.

1.2.1 Separability

A fire prediction algorithm needs to actually predict where fires and non-fires will be. A bail algorithm needs to be able to distinguish between those who are flight risks and those who are not. Simply training a classifier will not do. The success of a classifier fundamentally lies in if the classes are *separable* based on its input variables – basically the idea that the algorithm can learn from the underlying data and distinguish one class from another. Successfully distinguishing between classes can be measured through various forms of accuracy, but fundamentally is dependent on finding as many *True Positives* (TP) and *True Negatives* (TN) as possible.

Separability manifests itself differently depending on the data type. When the inputs are discrete variables, the goal is to find a set of variables that can best describe the target. In the example below, we analyze how building violations are related to the target fires. The ideal situation is to place the mass of sample along the diagonal with TNs (no fire - no violation) and TPs (fire - violation). In this case, whether a building had a building violation can accurately distinguish 80% of observations.

Table 1: Separability of two classes given a separable discrete variable.

	No Violation	Violation - T
Fire	10.9	22.6
No Fire	59.7	6.8

Alternatively, a low separable case may will have a large proportion of observations along the other diagonal, indicating many *False Positive* (Type I error) and *False Negatives* (Type II error). We can see that the majority of cases have been incorrectly classified. Classification problems are a balancing act: it is unlikely that a model will yield 100% accuracy everytime and often times, we will need to weigh accuracy against classification error. As we will see later in the chapter, more True Positives may also lead to more False Positives and how we balance these are very much a policy decision.

Table 2: Separability of two classes given an non-separable discrete variable.

	Elevator Apartment - F	Elevator Apartment- T
Fire	2.3	31.2
No Fire	6.8	59.8

When inputs are continuous variables, separability can be assessed in terms of means and dispersion of the target classes. For example, a low separability scenario (left) would be one where the input variable's distribution for each class substantially overlap, which suggests an absence of any distinguishing information. High separability (middle), in contrast, would have means that are significantly different from one another and the distributions themselves overlap minimally. But perhaps the neatest thing is perfect separability (right). In this case, we basically can produce a perfect definition of the outcome in question, in which case we do not need a classifier to model the relationship – just a threshold to serve as the definition.

The bottom line about separability is that a good data scientist will check their assumptions. In policy, there will be prevailing theories and conventional wisdom that dictate how certain factors influence a phenomenon. But what may sound good may in actuality have little separability and in turn offer little predictive power. In those cases, it is worth revisiting and revising the assumptions. After all, the goal of a classifier is to classify.

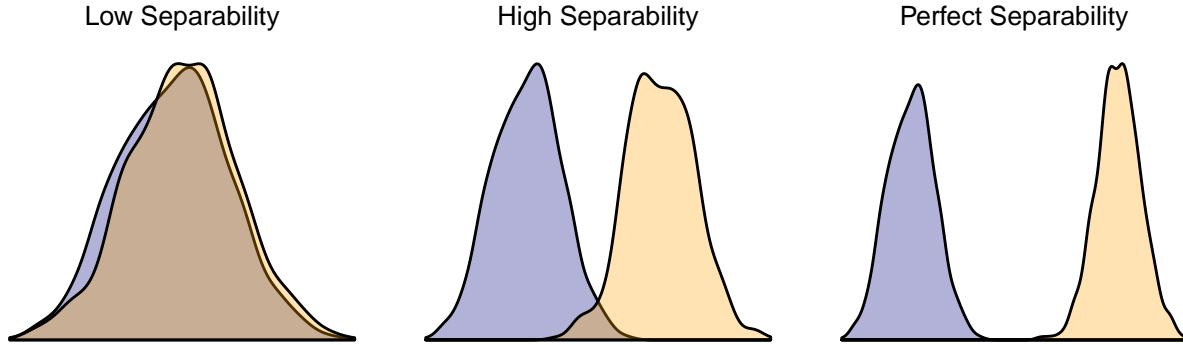


Figure 3: Separability of two classes given a continuous variable.

1.2.2 Measures of accuracy

Accuracy is central to evaluating if a classifier's results are useful. This requires understanding of a minimum of two ideas: what does probability have to do with accuracy and what exactly is the substance of accuracy.

Predicted Probability. Most classifiers output a form of conditional probability. Accuracies are derived from probabilities produced by the classifier, indicating whether a given observation is predicted to belong to a given class (e.g. fire vs. no fire, Yankee fan vs. Red Sox fan, Sith vs. Jedi, Republican v. Democrat). In order to convert probabilities into a prediction of a class, we need to set a threshold. So, what is the classification threshold?

Short answer: It depends on the sample.

For balanced samples in which the outcome variable's classes are represented in approximately equal proportion, the threshold is $\text{pr}(Y = 1) \geq 0.5$ as the objective probability of occurrence is close to 0.5. In reality, this may not be all that common scenario due to *class imbalance* and *rare events*. Class imbalance is the case in which the *minority class* – the class with fewer observations – is proportionally less prevalent in the sample than other classes. For example, political campaigns often deal with voter data in districts that are overwhelming in favor of one party over another. Rare events are infrequent events in which there is a super minority of observations. Fires and disasters tend to fall into this category. In some cases, fewer than $n = 200$ observations may lead to biased estimates.¹⁰ We will revisit class imbalance and rare events later in this chapter, but for now, let's assume that the classification threshold is $\text{Pr}(Y = 1) \geq 0.5$.

Accuracy. Classifier accuracy measures rely on metrics derived from the *confusion matrix*, or a $n \times n$ table where the rows represent actual classes and columns represent predicted classes. For a two class problem, the confusion matrix is a 2×2 table.

Table 3: Structure of confusion matrix

	Predicted (-)	Predicted (+)
Actual (-)	True Negative (TN)	False Positive (FP)
Actual (+)	False Negative (FN)	True Positive (TP)

Each cell of the confusion matrix is a building block required to calculate accuracy:

- The True Positive (TP) is the number of cases in which the actual positive observations ($Y = 1$) are correctly predicted (e.g. model predicts a fire and a fire actually occurs). Note that P is used to denote the total number of positive records.
- The True Negative (TN) is the number of cases where the actual negative observation ($Y = 0$) are correctly predicted. Note that N is used to denote the total number of negative records.
- The False Positive (FP) is number of cases where the actual label was $Y = 0$, but the model classified a record as $\hat{Y} = 1$. This is also known as *Type I error*.
- The False Negative (FN) is number of cases where the actual label was $Y = 1$, but the model classified a record as $\hat{Y} = 0$. This is also known as *Type II error*.

In a perfectly balanced sample of $n = 100$, we would expect the $P = 50$ and $N = 50$. Likewise, a perfect predictions should yield $TP = P$ and $TN = N$ along the diagonal. This is a rare case, but it stands that the goal is to check if the diagonal captures the majority of observations.

In contrast, a model with little predictive power will have the majority of observations in the upward diagonal. Below, the matrix suggests the trained model has little predictive power as it classified the vast majority of records as negative even when some should have been positive.

Table 5: Confusion matrix for poor predictions.

	Predicted (-)	Predicted (+)
Actual (-)	40	10
Actual (+)	45	5

Despite having well-defined building blocks of accuracy, the tricky thing is that there is not just one measure of accuracy or error. In fact, there are many measures of accuracy that emphasize very different predictive qualities.

^ADD EXAMPLE HERE

Measure	Formula	What It Answers
Individual Measures		
True Positive Rate (TPR), Sensitivity, or Recall	$TPR = \frac{TP}{TP+FN}$	What proportion positive cases were correctly identified?
True Negative Rate (TNR) or Specificity	$TNR = \frac{TN}{TN+FP}$	What proportion negative cases were correctly identified?
False Positive Rate (FPR)	$FPR = \frac{FP}{FP+TN}$	What proportion of negative cases were incorrectly predicted as positive? This is also known as the false detection rate or Type I error.
False Negative Rate (FNR)	$FNR = \frac{FN}{TP+FN}$	Proportion of positive cases that were incorrectly predicted negative. This is also known as the false alarm rate or Type II error.
Predicted Positive Value (PPV) or Precision	$PPV = \frac{TP}{TP+FP}$	Proportion of positive cases that were incorrectly predicted negative.
Overall Measures		
Accuracy (ACC)	$ACC = \frac{TP+TN}{n}$	What proportion of records were correctly classified?
F1-Score (F1)	$F1 = \frac{2}{\frac{1}{TPR} + \frac{1}{FPR}}$	Alternative method of calculating accuracy using a harmonic mean
Area Under the Curve (AUC) or Concordance statistic	$AUC = \int_0^1 ROC(u)du$	The AUC is the area under the Receiving-Operating Characteristic (ROC) Curve. The ROC is derived by plotting values of TPR (Sensitivity) and FPR for varying probability cutoffs c.

1.2.3 Interpretability versus prediction

When selecting models for classification tasks, data scientists often times balance interpretability and predictive accuracy. Some situations require a narrative to communicate the insights. In the social and natural sciences, classification methods like logistic regression are favored for their use in *parameter estimation* to infer

relationships between variables captured in coefficients. These coefficients are directly interpretable, showing how a variable can contribute or detract from the probability of an outcome holding all else constant. For example, a building built X decades ago is Y-times more likely to catch fire than a building built this decade holding all else constant. This simple factoid can facilitate narratives that communicate insight. But there is a tradeoff. While interpretable methods may extract the gist of the relationships, they may at times miss the finer variations in the data necessary for a reliable prediction.

Alternatively, some audiences are far less interested in the story, but just want an accurate number. In other words, predictive accuracy is prized in which case it may be worth considering a body of exciting methods that have arisen from statistics and computer science that are optimized for the task. These other methods may be more versatile, adapt to scenarios where there are more variables than observations, find interactions between variables and nonlinear patterns, and optimize for robustness. In the tech sector, for example, data science pursuits are often a matter of how well classifiers can scale to service their customers and drive sales. Understanding which variables definitively drive the predictions can be useful for teams focused on communicating insights, but not as much for the technology side of the house.

The bottom line is that one's choice of classifier depends on how much one values interpretability versus accuracy.

1.2.4 Fairness and transparency

In the area of criminal justice, algorithms such as COMPAS are actively being used to predict the chance that someone accused of a crime will recidivate within two years. While classifiers may offer a degree of efficiency when sifting through a multitude of data and are relied upon by judges in passing sentences, there are significant concerns regarding their ethical use. In 2016, ProPublica conducted research that examines the fairness of risk ratings from COMPAS, finding that while the algorithm is able to correctly predict recidivism 61% of the time, black defendants are 1.9-times more likely than white counterparts to be labeled high risk but not actually re-offend.¹¹ In other words, if two defendants with the same criminal history were scored by COMPAS and the only difference is their race, the two defendants would be treated differently.

It is unsurprising that *fairness* and *transparency* are emerging areas of focus in the use of classifiers. Fairness can take on many definitions, but for the purpose of this text, we define it as whether two or more subpopulations receive equal treatment. Transparency means that the provenance of the underlying information and its method of producing predictions are available for review – not that it fits a neat normative, but that it can be traced and scrutinized.

1.2.5 Common techniques

There are hundreds of classification algorithms that are easily accessible using modern data science software, but six algorithms cover much of the use cases that are relevant for public policy. Each method has its strengths and weaknesses, and each is more appropriate in certain contexts than others.

Table 7: Overview of classifiers

Method	Common Uses	Example Data
k-Nearest Neighbors	Imputation of missing data; Relating data that can be treated like coordinates.	Hurricane response
Logistic Regression	Interpretation and narrative building; Hypothesis testing.	Health care coverage
LASSO logistic	Sift through large number of variables to surface likely covariates; Estimate logistic regression when there are more variables than observations.	Health care coverage

¹¹<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

Method	Common Uses	Example Data
Decision Tree	Produce predictions that can interact or cross variables together to see joint effects; Resulting nodes can be used as a “profiles”.	Activity trackers
Random Forest	Produce predictions that can interact or cross variables together to see joint effects; Reasonable for cases when straight forward explanations are foregone.	Activity trackers

1.3 K-Nearest Neighbors (KNN)

As hurricanes become more intense and leave a trail of destruction, city governments will need to be able to more efficiently triage requests for help. Let’s take the example of Hurricane Sandy and its effect on NYC. One of the main services offered by cities is the management and care of its trees. A downed tree can cause property damage, bodily harm and traffic disruptions. Due to the high wind and lush foliage during Sandy, many trees fell.

In NYC, the Department of Parks and Recreation is responsible for tree removal. When a resident makes a call to the city’s services hotline 311, a work order is created and a tree removal team is dispatched. This may be a transactional process: one call for tree removal, one tree is then removed. As it takes time for crews to move and set up, a first-in/first-out queuing process can be inefficient. Imagine if 20 of 100 blocks in a neighborhood were flagged for tree removal. It would make sense to use that call data to identify other blocks that may also have downed trees.

We would expect that downed trees are more likely to occur in *pockets* and proximity is the best indicator of activity. As the city knows where residents call for tree- and non-tree-related issues, we can use the location of the calls to triangulate on likely problem areas as well as anticipate pockets of yet-to-be-reported downed trees, or at least serve that is a reasonable working theory.

For this task of predicting based on proximity, k-Nearest Neighbors (KNN) can help.

1.3.1 Under The Hood

K-nearest neighbors (KNN) is a non-parametric, instance-based algorithm that is based on a simple idea: *observations that are closer together are more likely to be similar*. The method is non-parametric as it does not directly use its inputs x to determine the value of y , but rather assumes that inputs serve as ways to quantify how close points are. It is instance-based as each prediction is determined on a case-by-case basis using its surrounding neighbors.

The technique is simple. For each case y_i :

1. *Distance.* First, we calculate the distance d to all other records with known outcomes. Distance most commonly takes the form of Euclidean distance, which is appropriate with continuous values. For cases where the underlying data are boolean or binary, Manhattan distance may be more appropriate. In effect, the input variables X serve as sets of coordinates to triangulate which points are closer to a given example. Note that treating variables like coordinate sets implies two things. First, all variables have equal importance – no single variable should weigh in on distance more than any other. Second, the scale of the observations would need to be the same.

$$d = \text{distance} = \sqrt{\sum_{i=1}^n (x_i - x_0)^2}$$

2. *Voting.* For the k nearest observations to a given observation, calculate the proportion of observations in each class j in Y . This procedure yields a conditional probability for each observation, which is converted into a predicted class through *majority voting* – assign an observation to the class that is

most represented in the neighborhood. There are various flavors of the voting calculation that account for distance from a given observation.

Voting Type	Formula	Interpretation
Rectangular	$Pr(Y = j) = \frac{1}{k} \sum_{i=1}^k I(y^i = j)$	Calculate the proportion of j based on k nearest neighbors.
Inverse	$Pr(Y = j) = \sum_{i=1}^k w(d)(y^i = j)$ where $w(d) = \frac{1}{d_i \sum_{i=1}^k (\frac{1}{d_i})}$	Calculate the weighted proportion of j based on the inverse distance to k nearest neighbors.

3. *Tuning.* The method is sensitive to the value of k , requiring tuning – or testing different values of k . When $k = 10$, the conditional probability for y_i reflects the 10-nearest neighbors. When $k = n$, the conditional probability is the sample mean.

The above process yields the results for just one value of k . *Is that value of k the right one?* Like many other algorithms, KNNs are an iterative procedure, requiring tuning of *hyperparameters* – or values that are starting and guiding assumptions of a model. In the case of KNNs, k is a hyperparameter and we do not precisely know the best value of k . Tuning hyperparameters involves a grid search in which a range of possible hyperparameters at equal intervals are tested and compared to find the optimal.

To illustrate how tuning works, we have extracted a raster data from the US Department of Agriculture (USDA) CropScape data layer derived from satellite imagery.¹² In particular, we focus on a small area in Kansas that grows corn (yellow) and soybeans (green). More often than not, data will only contain a sample of the full picture. In the case of spatial data, KNNs are quite useful for filling in the full picture. Suppose we only have a 10% sample of the farmland – how important is k ? We can see that as the value of k grows, the corn fields become overly represented as it is the majority class in the sample. The TPR is highest when the TNR is highest, demonstrating that there is a balancing act when choosing the value of k .

1.3.2 Tips of the trade

Tuning. As the accuracy of a KNN model is dependent on finding the optimal value of k , it makes the process far simpler to think about the tuning process as a grid search. A systematic way of testing for k is try test all values from $k = 1$ to $k = \sqrt{n}$ in multiples of one's choosing. The idea is to get a ballpark sense of what works, then hone in on the best value of k . Log the values of k that are tested and compare the accuracies to find the optimum. For example, searching for values of k that provide the best prediction in the cropland example suggests that TPR is near its max value at $k = 1$ and the FPR is minimized at the same value, thus choosing the minimum value of k would likely be adequate.

There are a number of other factors that influence the performance of the algorithm:

- *Scale.* A trained KNN is influenced by the scale of its inputs. If one input has a scale from 0 to 10,000 and another ranges from 0.1 to 0.3, a KNN will lean more heavily on the latter variable. To ensure equal weights, it is necessary to transform variables into a standardized scale:

$$scaled = \frac{x_i - \mu}{\sigma}$$

where the result transformed variable is mean centered and standardized. The new unit is standard deviations.

- *Grids.* Similar to the scale issue, KNNs are particularly effective in data that are distributed on a grid – measurements along a continuous scale at equal increments, but may be a poor choice when the data are mixed data formats such as integers and binary.

¹²<https://nassgeodata.gmu.edu/CropScape/>

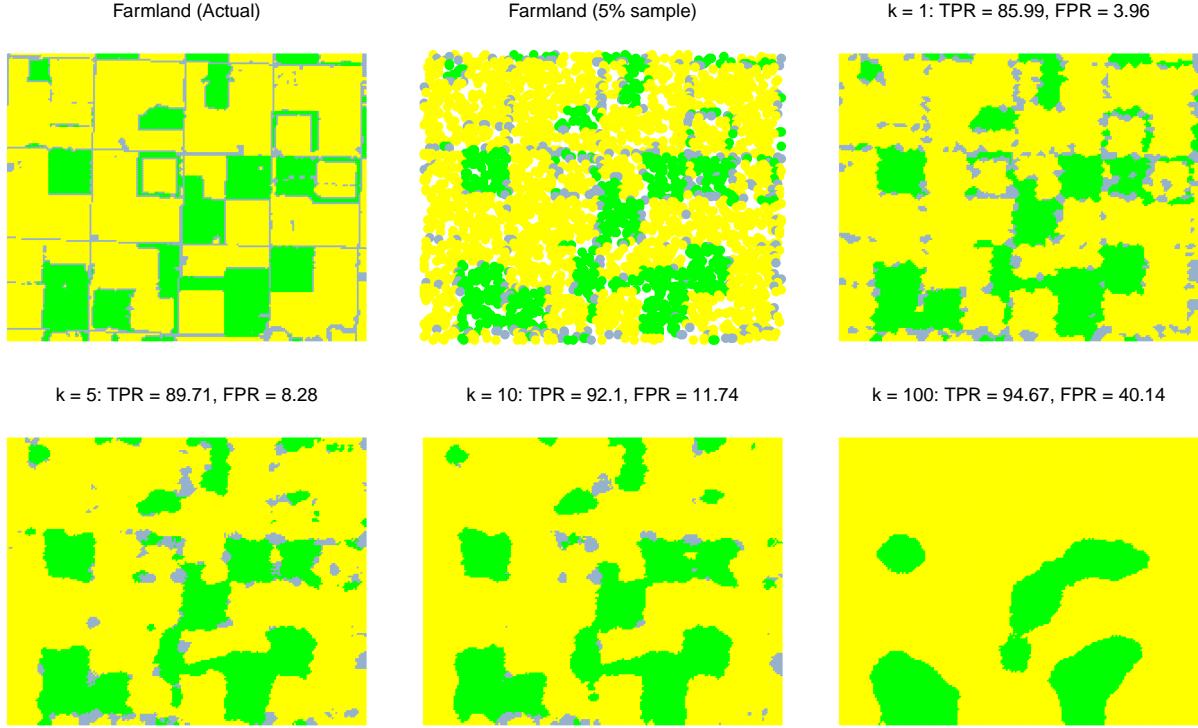


Figure 4: Comparison of prediction accuracies for various values of k .

- *Symmetry.* It's key to remember that neighbors around each point will not likely be uniformly distributed. While kNN does not have any probabilistic assumptions, the position and distance of neighboring points may have a skewing effect.

Usage. KNNs are great in some cases; Not so much in others.

KNNs are commonly associated with imputation of missing values and scenarios where proximity of observations has some bearing on the predictive accuracy. But setting up the KNNs requires some care.

As scale matters, data sets with mixed data types (discrete, continuous) need to be transformed into the same units. Discrete variables can be converted into a dummy variable matrix. Continuous variables can be binned into discrete levels, then converted into a dummy variable matrix as well. This effectively means that all variables are in terms of 0/1 and a Minkowski distance may be more appropriate to relate distances than Euclidean.

KNNs are best used when data sets are relatively smaller with fewer variables as each distance calculation is computationally taxing. Furthermore, as more variables are added, the importance of any one variable is diluted – it may be worth trying another algorithm to sift through the data.

Lastly, KNNs are not interpretable as it is a nonparametric approach. It should be instead be viewed as a processing method to fill in the gaps.

Table 9: The good and ugly of KNNs.

Useful Properties	Challenges
Efficient and timely when there are relatively few variables.	Mixed data types require convert all data into dummy matrices
Effective in capturing patterns in cases where proximity matters.	Does not offer an interpretation.

| Useful Properties | Challenges |
| Common choice for imputing missing values. | |

1.3.3 DIY: Anticipating the extent of damage from a storm

We now turn back to our motivating example. Suppose the location of all calls for non-emergency help from the day of Hurricane Sandy are captured in NYC's 311 system, yet there are still neighborhoods that likely have downed trees but have not reported it. *From what we know, how can we guess the disposition of other parts of the city to have a fuller picture?*

Prepare the data. Some calls for help are associated with downed trees while others may concern non-emergency issues. By this logic, we could assume that a downed tree would be called in if a call were made at all.

Our data set contains $n = 7513$ observations, each of which is a 1000 foot by 1000 foot area. When plotted, the data set captures the outline of NYC.

```
#Load data (need persistent link)
nyc <- read.csv("data/sandy_trees.csv")

#Plot
plot(x = nyc$xcoord, y = nyc$ycoord,
      xlab = "x coordinate", ylab = "y coordinate",
      cex = 0.3, pch = 15, asp = 1, col = "grey")
```

For simplicity, we will focus on the two largest and geographically connected boroughs, Brooklyn (BK) and Queens (QN). The remaining boroughs of the city are separated by rivers, which may be better represented using separate models. We subset our data using the `boro` variable to focus, scale the geographic coordinates using the `scale` function, and map the resulting subset.

```
#Subset
nyc <- subset(nyc, boro %in% c("QN", "BK"))

#Standardize input variables
nyc$xcoord <- scale(nyc$xcoord)
nyc$ycoord <- scale(nyc$ycoord)

#Plot
```

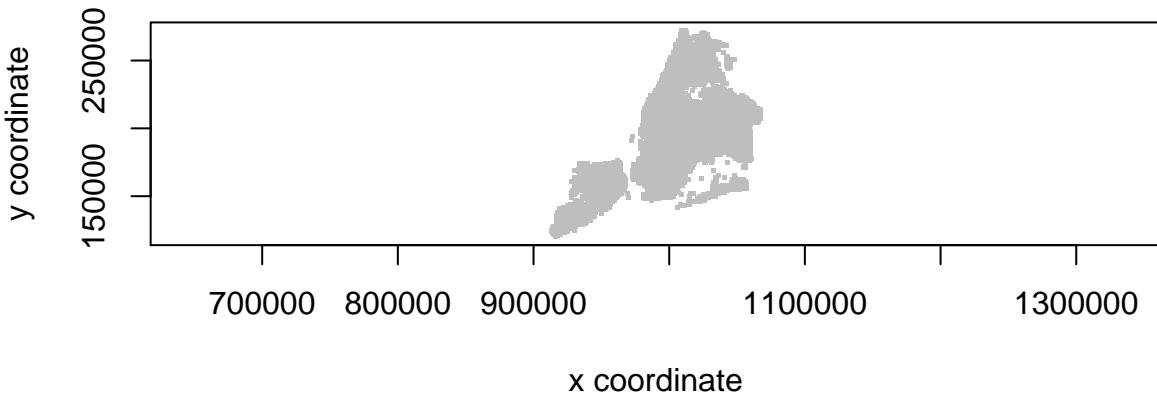


Figure 5: Plot of X-Y coordinates for NYC.

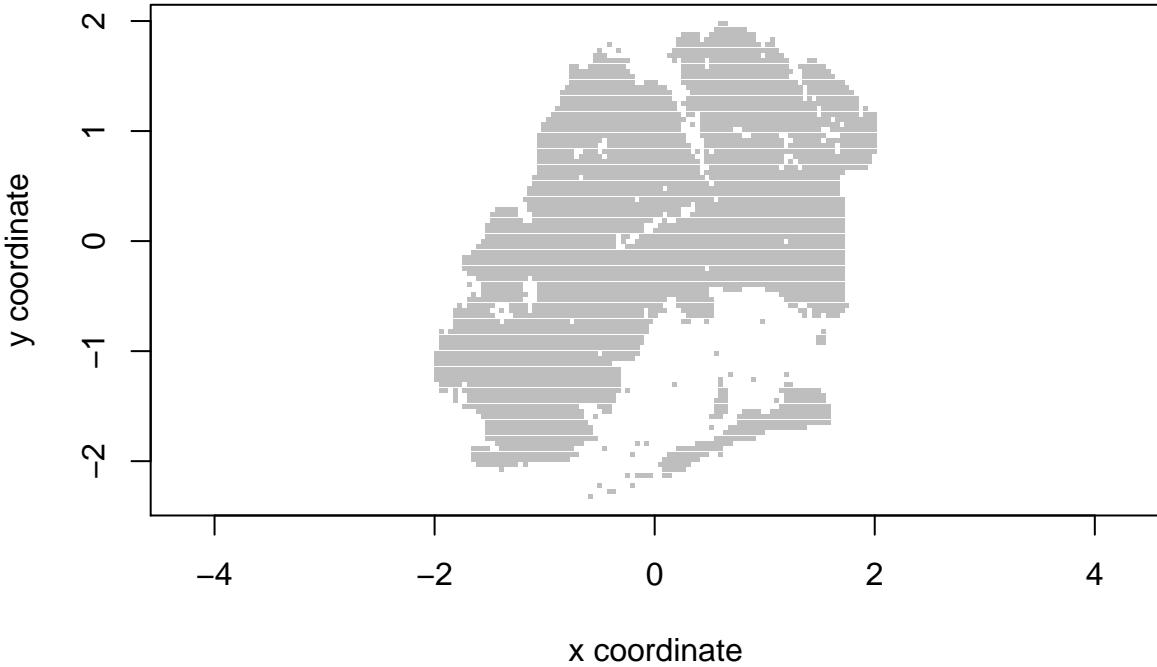


Figure 6: Plot of X-Y coordinates for Brooklyn and Queens.

```
plot(x = nyc$xcoord, y = nyc$ycoord,
      xlab = "x coordinate", ylab = "y coordinate",
      cex = 0.3, pch = 15, asp = 1, col = "grey")
```

Train. While this is a retrospective analysis, we simulate the process of producing the complete map as if we had partial information. The `nyc` data frame is split into a `train` set, keeping only locations where the target variable `tree.sandy` are available. A quick tabulation shows that the $n = 1550$ of the $n = 1946$ training set observations have a downed tree reported.

```
#Subset training sample
train <- subset(nyc, !is.na(tree.sandy),
                 select = c("ycoord", "xcoord", "tree.sandy"))

#Split out
table(train$tree.sandy)

##
##     0     1
## 396 1550
```

The test set is the entirety of Brooklyn and Queens. The `tree.next7` variable flags any location that had a report of a downed tree over the seven days after the hurricane.

```
test <- subset(nyc,
                select = c("ycoord", "xcoord", "tree.next7"))
```

With the data in the right shape, we load the `kknn` library:

```
#Call "class" library
library(kknn)
```

The KNN algorithm needs to be calibrated for the best k using the training set, then applied to a test set.

To do this, we will use the `kknn` library. The training portion uses the `train.kknn()` function to conduct k-folds cross validation, then the scoring uses `kknn()`. While both functions can be fairly easily written from scratch (and we encourage new analysts to write their own to intimately understand the assumptions), we will plow forth with using the library.

In order to find the optimal value of k , we will execute the `train.kknn()` function, which accepts the following arguments:

```
train.kknn(formula, data, kmax, kernel, distance, kcv)
```

- `formula` is a formula object (e.g. “`no.coverage ~ .`”).
- `data` is a matrix or data frame of training data.
- `kmax` is the maximum number of neighbors to be tested
- `kernel` is a string vector indicating the type of distance weighting (e.g. “rectangular” is unweighted, “biweight” places more weight towards closer observations, “gaussian” imposes a normal distribution on distance, “inv” is inverse distance).
- `distance` is a numerical value indicating the type of Minkowski distance. (e.g. 2 = euclidean, 1 = binary).
- `kcv` is the number of partitions to be used for cross validation.

The flexibility of `train.kknn()` allows for test exhaustively and find the best parameters. Below, we conduct 20-folds cross validation testing between $k = 1$ and $k = 100$ neighbors using two kernels (rectangular and inverse) that impact the voting step. This simple command does much of the hard work by running the KNN algorithm 2000 times (20 cross-validation models for each k and `kernel` combination), then surfaces the best parameters. We store the results in `fit.cv`.

```
#Set seed to ensure cross validation is replicable
set.seed(100)

#Run with 20-folds cross validation
fit.cv <- train.kknn(tree.sandy ~ ycoord + xcoord ,
                      data = train,
                      kcv = 20,
                      distance = 1, kmax = 100,
                      kernel = c("rectangular", "inv"))
```

Within `fit.cv` is a `best.parameters` element that KNNs perform the best when $k = 59$ using an inverse distance kernel.

```
#Plot Cross Validation
plot(fit.cv)

#Retrieve best parameters
best <- fit.cv$best.parameters
```

With the KNN algorithm tuned, we can now proceed to scoring the test set using the `kknn()` function. The function syntax is as follows:

```
kknn(formula, train, test, k, kernel, distance)
```

- `formula` is a formula object (e.g. “`no.coverage ~ .`”).
- `train` is a matrix or data frame of training data.
- `test` is a matrix or data frame of test data.
- `k` is the number of neighbors.
- `kernel` is the type of weighting of distance (e.g. “rectangular” is unweighted, “biweight” places more weight towards closer observations).
- `distance` is a numerical value indicating the type of Minkowski distance. (e.g. 1 = binary, 2 = euclidean,).

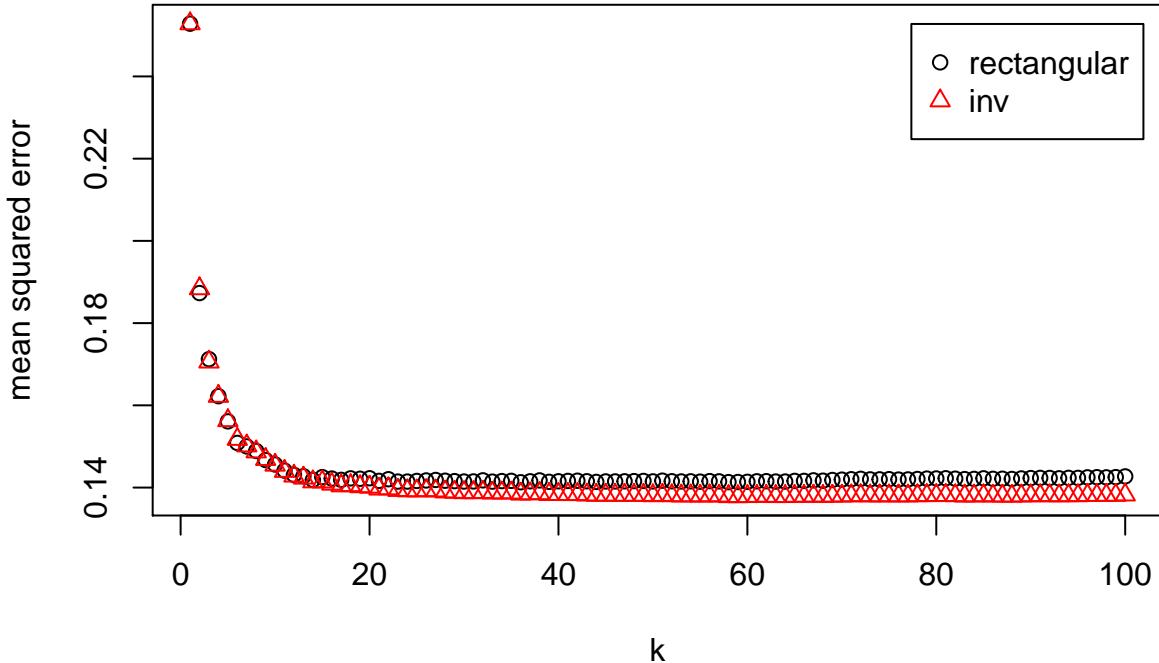


Figure 7: 20-fold cross validated errors for $k = 1$ to $k = 100$

Notice that in the following code block, we train the KNN and apply it to the test sample all in one step as the KNN itself does not learn patterns, but just applies a simple calculation following a pre-specified routine. This is a marked difference compared with other algorithms covered in this chapter.

```
#Apply tune KNN parameters
fit <- kknn(tree.sandy ~ ycoord + xcoord,
            train = train,
            test = test,
            k = best$k,
            kernel = best$kernel)

#Produce
test$prob <- fit$fitted.values
test$tree.next7[is.na(test$tree.next7)] <- 0
```

Evaluate. With all the right pieces computed, we can examine how closely the predictions based on tree downing patterns on the day of Hurricane Sandy compare with where trees were reported to have fallen over the 7 days that followed. During the storm, approximately 43% percent of the focus area made a call, of which 80% reported a downed tree. This appears as a cloud of points capturing the gist of the downed tree pattern.

```
par(mfrow = c(1,3))

plot(train[,2:1], main = "(1) Calls during storm",
      col = rgb(train$tree.sandy , 0, 1 - train$tree.sandy, 1),
      cex = 0.4, pch = 16, asp = 1)

plot(test[,2:1], main = "(2) Predicted probabilities",
      col = rgb(test$prob, 0, 1 - test$prob, 1),
      cex = 0.4, pch = 16, asp = 1)
```

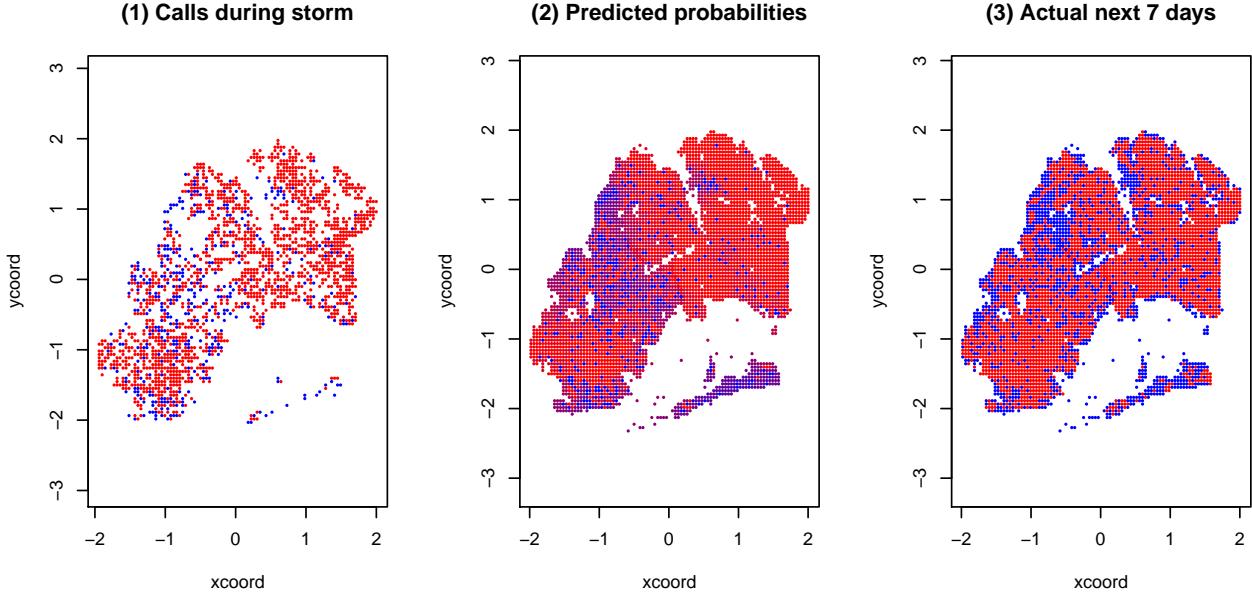


Figure 8: Graphical comparison of actual and predicted areas with reported downed trees. Red indicates at least one tree was reported in a given 0.359 square-mile area

```
plot(test[,2:1], main = "(3) Actual next 7 days",
      col = rgb(test$tree.next7, 0, 1 - test$tree.next7, 1),
      cex = 0.4, pch = 16, asp = 1)
```

Using the predicted probabilities for the test sample, we calculate the TPR and FPR using both a naive cutoff threshold ($p = 0.5$), finding a high TPR but a FPR of almost similar magnitude. Using such prediction would not provide any insight to field crews.

```
tab <- table(test$prob >= 0.5, test$tree.next7)
tpr <- tab[2,2]/sum(test$tree.next7)
fpr <- tab[2,1]/sum(test$tree.next7 == 0)
print(paste0("TPR = ", tpr, ", FPR = ", fpr))

## [1] "TPR = 0.881461969048403, FPR = 0.777083333333333"
```

However, a more informed cutoff based on the first day's probability of a downed tree ($p = 0.8$) yields slightly more balanced results – sacrificing some true positives for far fewer false positives.

```
tab <- table(test$prob >= 0.8, test$tree.next7)
tpr <- tab[2,2]/sum(test$tree.next7)
fpr <- tab[2,1]/sum(test$tree.next7 == 0)
print(paste0("TPR = ", tpr, ", FPR = ", fpr))

## [1] "TPR = 0.71188673032598, FPR = 0.390277777777778"
```

The test model accuracy can also be calculated by taking the Area Under the Curve (AUC) of the Receiving-Operating Characteristic. The ROC calculates the TPR and FPR at many thresholds, that produces a curve that indicates the general robustness of a model. The AUC is literally the area under that curve, which is a measure between 0.5 and 1 where the former indicates no predictive power and 1.0 indicates a perfect model.

In order to visualize the ROC, we will rely on the `plotROC` library, which is an extension of `ggplot2`. We will create a new data frame `input` that is comprised of the labels for the test set `ytest` and the predicted probabilities `test.prob`.

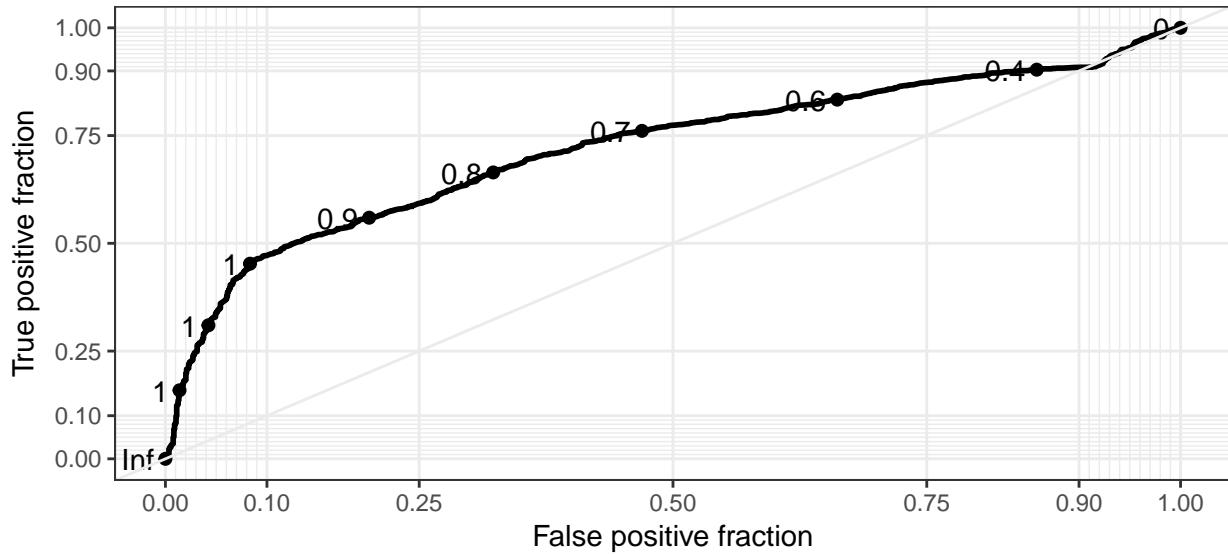


Figure 9: ROC curve out of sample

```
#Load libraries
library(ggplot2)
library(plotROC)

#Set up test data frame
input <- data.frame(ytest = test$tree.next7,
                     prob = test$prob)
```

We then will first create a ggplot object named `base` that will contain the labels (`d` =) and probabilities (`m` =), then create the ROC plot using `geom_roc()` and `style_roc()`. A ROC curve for a well-performing model should sit well-above the the 45 degree diagonal line, which is the reference for an AUC of 0.5 (the minimum expected for a positive predictor). However, as the curve is below the 45 degree line, we may have a seriously deficient model.

```
#Base object
roc <- ggplot(input, aes(d = ytest, m = prob)) +
  geom_roc() + style_roc()

#Show result
roc
```

As estimated using `calc_auc()`, the out-of-sample AUC is 0.721, which is not a bad start. While we are able to fill impute the status of downed trees in other parts of the city, it is helpful to remember that the output of the KNN needs to match the intended use. If a limited number of field crews are deployed, then it may more sense to use the probabilities to prioritize neighborhoods. Otherwise, if additional resources could be hired, then knowing the total number of likely affected areas could inform how much to budget for the downed trees effort.

```
calc_auc(roc)$AUC
## [1] 0.7205748
```

Despite the promising result, we should be cognizant that KNNs generally are not the algorithm of choice of modelers unless there is relatively little data. We should thus ask: *Is there a better classifier?*

1.3.4 Practice Exercises

The US Census Bureau's American Community Survey provides an in-depth view of life in America. One of the many features that are captured in the survey is healthcare coverage. Apply the above methods to predict healthcare coverage in the US State of Georgia in the year 2009.

The data can be obtained here [link to go here]:

1. Randomly split the sample into a 50% training and 50% test set.
2. Predict healthcare **coverage** using continuous variables such as age (**age**) and **wage**.
3. Calculate the performance on the test sample.

1.4 Logistic Regression

NEED A NEW BEGINNING For much of the natural and social sciences, the goal of classification is inference. Inference of how much specific factors are associated with an observed phenomenon – not just prediction. The association typically are furnished with probabilistic qualities that allow an analyst to gauge how certain the pattern is. The output, in turn, lend themselves to building narratives.

The statistically-driven narrative are part of our daily lives. Nowadays, it would not surprise one to hear that a smoker has X-times higher chance of developing cancer than a non-smoker.¹³ ^^MORE STATS EXAMPLES NEEDED HERE^^

These short empirical tid bits are rooted in a method known as *logistic regression*. Like ordinary least squares, logistic regressions are the workhorse of the social and natural sciences for inferring the marginal effects of input factors holding all else constant.

1.4.1 Under The Hood

Before we dive into the particulars, we will cut straight to the chase: logistic regression is best suited for cases where we believe that the decision plane between two or more classes is a straight line. It imposes strong linear assumptions on a problem, meaning that the method may gloss over any pattern that is non-linear or discontinuous. That aside, logistic regression has the ability to contribute to statistically-driven narratives and has a long tradition of probabilistic inference backing it.

Let's assume that the classes of a target variable y can be distinguished using some linear combination of input variables x_1 and x_2 . Upon graphing the features and color coding using the labels, you see that the points are clustered such that purple points represent $z = 1$ and gold points represent $z = 0$.

We can express the relationship between z , x_1 , and x_2 as a linear model similar to OLS:

$$z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

where z is a binary outcome and, like OLS, β_k are coefficients that are learned using a *Maximum Likelihood Estimation* or MLE. The simple idea of MLE is that coefficients can be adjusted so the coefficients can work together to maximize the likelihood of mimicking the target. While the innards of MLE are beyond the scope of this text, for more in-depth treatment on the topic, refer to *Elements of Statistical Learning*¹⁴ or *Introduction to Statistical Learning*.¹⁵

The challenge with binary outcomes in a linear framework is that we run the risk that \hat{y} would exceed the binary bounds of 0 and 1 and would thus make little sense. Imagine if the predicted \hat{y} were -103 or +4 –

¹³https://www.cdc.gov/cancer/lung/basic_info/risk_factors.htm

¹⁴Ref needed

¹⁵Ref needed

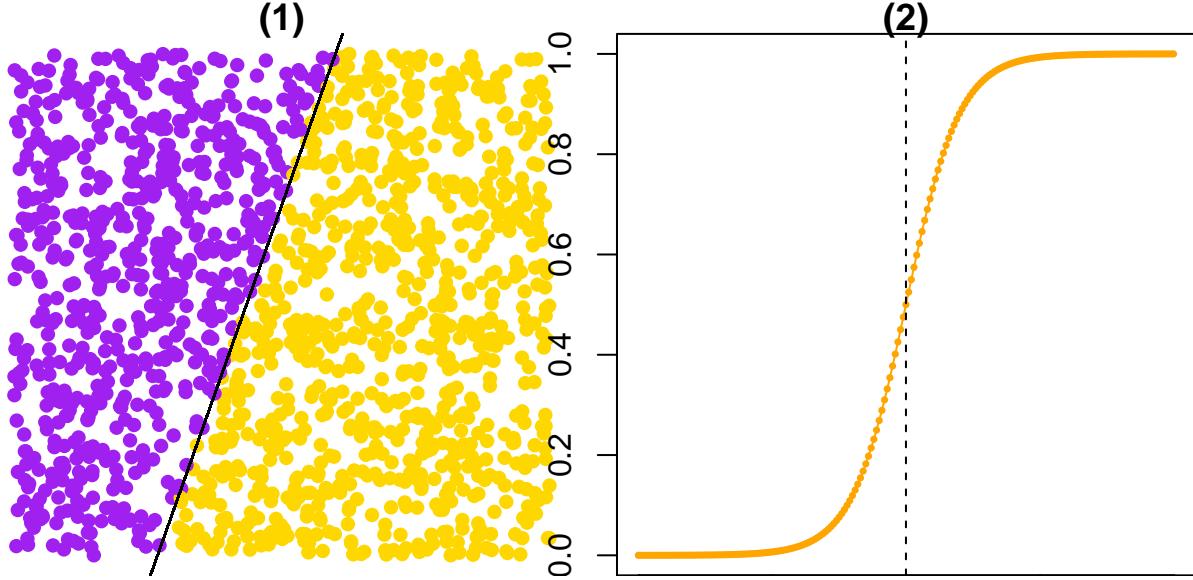


Figure 10: (1) A linearly separable problem. (2) A sigmoid.

What would that mean in the case of a binary variable? How would we interpret a positive value versus a negative? This is an obvious shortcoming of a linear model. Statisticians have cleverly solved the bounding problem by inserting the predicted output into a logistic function:

$$F(z) = \frac{1}{1 + e^{-z}}$$

For a random variable x that takes on values from -10 to +10, the logit transformation converges to +1 where $x > 0$ and to 0 where $x < 0$. This S-shaped curve, or *sigmoid*, bounds \hat{y} to a 0/1 range. By substituting the linear model output z into the logistic function, we bound the output between 0 and 1 and interpret the result as a conditional probability:

$$p = Pr(Y = 1|X) = F(z) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}}$$

This may seem to be a convoluted set of formulas, but it serves a convenient purpose. Unlike many of the techniques in this book, logistic regression is directly interpretable so as long as we believe that the decision boundary is linear. To be able to state that, for example, smokers have a 15 to 30-times higher chance of lung cancer than non-smokers, we use coefficients as a means to contextualize relative effects of input variables. This requires some basic math. The odds of an event are defined as the following:

$$odds = \frac{p}{1-p} = \frac{F(z)}{1-F(z)} = e^z$$

In its purest form, probability p can be calculated without a model, but to hold all covariates constant, we can fit the output of a logistic regression into this framework where $F(z)$ is a probability of some event $z = 1$ and $1 - F(z)$ is the probability of $z = 0$. The odds can be re-arranged as:

$$pr(success) = \frac{e^{(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}}{1 + e^{(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}}$$

$$pr(failure) = \frac{1}{1 + e^{(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}}$$

Typically, we deal with *odds* in terms of *log odds* as the exponentiation may be challenging to work with:

$$\log(\text{odds}) = \log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

where *log* is a natural logarithm transformation. This relationship is particularly important as it allows for conversion of probabilities into odds and vice versa.

The underlying coefficients of the logistic regression can be interpreted using *Odds Ratios* or *OR*. Odds ratios essentially express a marginal unit comparison. Since $\text{odds} = e^z = e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2}$, then we can express an odds ratio as a marginal 1 unit increase in x_1 comparing $\text{odds}(x+1)$ over $\text{odds}(x+0)$:

$$OR = \frac{e^{\beta_0 + \beta_1(x_1+1) + \beta_2 x_2}}{e^{\beta_0 + \beta_1(x_1+0) + \beta_2 x_2}} = e^{\beta_1}$$

After some arithmetic, it turns out the OR is simply equal to e^{β_1} , which can be interpreted as a multiplicative effect or a percentage effect ($100 \times (1 - e^{\beta_1})\%$). More simply, this means that one can obtain the ballpark effect of a regression coefficient by exponentiating it. For example, if a logistic regression were trained to relate wages and citizenship to whether people have health care insurance, the result may look as follows:

$$y(\text{no coverage}) = 0.468 - 0.048 \times \text{wage} + 0.372 \times \text{non-citizen}$$

The coefficients provide little information other than the fact that a positive coefficient is associated with an increase in the chance of y . However, by exponentiating the coefficients, the odds of coverage are as follows for each variable:

- $OR_{\text{wage}} = e^{-0.048} = 0.953$ translates to -4.68% lower chance of not having health coverage for every \$1000 increase in wages. Otherwise stated, more that one earns, better the chance of having health coverage.
- $OR_{\text{non-citizen}} = e^{0.372} = 0.451$ translates to 45% higher chance of not having health coverage among non-citizens.

1.4.2 Tips of the trade

Training a logistic regression is easy to do using a statistical analysis software, but is also easily used incorrectly. Tuning is a matter of finding the optimal combination of variables and generally starts from a hypothesis of how variables are related to a target. As tuning process unfolds, it invariably becomes a series of trial and error tests. While this forces data scientists to start from a theory and hypothesis, it also highlights some of the core challenges of modeling.

Suppose an analyst is tasked with understanding what drives health care insurance coverage in the United States and settles on a four-variable specification that lends itself to a compelling narrative. On the surface, the results sounds promising. However, it may seem less so when considering that the four variables were selected from among 100 possible variables – the results reflect only one of 3.9 million possible four variable specifications. *How does one know if the specification is the best?*

The possibilities of tuning an accurate model are seemingly endless, but are far more manageable when put in the context of the policy objective. As logistic regression is one of the few truly interpretable statistical learning algorithms, the goals can be broadly divided between *interpretability* (focus on $\hat{\beta}$) and *prediction* (focus on \hat{y}).

Table 10: The good and ugly of logistic regression.

Useful Properties	Challenges
Interpretable in terms of odds ratios.	Collinear input variables may cause coefficients to return odd, inexplicable results.
Probabilistic properties allow for hypothesis testing that help build statistically-rooted arguments.	Unable to handle cases where $k > n$
Problems where the distinction between classes is linear	When k is large, finding optimal specification is challenging. Class imbalance makes prediction a non-trivial task.

Interpretability. As we have seen previously, the coefficients lend themselves to articulating how $X \rightarrow Y$ as the ultimate goal is causal inference. Causality is a challenging thing to prove, thus the process of producing reliable and trustworthy parameter estimates takes into account the basic assumptions of logistic regression. Failing to account for the basic assumptions of the technique will lead to misleading results. There are many assumptions, but we focus on a few key considerations:

In virtually all introductory texts that cover logistic regression, *multicollinearity* is flagged and rightfully so. As a refresher, the condition is one in which two or more input variables are not only correlated with the target variables, but amongst each other as well. The consequence is that coefficients will behave oddly. The magnitude of the coefficients may seem unreasonably large. At times, the direction of the relationship is not what one may expect. The answer lies within what the coefficients represent: they are the average effect of x on y , partially isolated holding all else constant. Thus, if two or more variables have identical or very similar information, the algorithm is not able to precisely distill each variable's effects when coefficients are learned during maximum likelihood estimation process. This explains the odd behavior in coefficients and makes coefficients invalid for interpretation. Interestingly, the predictions \hat{y} will still be usable. The best option is to conduct variable selection in advance to minimize double counting of signal. Note that while multicollinearity may interfere with the coefficients, the prediction \hat{y} is still valid.

In information-rich environments, the *ill-posed problem* becomes more common: there may be more variables than observations. Like ordinary least squares, logistic regression cannot be solved when $k > n$. In the social sciences, subject matter intuition is the guiding force in choosing a parsimonious specification. More recently, regularization methods such as Least Absolute Shrinkage and Selection Operator (LASSO) and Ridge Regression have been relied upon to efficiently surface correlated variables. These methods are not without their problems, however. LASSO, for example, can force coefficient of noisy variables to zero, but coefficients do not have standard errors – thus for causal inference is a challenge.

For more in-depth treatment on the subject, consider reviewing [Need recommendations here].

Prediction. Unlike interpretation, prediction is purely focused on accuracy. Automated techniques are often employed to conduct variable selection, but there is one core modeling consideration that matters significantly for logistic regression: *class imbalance*. The idea of class imbalance is simple – when the classes in target y are not represented in relatively even proportions.¹⁶ We can imagine that many problems are imbalanced:

- fire-related fatalities (only 3,400 fire deaths in 2017¹⁷ vs. 2.7 million total deaths in the US in 2016¹⁸)
- homelessness (approximately 0.17-percent of the US population¹⁹)
- health care coverage (8.7% of US population is uninsured in 2017²⁰)

The smaller, less represented class is referred to as the *minority class*. While a logistic regression can be estimated, the size of the predicted values \hat{y} will typically hover around average probability that an event

¹⁶ref needed

¹⁷<https://www.usfa.fema.gov/data/statistics/>

¹⁸<https://www.cdc.gov/nchs/data/hus/2017/019.pdf>

¹⁹<https://www.bbc.com/news/world-us-canada-42248999>

²⁰<https://www.census.gov/data/tables/time-series/demo/health-insurance/hic.html>

will occur. Thus, if a model is trained to predict homelessness, a representative sample of the US will yield $\hat{y} = 0.17$ as the intercept in the linear model will have shrunk.²¹ This, in turn, makes evaluating accuracy more challenging to evaluate as the classification threshold is no longer $c = 0.5$. There is not one solution to addressing the class imbalance problem, but rather many possible avenue. Here are a few common strategies:

- *AUC*: If the goal is to rank observations (e.g. prioritization of risky buildings), consider using the AUC derived from the ROC curve as it summarizes the overall accuracy across multiple thresholds.
- *Adjusted cutoff*: If the goal is to definitively flag observations for closer inspection, then it is worth considering to adjust the classification threshold c to a value that optimizes desired accuracy. For example, setting the threshold to the average probability can better balance the TPR and FPR. In these cases, it also becomes necessary to assign value to how much more or less a false positive is worth compared with one additional true positive prediction.
- *Over-sampling/Under-sampling*: If the goal is to produce a balanced prediction, there is a whole host of methods for adjusting the composition of the sample. The simplest is over-sampling in which the minority class observations are replicated until it is of equal proportion as the majority class. The inverse is also done: under-sampling the majority class so that it is of equal proportion to the minority. In either case, the threshold c can remain at 0.5, but there is a higher chance of overfitting.²² Alternatively, Synthetic Minority Over-Sampling Technique (SMOTE) produces new observations that look like the minority class that add fresh signal to the sample while minimizing overfitting.²³

1.4.3 DIY: Expanding Health Care Coverage

Background. Universal healthcare has become a basic human right in many countries. In the United States, this is not currently a guarantee, shrouded in heated political debate and controversy whether its a matter of human rights or a matter in which an individual may choose his or her fate. Regardless of the politics, there is data on healthcare coverage.

According to the American Community Survey ACS, an annual survey of approximately 3.5% of the US population as conducted by the US Census Bureau, over 22.4% of residents of the U.S. state of Georgia were without healthcare coverage in 2009. That is a fairly sizable proportion of the population – for every ten people, between two to three did not have coverage. If you read the news in 2010, a new law to provide affordable healthcare came into effect to help the uninsured.

Imagine that you are hypothetically tasked with getting the word out and drive recruitment in the state of Georgia. There is a hiccup, however. While commercial registries exist with people's demographic and personal contact information, most statistics on coverage are based on surveys, thus we do not precisely know *who* does not have insurance. A brute force approach could be to reach out to everyone under the sun though we can easily infer a wasted effort as 776 of every 1000 people are already covered. *How do we get to the 224 people who are not already insured?* For marketers, this is a classic targeting problem.

Data needs to enable the prediction and classification of a population into two classes: covered and not covered. By correctly classifying people as covered and not covered, decision makers and outreach staff can mobilize targeted outreach. From a data science perspective, the real objective is to be able to identify and replicate re-occurring patterns in the training data, then generalize the insights onto a sample or population that is not contained in the sample.

Given the label $y(\text{Coverage})$, we can use logistic regression to not only infer what is associated with coverage, but also train a model to prioritize who should be contacted about receiving coverage:

$$y(\text{Coverage}) = f(\text{Sex, Age, Education, Marital Status, Race, Citizenship})$$

²¹Ref from Gary King from Harvard?

²²ref required

²³<https://jair.org/index.php/jair/article/view/10302>

Data. Based on the Census American Community Survey (ACS), we will illustrate how to construct a logistic regression. First, we will load the ACS sample, which has been balanced – meaning that both covered and non-covered survey respondents are represented in equal proportions in the sample.

```
# Load ACS health care data
load("data/acs_health.Rda")
```

Next, we will simulate a train and test to determine how well our models can do when applied in the wild.

```
# Randomly assign
set.seed(100)
rand <- runif(nrow(health)) > 0.5

# Create train test sets
train <- health[rand == T, ]
test <- health[rand == F, ]
```

Function. Training a logistic regression can be easily performed using the `glm()` function, which is a flexible algorithm class known as Generalized Linear Models. Using this one method, multiple types of linear models can be estimated including ordinary least squares for continuous outcomes, logistic regression for binary outcomes and Poisson regression for count outcomes.

At a minimum, three parameters are required:

```
glm(formula, data, family)
```

where:

- `formula` is a formula object. This can take on a number of forms such as a symbolic description (e.g. $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$ is represented as $y \sim x_1 + x_2$).
- `data` is a data frame containing the target and inputs.
- `family` indicates the probability distribution used in the model. Distributions typically used for GLMs are *binomial* (binary outcomes), *poisson* (count outcomes), *gaussian* (continuous outcomes - same as OLS), among others.

The `family` refers to the type of probability distributions that underlie each estimation method. In the case of logistic regression, the probability family is *binomial*.

Focusing on interpretation. In the social sciences and in public policy, the focus of regression modeling is typically placed on identifying an effect or an associated relationship that describes the process being studied. To tease out the effects, logistic regression analyses tend to center around what is colloquially known as a *build up*. Blocks of conceptually-related variables are tested in order to infer their important and effect on the target outcome. As we are focused on insurance coverage, we specify four models:

- *Personal:* $coverage = f(\log(age) + race + sex)$
- *Economic:* $coverage = f(wage + employment)$
- *Social:* $coverage = f(citizenship + marital + schooling)$
- *Combined:* $coverage = f(\log(age) + race + sex + wage + employment + citizenship + marital + schooling)$

It is worth highlighting that the inclusion of certain characteristics, such as race and sex, can improve the accuracy of a model. However, when used for prediction – particularly for targeting – one needs to consider the ethics and fairness of including these variables in application.

```
# Estimated GLM models
glm_pers <- glm("no.coverage ~ log(age) + race + sex",
                 data = train, family = binomial)

glm_econ <- glm("no.coverage ~ wage + esr",
                  data = train, family = binomial)
```

```

glm_soc <- glm("no.coverage ~ cit + mar + schl",
                data = train, family = binomial)

glm_all <- glm("no.coverage ~ log(age) + wage + schl + esr + cit + mar + race + sex",
                 data = train, family = binomial)

```

The regression outputs chronicles every pertinent aspect of each model, such as the direction of the relationships (e.g. positive or negative weights), their statistical significance (e.g. p-value or t-statistics), and the relative fit of the model (e.g. the lowest Akaike Information Criterion or AIC provides *relative* model fit comparison). For example, an analyst may point out that education has an effect on coverage by interpreting the coefficient point estimates. In the combined model, education attainment coefficients are estimated relative to people who hold a graduate degree, thus indicating that people who :

- did not finish high school have a *6.58-times* higher chance of not having health coverage ($e^{w = 1.884} = 6.58$)
- hold a high school degree have a *4.91-times* higher chance of not having health coverage ($e^{w = 1.592} = 4.91$)
- hold a college degree are relatively better off than the previous two groups with a *1.79-times* higher chance of not having health coverage ($e^{w = 0.584} = 1.79$)

All coefficients are statistically significant. While it is valid to evaluate models on this basis, it is necessary to remember that this is not the same as evaluating a model for predictive use cases as predictive accuracy is not assessed on the basis of coefficients.

```
% Table created by stargazer v.5.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu
% Date and time: Fri, Oct 26, 2018 - 21:33:56
```

```
% Table created by stargazer v.5.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu
% Date and time: Fri, Oct 26, 2018 - 21:33:56
```

Focus on prediction. Prediction, in contrast, emphasizes absolute generalizable accuracy and requires model validation techniques like k-cross validation. We can rely on the `boot` library to generate cross-validated accuracy estimates through the `cv.glm()` function:

```
cv.glm(data, glmfit, cost, K)
```

where:

- `data` is a data frame or matrix.
- `fit` is a `glm` model object.
- `cost` specifies the cost function for cross validation.
- `K` is the number of cross validation partitions.

Note that the cost function needs to take two vectors. The first is the observed responses and the second is the predicted responses. For example, the cost function could be the overall accuracy rate (*ACC*):

$$ACC = \frac{FP + FN}{TP + FP + TN + FN}$$

as translated into R code below, assuming that the threshold is the sample probability of being uncovered.

```
costAccuracy <- function(y, yhat){
  #
  # Returns vector of accuracy, TPR, and FPR using sample
  #
  # Args:
  #   y = vector of actual target
  #   yhat = vector of predicted probabilities
  #
```

Table 11:

	<i>Dependent variable:</i>			
	Personal	Economic	Social	All
	(1)	(2)	(3)	(4)
log(age)	-1.216*** (0.032)			-0.698*** (0.039)
raceAmer. Ind.	10.285 (80.815)			10.406 (82.664)
raceAsian	9.789 (80.814)			9.455 (82.664)
raceBlack	9.916 (80.814)			10.217 (82.664)
raceNat. Hawaiian/Pac. Isl.	9.210 (80.817)			9.323 (82.666)
raceOther	10.860 (80.814)			10.490 (82.664)
raceWhite	9.397 (80.814)			9.872 (82.664)
sex	-0.327*** (0.028)			-0.278*** (0.029)
wage		-0.001*** (0.0001)		-0.001*** (0.0001)
esr		0.021*** (0.006)		-0.012* (0.007)
cit			0.387*** (0.011)	0.363*** (0.013)
mar			0.326*** (0.008)	0.201*** (0.010)
schl			-0.100*** (0.004)	-0.108*** (0.004)
Constant	-6.171 (80.814)	-1.527*** (0.022)	-1.248*** (0.070)	-7.681 (82.664)
Observations	37,142	37,142	37,142	37,142
Log Likelihood	-16,410.350	-17,632.960	-15,759.410	-15,230.020
Akaike Inf. Crit.	32,838.710	35,271.920	31,526.820	30,488.050

Note:

*p<0.1; **p<0.05; ***p<0.01

Table 12:

Coefficient table of four alternative logistic regression specifications.

```
#Components
ymean <- mean(y)
p <- sum(y == 1)
n <- sum(y == 0)
tp <- sum((y == 1) & (yhat >= ymean))
fp <- sum((y == 0) & (yhat >= ymean))
tn <- sum((y == 0) & (yhat < ymean))
tot <- length(y)

#scores
acc <- (tp + tn) / tot
tpr <- tp / p
fpr <- fp / n

return(c(acc, tpr, fpr))
}
```

Using `cv.glm()`, we test each of the specifications using 10-folds cross validation to obtain an accurate estimate on model performance

```
# Load boot library
library(boot)

# Estimate k-folds
pers <- cv.glm(data = train, glmfit = glm_pers, cost = costAccuracy, K = 10)
econ <- cv.glm(data = train, glmfit = glm_econ, cost = costAccuracy, K = 10)
soc <- cv.glm(data = train, glmfit = glm_soc, cost = costAccuracy, K = 10)
all <- cv.glm(data = train, glmfit = glm_all, cost = costAccuracy, K = 10)
```

The accuracy results are stored in the `cv.glm()` objects, specifically in the first three elements in the `delta` component (e.g. `pers$delta`). The order of the model accuracies are not different than what can be inferred from the AICs in the regression table: The fully loaded specification yields the best overall predictive accuracy.

Table 13: Cross-validated prediction accuracies.

Model	Accuracy	TPR	FPR
Personal	0.6602499	0.6126950	0.3282049
Economic	0.4857574	0.5724453	0.5295883
Social	0.6680039	0.6718639	0.3326993
Fully Loaded (All)	0.7187550	0.6588178	0.2675087

Scoring. In order to obtain the score the test set, we use `predict()`:

```
predict(object, newdata, response)
```

where:

- `object` is a GLM model object.
- `newdata` is a data frame. This can be the training data set or the test set with the same format and features as the training set.
- `response` indicates the type of value to be returned, whether it is the untransformed “link” or the probability “response”.

We will now apply `predict()` to score the responses for each `train` and `test` samples.

```

pred.glm.train <- predict(glm_all, train, type = "response")
pred.glm.test <- predict(glm_all, test, type = "response")

```

Lastly, to calculate the prediction accuracy, we will once again rely on the combination of `ggplot2` and `plotROC` libraries for the AUC. Interestingly, the test set AUC is greater than that of the train set. This occurs occassionally and is often times due to the luck of the draw. ^^NEED A BIT MORE TO CLOSE OUT THE SECTION

```

#load libraries
library(plotROC)
library(ggplot2)

#Set up ROC inputs
input.glm <- rbind(data.frame(model = "train",
                               d = train$no.coverage,
                               m = pred.glm.train),
                     data.frame(model = "test",
                               d = test$no.coverage,
                               m = pred.glm.test))

#Graph all three ROCs
roc.glm <- ggplot(input.glm, aes(d = d, model = model, m = m, colour = model)) +
  geom_roc(show.legend = TRUE) + style_roc()

#AUC
calc_auc(roc.glm)[,2:3]

##   group      AUC
## 1     1 0.7690835
## 2     2 0.7716269

```

1.4.4 Practice Exercises

1. Can logistic regression be applied to the downed tree problem from the KNN section? Apply the method to the downed trees data. How do the accuracies compare and why?

1.5 LASSO regression

[expand]

1.6 Decision Tree Learning

Mobile technologies have lowered the bar to using lightweight sensors that measure the physical world and have opened new applications of data in daily life. From a smart phone's accelerometer, it's possible to track distinct patterns in one's activity based on the fluctuations in acceleration ($\frac{m}{s^2}$). In fact, many of these technologies have become commonly available, enabling physical fitness activity monitoring to characterize transportation quality. Below is a set of exercise measurements from an smartphone accelerometer that lasted approximately 6.5 minutes and graphed at a frequency of 5 hertz (five readings per second).

Can you visually identify distinct patterns? What makes those patterns distinct?

It becomes immediately apparent that the methods covered thus far are not suitable for the task at hand. If we manually extract samples from these periods, we can quantify the patterns in terms of their central

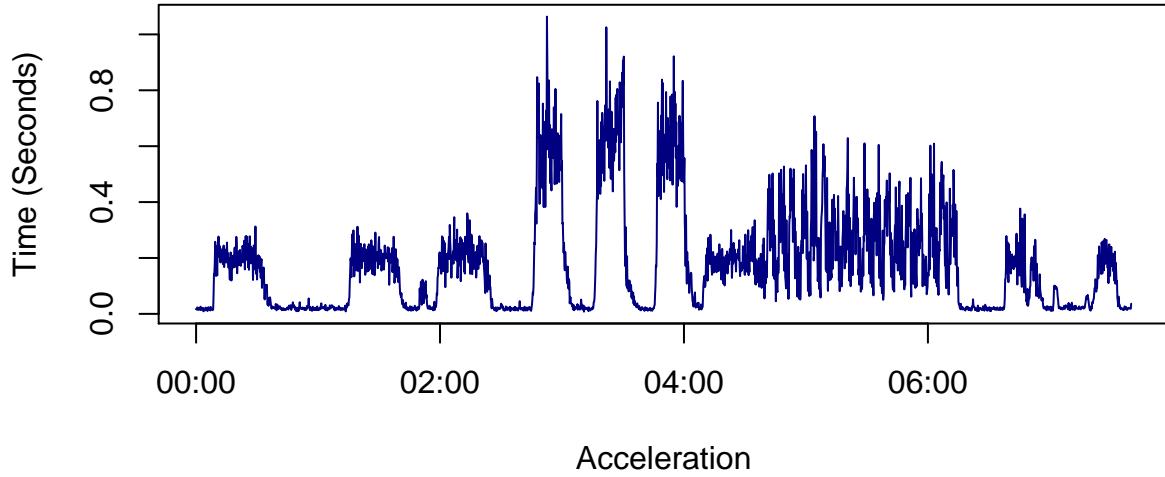


Figure 11: Accelerometer data collected from an iPhone.

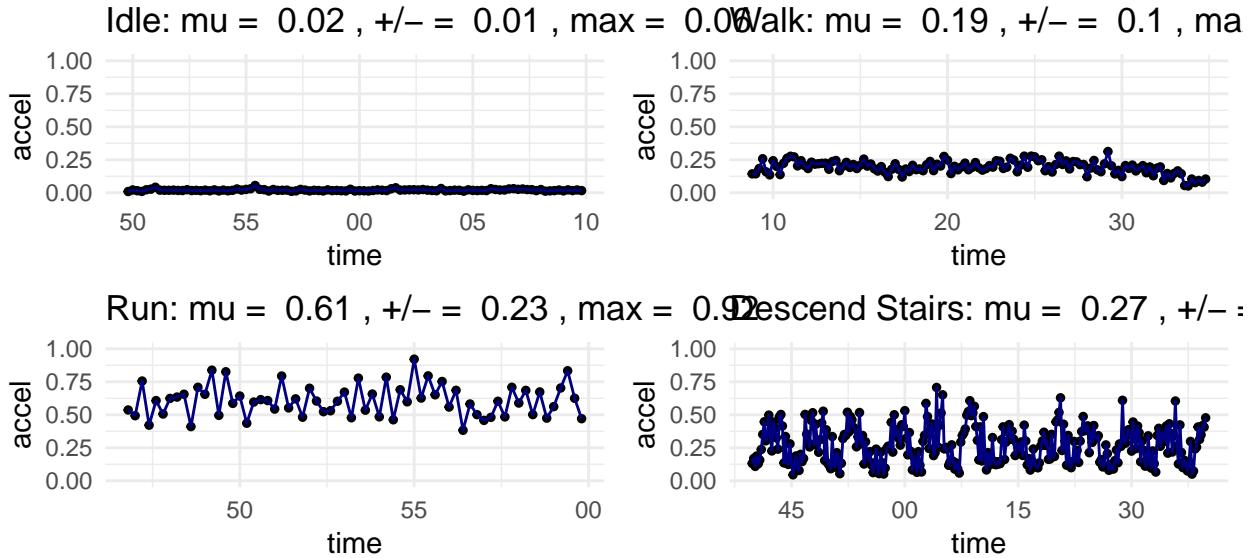


Figure 12: Acceleration by type of movement.

tendencies. Idle periods have near zero acceleration, walking periods have acceleration around 0.2 with tight dispersion, running periods hover around 0.6 +/- 0.2, and descending stairs vary widely. Using simply the level of acceleration may not accurate as at least two types of motion have overlapping distributions.

Decision tree learning can help bring clarity. Trees are designed to look at inputs and partition the sample into smaller more homogeneous cells with respect to the target. This recursive partitioning allows a tree to resemble an inverted tree: moving away from the base of the tree, the tree trunk splits into two or more large branches, which then in turn split into even smaller branches, eventually reaching even small twigs with leaves.

Decision trees use recursive partitioning to learn patterns, doing so using central concepts of *information theory*. There are a number of decision tree algorithms that were invented largely in the 1980s and 1990s, including the ID3 algorithm, C4.5 algorithm, and Classification And Regression Trees for Machine Learning (CART). All these algorithms follow the same framework that includes the following elements: (1) nodes and edges, (2) attribute tests, and (3) termination criteria.

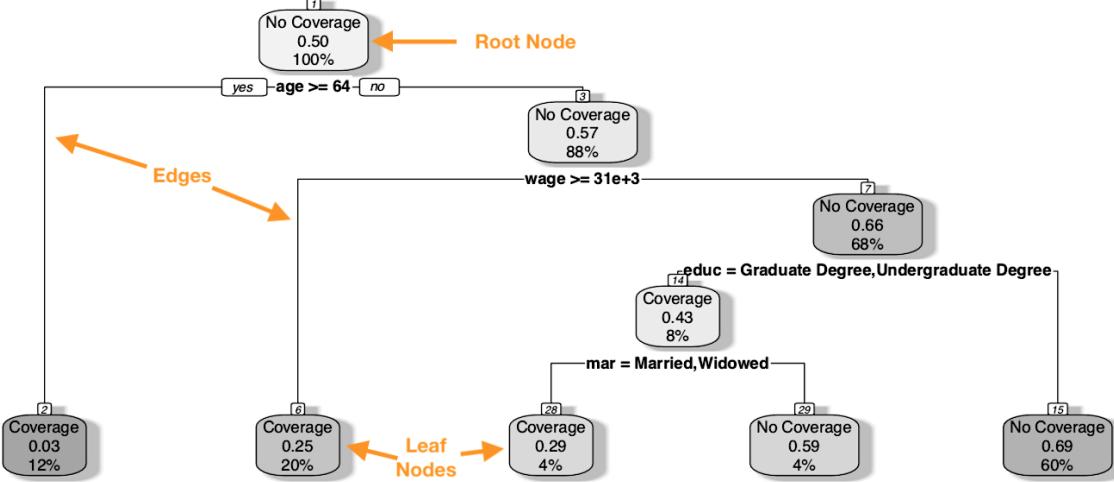


Figure 13: Anatomy of a Decision Tree.

1.6.1 Under the hood

Anatomy of a decision tree. The tree is comprised of nodes and edges. Nodes (circles) contain records. Edges (lines) show dependency between nodes and is the result of an *attribute test* – or a process that finds the optimal criterion to subset records into more homogeneous groups of the target variable. The node at the top of the tree is known as the *root* and represents the full population. Each time a node is split, the result is two nodes – each of which is referred to as a *child node*. A node without any child nodes is known as a *leaf*. The node is labeled using majority voting based on whichever class is most represented. The goal is to grow a tree from the root node into as many smaller child nodes that contain more of one class than another.

Decision trees split nodes based on finding thresholds along the input variables. There can be seemingly infinite number of potential variable-threshold combinations – which is best? Drawing from *information theory*, we can apply an *information gain* formula to evaluate all candidate splits and find one that provides the most information. This optimal split yields to more homogeneous child nodes, which in turn can be split even further. The search for the best threshold is known as an *attribute test*.

As we can see in the example decision tree for health care insurance, each node is connected to at least one other node. Starting at the root node, we can see that overall, the population is labeled “no coverage” based on the decimal percentage 0.5. The 100% indicates the proportion of the sample that is contained at the node. Below is `age >= 64`, which is the most informative attribute test that is used to split. To the left, the edge leads to another node at the bottom left corner of the diagram, which contains people who are age 64 or older. While the leaf node only contains 12% of the entire sample, it is almost exclusively people who have health care coverage. To the right, the remaining 88% of the sample, which is further split by wage and other variables. Each leaf node is defined as the intersection of multiple binary criteria, giving way to profiles of users that can be easily segmented.

Growing trees. There are a number of decision tree algorithms that were invented largely in the 1980s and 1990s, including the ID3 algorithm, C4.5 algorithm, and Classification And Regression Trees for Machine Learning (CART). The process is fairly clear cut and iterative:

1. **Base Cases.** The process starts with checking for “base cases” at the root node, the idea being that it might not be worth exerting effort to grow the tree if the data do not support it. The algorithm will first check to see if (a) all values of the target are of one class, and (b) none of the input variables offer any useful information. There are other base cases to consider depending on the algorithm, but any true base case will result in stopping the algorithm and returning only the root node.

- 2. Recursive Partitioning.** If none of the base cases are true, the algorithm proceeds to attribute testing using either *Information Gain* or *Gini Impurity* as will be covered in the following section. At the root node, if there are 10 variables with each 30 possible thresholds, attribute tests are applied 300 times choosing the candidate threshold that yields the most homogeneous child nodes. Upon splitting, attribute tests are applied to each child node, making this a recursive partitioning procedure.
- 3. Stopping Criteria versus Pruning.** At some point, the algorithm needs to stop. The question is *when*? One way is to grow the tree until some *stopping criteria* are met, such as if a leaf has fewer records than a pre-specific threshold, the purity or information gain falls below a pre-specified level, or if the tree has grown to n-number of levels (e.g. number of rows of splits). While stopping criteria are useful, the results in some studies indicate their performance cap the tree from reaching its full predictive potential.²⁴ The alternative approach involves growing a tree to its fullest, then comparing the prediction performance given tree complexity (e.g. number of nodes in the tree) using cross-validation. In the example graph below, model accuracy degrades beyond a certain number of nodes. Thus, optimal number of nodes is defined as when cross-validation samples (e.g. train/test, k-folds) reaches a minimum across samples. Upon finding the optimal number of nodes, the tree is *pruned* to only that number of nodes.

Attribute Tests. Information gain is a form of *entropy* that measures the consistency of information. Based on these distinct states of activity, entropy is defined as:

$$\text{Entropy} = \sum_{i=1}^k -p_i \log_2(p_i)$$

where i is an index of states, p is the proportion of observations that are in state i , and $\log_2(p_i)$ is the Base 2 logarithm of the proportion for state i . Information Gain (IG) is variant of entropy, which is the entropy of the root node *less* the average entropies of the child nodes.

$$\text{IG} = \text{Entropy}_{\text{root}} - \text{Avg Child Entropy}$$

How does this work in practice? Starting from the root node, we need to calculate the root entropy, where the classes are based on the classes of the target `usership`.

$$\begin{aligned} \text{Entropy}_{\text{usership}} &= (-p_{\text{user}} \log_2(p_{\text{user}})) - (-p_{\text{non-user}} \log_2(p_{\text{non-user}})) \\ &= \left(-\frac{6}{12} \log_2\left(\frac{6}{12}\right)\right) + \left(-\frac{6}{12} \log_2\left(\frac{6}{12}\right)\right) \\ &= 1.0 \end{aligned}$$

Then, the attribute test is applied to the root node by calculating the weighted entropy for each proposed child node. Using the `income` feature, the calculation is as follows:

- Split the root node into two child nodes using the `income` class. This yields the following subsamples as shown in the table below:

	< \$20k	> \$20k
No	0	6
Yes	5	1
Total	5	7

- For each child node (the columns in the table), calculate entropy:

$$\begin{aligned} \text{Entropy}_{\text{income} < 20k} &= (-p_{\text{user}} \log_2(p_{\text{user}})) - (-p_{\text{non-user}} \log_2(p_{\text{non-user}})) \\ &= -\frac{5}{5} \log_2\left(\frac{5}{5}\right) = 0 \end{aligned}$$

²⁴

$$\begin{aligned}\text{Entropy}_{\text{income} > 20k} &= (-p_{\text{user}} \log_2(p_{\text{user}})) - (-p_{\text{non-user}} \log_2(p_{\text{non-user}})) \\ &= -\frac{6}{7} \log_2(\frac{6}{7}) + -\frac{1}{7} \log_2(\frac{1}{7}) = 0.5916728\end{aligned}$$

- Calculate the weighted average entropy of children:

$$\text{Entropy}_{\text{income split}} = \frac{5}{12}(0) + \frac{7}{12}(0.5916728) = 0.3451425$$

- Then calculate the information gain:

$$\begin{aligned}\text{IG}_{\text{income}} &= \text{Entropy}_{\text{root}} - \text{Entropy}_{\text{income split}} \\ &= 1 - 0.3451425 = 0.6548575\end{aligned}$$

- We then can perform the same calculation on all other features (e.g. employment, part of town) and compare results. The goal is to *maximize* the IG statistic at each decision point. In this case, we see that income is the best attribute to use for splitting. This split is easily interpretable: “The majority of users of health services can be predicted to earn less than \$20,000.”

Measure	IG
Employment	0.00
Income	0.6548575
Area of Town	0.027119

Gini Impurity is closely related to the entropy with a slight modification.

$$\text{Gini Impurity} = \sum_{i=1}^k p_i(1-p_i) = 1 - \sum_{i=1}^k p_i^2$$

Using Gini Impurity as an attribute test is also similar to Information Gain.

$$\text{Gini Gain} = \text{Gini}_{\text{root}} - \text{Weighted Gini}_{\text{child}}$$

1.6.2 Tips of the trade

Like any technique, decision trees have strengths and weaknesses. Unlike logistic regression and KNN, decision trees can conduct automated variable selection on any type of data type. In addition, the recursive partitioning produce a tangible definition for each subpopulation represented by a node. The splitting mechanism make it possible to capture interactions and non-linearities that are otherwise not easily accounted for in the previous methods. The implications are that two or more input variables can be blended together to find cells of activities that are otherwise overlooked.

There are detractors, however. Trees can be grown so deeply that there are too many subpopulations to articulate. If left unpruned, terminal leafs may give a false impression of accuracy and precision – the small samples may give a false impression. Nonetheless, decision tree learning is an important contribution to classification problems and form the basis of many other algorithms.

Table 16: The good and ugly of decision trees.

Useful Properties	Challenges
Rules (e.g. all the criteria that form the path from root to leaf) can be directly interpreted.	Data sets with large number of features will have overly complex trees that, if left unpruned, may be too voluminous to interpret.

Useful Properties	Challenges
Method is well-suited to capture interactions and non-linearities in data.	Trees tend to overfitted at the terminal leafs when samples are too small.
Technique can accept both continuous and continuous variables without prior transformation.	Feature selection is conducted automatically.

1.6.3 DIY: Predicting activity from phone accelerometers

Overview. Accelerometer data may become more common place in the future. In many respects, one's regular movements like walking and running can be viewed as a finger print. Thus, accelerometer data has many uses from measuring and improving health outcomes to multi-factor authentication²⁵ and roadway quality²⁶

It is quite easy to record and use the data. For this DIY, we use data collected on an iPhone at 20 hertz, or 20 recordings per second and will attempt to build our own activity classification algorithm that are calibrated to a specific person's motions. There are two pre-partitioned data frames – one for training ($n = 311174$) that was captured under controlled conditions and an extended test sample taken in non-controlled conditions ($n = 35486$).

```
#Load data
load("data/accel_activity.Rda")
```

Our primary target is the `activity` variable that includes labels for *idle*, *run*, *stairs*, and *walk*. For a real-world application, the data should be augmented with other activities as algorithms only know what they have previously seen.

```
#Quick summary
table(train$activity)

##
##    idle    run stairs   walk
## 193408  21063   7733  88970
```

Smartphone accelerometers captures data along three axes. As the position of the phone can shift, data from any one axis may not be useful without considering measures on orientation. While rotation is captured by the sensors, we simply acceleration into total acceleration $accel = \sqrt{accel\ x^2 + accel\ y^2 + accel\ z^2}$. But if we recall the acceleration chart, we could deduce that a single value of acceleration does not likely have enough context to make for an informed prediction. We augment this data by engineering a few additional variables, namely the average, maximum, minimum, and standard deviation of total acceleration for the last 2.5 seconds (or 50 measurements).

Training. There are a number of R implementations of decision trees, the most popular of which is the `rpart` library:

```
library(rpart)
```

The main function within the library comes with flexible capabilities to induct decision trees:

```
rpart(formula, method, data, cp, minbucket, minsplits)
```

where:

- `formula` is a formula object. This can take on a number of forms such as a symbolic description (e.g. $y = f(x_1, x_2, \dots)$) is represented as "y ~ x1 + x2".

²⁵<https://ieeexplore.ieee.org/document/8090407>

²⁶<https://www.wired.com/insights/2014/03/potholes-big-data-crowdsourcing-way-better-government/>

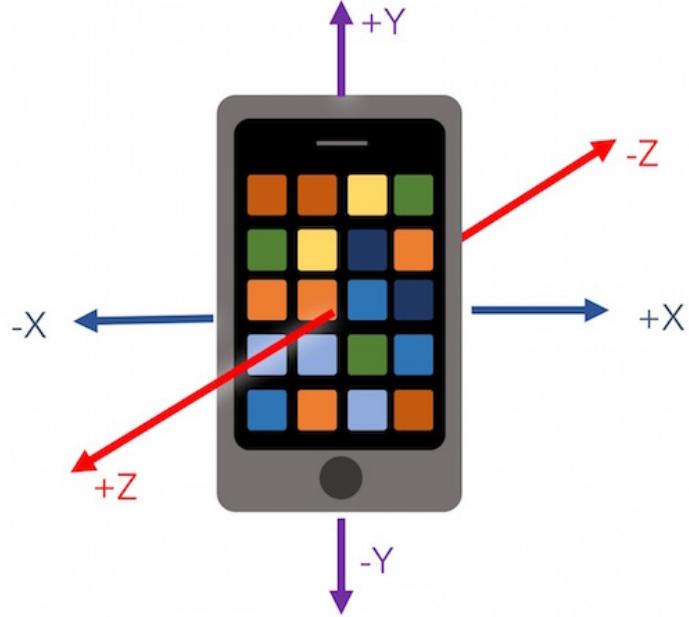


Figure 14: Accelerometer and its three axes.

- `method` indicates the type of tree, which are commonly either a classification tree “class” or regression tree “anova”. Split criteria can also be custom written.
- `data` is the data set in data frame format.
- `cp` is a numeric indicates the complexity of the tree. `cp = 1` is a tree without branches, whereas `cp = 0` is the fully grown, unpruned tree. If `cp` is not specified, `rpart()` defaults to a value of 0.01.
- `minbucket` is a stopping criteria that specifies the minimum number of observations in any terminal leaf.
- `minsplit` is a stopping criteria that specifies the number of observation in a node to qualify for an attribute test.

As a first pass, we’ll run `rpart()` setting `cp = 0`, meaning that the tree will be fully grown without any stopping criteria applied.

```
fit <- rpart(factor(activity) ~ accel + avg50 + max50 + min50 + sd50,
              method = "class", data = train,
              cp = 0)
```

The `fit` object captures all of the inner workings of the decision tree. For example, just plotting the `fit` object will show the full depth of the tree. More importantly is the cross validation results collected at each level of additional complexity. Using the `printcp()` function, we can extract the *CP table*, which contains various accuracy measures associated with each value of the tree complexity value `cp`, including:

- the number of splits `nsplit`,
- the prediction error in the training data `rel error`,
- the cross-validation error `xerror`, and
- the standard error `xstd`.

```
printcp(fit)
```

```
##          CP nsplit rel error     xerror      xstd
## 1 0.6944788819      0 1.0000000 1.0000000 0.0022973417
## 2 0.1768252297      1 0.3055211 0.3056315 0.0015149443
```

```

## 3 0.0021058710      2 0.1286959 0.1288572 0.0010202061
## 4 0.0011803067      4 0.1244841 0.1248408 0.0010049824
## 5 0.0009170728      6 0.1221235 0.1228113 0.0009971820
## 6 0.0005745858      7 0.1212065 0.1220047 0.0009940607

```

Tuning. How do we find the optimal tree depth? First, find the lowest cross-validation `xerror`, then find the tree that has the lowest number of splits that is still within one standard deviation `xstd` of the best tree²⁷. The idea behinds this rule of thumb takes advantage of uncertainty: the true value lies somewhere within a confidence interval, thus any value within a tight confidence interval of the best value is approximately the same. In this first model, the best tree has `nsplit` = 1040 and `xerror` = 0.0284632236808587. By applying the rule, the upper bound of acceptable error is `xerror` = 0.028463 + 0.000489 = 0.0289521914772729. As it turns out, the tree with `nsplit` = 986 is within one standard deviation and is thus the best model.

In other words, the following function can extract the optimal `cp` value.

```

bestCP <- function(fit_obj){
  # Returns best CP val within 1 SD of lowest xerror
  #
  # Args:
  #   fit_obj: decision tree object
  #

  #Find lowest error and associated xstd
  best.error <- min(fit$cptable[,4])
  best.sd <- fit$cptable[,5][which(fit$cptable[,4]==best.error)]

  #Pull CP closest to lower bound
  lower.bound <- best.error + best.sd
  opt.select <- fit$cptable[,1][which(fit$cptable[,4] <= lower.bound)][1]

  return(opt.select)
}

```

Now, we can prune the tree using the optimal `cp` value, then score both the test set. As a comparison point, we will also apply the unpruned model as well.

```

#Get best CP
best.value <- bestCP(fit)

#Prune tree
fit.opt <- prune.rpart(fit, cp = best.value)

#Score, returning probabilities
pred.full <- predict(fit, test, type = 'class')
pred.opt <- predict(fit.opt, test, type = 'class')

```

Evaluating performance. One of the fascinating aspect of decision trees is that it is interpretable in a very different way than logistic regression. In lieu of a thorough review of the learned rules, we may rely on a measure of variable importance, that is defined as follows:

REVISIT AT SOME POINT

$$\text{Variable Importance}_k = \sum \text{Goodness of Fit}_{\text{split}, k} + (\text{Goodness of Fit}_{\text{split}, k} \times \text{Adj. Agreement}_{\text{split}})$$

Where *Variable Importance* for variable k is the sum of *Goodness of Fit* (e.g. Gini Gain or Information

²⁷Hastie et. al (2001)

Gain) at a given split involving variable k . In other words, a variable's importance is the sum of all the contributions variable k makes towards predicting the target. Below, we can see that the measure can be extracted from the `fit.opt` object. As may be expected, `accel` is not the main contributor to predictions, but rather measures of the maximum, mean and variability of acceleration. This also implies that the model could be further tuned by trying different windows for producing the engineered variables – perhaps shorter or longer windows could be even more important.

```
#Extract variable importance list from fit object
fit.opt$variable.importance
```

```
##    max50    avg50    sd50    min50    accel
## 144692.6 143318.8 139401.1 119687.2 117792.5
```

Unlike previous examples, there are four classes to predict, meaning that the AUC measure may be less effective in producing a singular summary. We consider the F1-statistic defined as:

$$F_1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

where $\text{precision} = \frac{TP}{TP+FP}$ and $\text{recall} = \frac{TP}{TP+FN} = \frac{TP}{\text{Actual}(+)}$. The score is out of 1, where a value of 1 is a very well-balanced prediction. To accommodate three or more classes, the Mean F1-Score is the weighted by each activity class k (e.g. run, walk):

$$\text{Mean}F_1 = \sum_{i=1}^k \left(\frac{n_k}{n} F_{1k} \right)$$

These formulae effectively boil down to the following function:

```
meanF1 <- function(actual, predicted){
  # Mean F1 score function
  #
  # Args:
  #   actual = a vector of actual labels
  #   predicted = predicted labels
  #

  #Loop through each class and calculate components
  classes <- unique(actual)
  results <- data.frame()
  for(k in classes){
    results <- rbind(results,
                      data.frame(weight = sum(actual == k, na.rm = T)/length(actual),
                                 precision = sum(predicted == k & actual == k)/sum(predicted == k, na.rm = T),
                                 recall = sum(predicted == k & actual == k, na.rm = T)/sum(actual == k, na.rm = T)))
  }
  results$score <- results$weight * 2 * (results$precision * results$recall) / (results$precision + results$recall)
  return(sum(results$score, na.rm=T))
}
```

When applied to the function to the predictions (`pred.opt` and `pred.full`), we find that the mean F1-statistics reached 0.8533023 and 0.8527223 – not bad for a first cut, but certainly can benefit from extra attention.

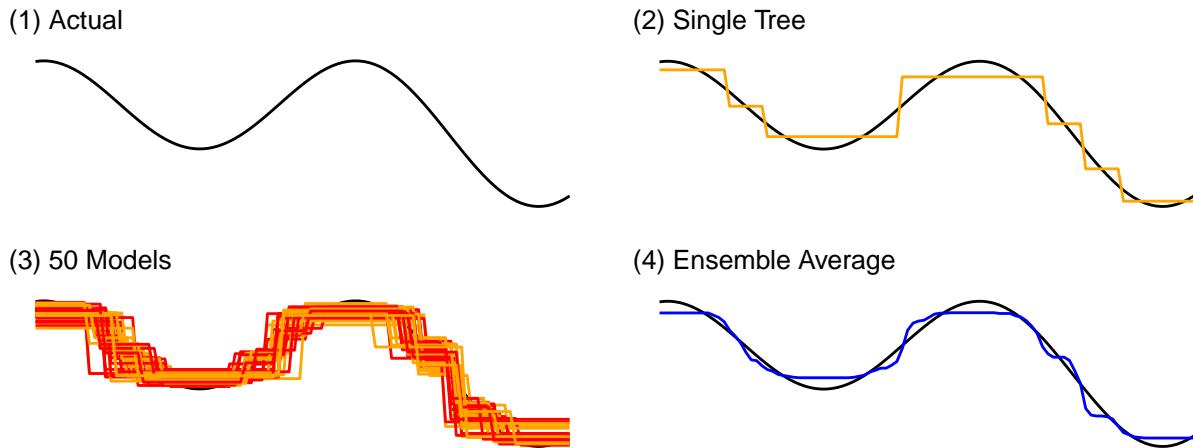


Figure 15: Comparison of results of applying a single model to fit a curve versus an ensemble of models.

1.6.4 Practice Exercise

Tune the decision tree to increase the mean F1-statistic by testing different windows for the average, maximum and standard deviation of `accel`.

1.7 Random Forests

How do we know anything for sure? Virtually every aspect of life has some uncertainty tied in. When a hurricane approaches the US Eastern Seaboard, forecasters often map the *cone of uncertainty* that provides the possible range of motion of a storm based on the results of many forecasted simulations. In presidential elections, often times the most polling results are ones that ensemble or average the results of many other similarly conducted polls. The reliance on predictions from a group of models with the same aims may well improve prediction accuracy. In statistical learning, average the results of multiple models is known as *ensemble learning* or *ensembling* for short.

Single models may impose biases on data and may be well-suited in specific situations. Ensemble methods combine the results of many models to obtain more stable results. For example, the curve in graph #1 can be approximated using a decision tree algorithm. The result of a single tree only loosely fits the curve in a jagged fashion (#2). That one tree may impose biases on the data, perhaps through how the tree is pruned or the assumption that the jagged approximation is appropriate, which may then translate into greater variance in predictions. One could imagine that the structure of that one tree may have happened by chance, and under different situations, the fit could be better.

Bootstrapping can help. Recall from elementary statistics that bootstrapping is defined as any statistical process that involves sampling records with replacement. By bootstrapping a sample, we treat a sample like a population, we can expose and characterize the qualities of an estimator under various scenarios already available in the data, which in turn produces an empirical probability distribution for predictions using the estimator. We can bootstrap the decision tree by (1) sampling the data with replacement up to the full size of the sample, then (2) run the decision tree. The result of repeating the process 50 times is (graph #3) produces a result that appears to be more organic and more accurate. This process of *bootstrapping* and *aggregating* the results is referred to as *bagging*.

Applying bagging to decision trees may not necessarily be enough to develop a well-balanced prediction. In the social sciences and public policy, it is generally assumed that a model's specification is a choice left to the analyst; However, it may also be a source of methodological bias.

Random forests can help. The technique is an extension of decision trees using a modified form of bootstrapping

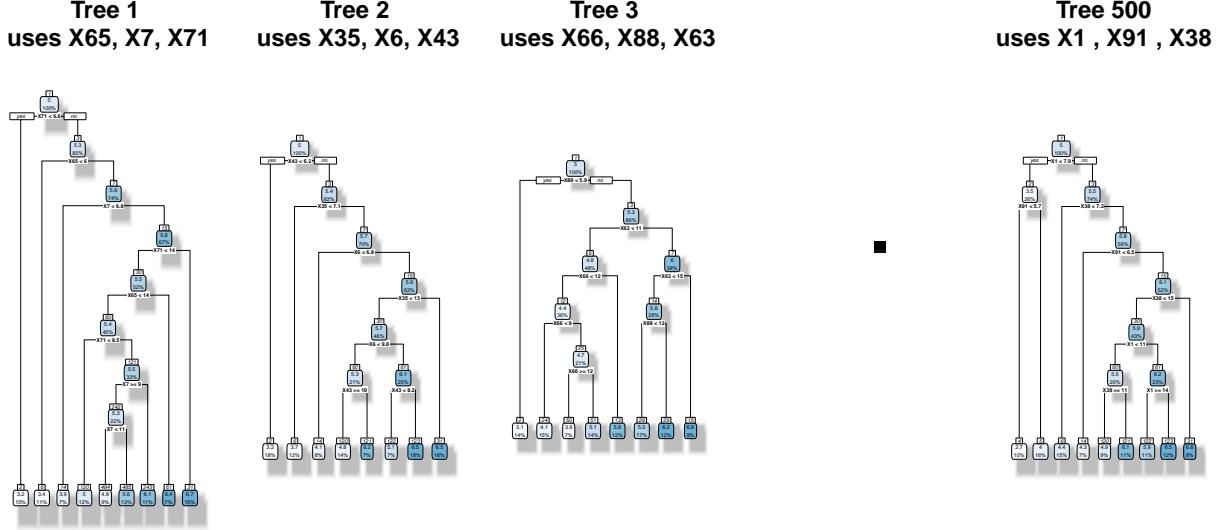


Figure 16: Random Forests construct hundreds of trees sampling from both observations and features, then combine the trees into one prediction through voting.

and ensemble methods to mitigate overfitting and bias issues.²⁸ Not only are individual records bootstrapped, but input features are bootstrapped such that if K variables are in the training set, then k variables are randomly selected to be considered in a model such that $k < K$. Each bootstrap sample is exhaustively grown using decision tree learning and is left as an unpruned tree. The resulting predictions of hundreds of trees are ensembled. The logic is described below.

Pseudo-code

Let S = training sample, K = number of input features

1. Randomly sample S cases with replacement from the original data.
2. Given K features, select k features at random where $k < K$.
3. With a sample of s and k features, grow the tree to its fullest complexity.
4. Predict the outcome for all records.
5. Out-Of-Bag (OOB). Set aside the predictions for records not in the s cases.

Repeat steps 1 through 5 for a large number of times saving the result after each tree.

Vote and average the results of the tree to obtain predictions.

Calculate OOB error using the stored OOB predictions.

The *Out-Of-Bag* (OOB) sample is a natural artifact of bootstrapping: approximately one-third of observations are naturally left un-selected, which can be used as the basis of calculating each tree's error and the overall model error. Think of it as a convenient built in test sample.

How about interpretation? Unlike decision trees, it is not a simple task to deduce rules or criteria that describe the target variable. Instead, random forests use *variable importance*, which, like for a decision tree, measures the contribution of a feature to the homogeneity of a classifier. Unlike decision trees, variable importance for a Random Forest is calculated as the mean decrease in the Gini coefficient of a split relative to the Gini coefficient of the root node. Gini coefficients measures homogeneity on a scale of 0 to 1, where 0 is perfect homogeneity and 1 is perfect heterogeneity. The Gini changes are summed for each variable and normalized.

²⁸Breiman (2001)

1.7.1 Tuning

Whereas methods like regression have a closed form solution, Random Forest require tuning as optimal models need to be searched for under different conditions. The principal tuning parameters include: Number of features and number of trees.

- *Number of input features.* As k number of parameters need to be selected in each sampling round, the value of k needs to minimize the error on the OOB predictions.
- *Number of trees* influences the stability the Variable Importance metric that is commonly used to infer variable influence in decision tree learning. More trees help to stabilize the Variable Importance estimate. To determine the number of trees, keep adding trees to a sample until the OOB error for a randomly select set of trees is approximately equal to that of the ensemble.

1.7.2 DIY: Revisiting activity classification

Like decision trees, much of Random Forests rely on easy to use methods made available through the `randomForest` library. There are a couple of ways to run the algorithm, including:

```
randomForest(formula, data, method, mtry, ntree)
```

where:
- `formula` is an object containing the specification to be estimated. Note that - `data` is a data frame.
- `mtry` is the number of variables to be randomly sampled per iteration. Default is \sqrt{k} for classification trees.
- `ntree` is the number of trees. Default is 500.

Using the same formula as the `rpart()` function, we can train a naive Random Forest and check the OOB error. Approximately 75.6% of observations in the OOB sample were correctly classified using 2 randomly selected variables in each of the 500 trees.

```
#Load randomForest library
library(randomForest)

#Run Random Forest
spec <- as.formula("factor(activity) ~ accel + avg50 + max50 + min50 + sd50")
fit.rf <- randomForest(spec, data = train, mtry = 2, ntree = 500)
```

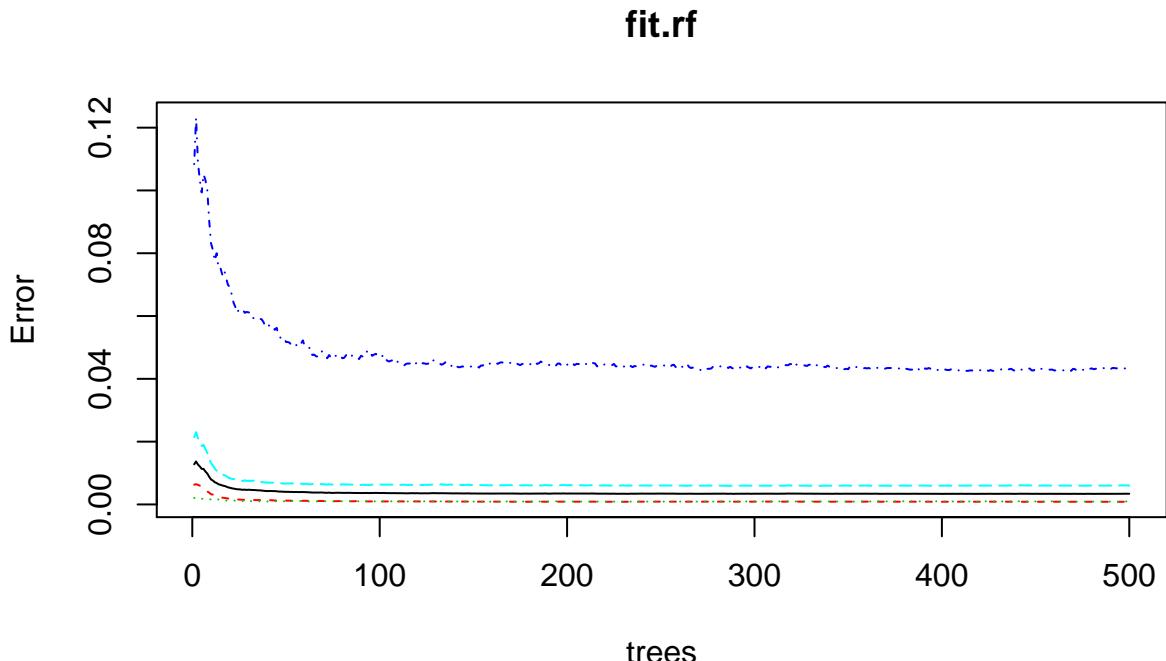
Using the `importance()` function, we can see the Mean Decrease Gini, which calculates the mean of Gini coefficients. However, the values themselves do not have any meaning outside of a comparison with other Gini measures.

```
importance(fit.rf)
```

```
##          MeanDecreaseGini
## accel      3902.912
## avg50     58664.594
## max50     40431.909
## min50     21605.409
## sd50      39314.035
```

By default, the `randomForests` library sets the number of trees to equal 500. By plotting the fit object, we can see how OOB error and the confidence interval converges asymptotically as more trees are added to the ensemble. Otherwise stated, more trees will help up to a certain point and the default is likely more than enough.

```
plot(fit.rf)
```



As we know that $n = 500$ trees is more than enough, we will now need to tune the tree for the number of variables. To tune the algorithm, we will use the `tuneRF()` method. The method searches for the optimal number of variables per split by incrementally adding variables. While it's a useful function, it is relatively verbose. In addition to the target and input features, a number of other parameters need to be specified:

```
tuneRF(x, y, ntreeTry, mtryStart, stepFactor, improve, trace, plot)
```

where: - `x` is a data frame or matrix of input features. - `ntreeTry` is the number of trees used in each iteration of tuning. - `mtryStart` is the number of variables to start. - `stepFactor` is the number of additional variables tested per iteration. - `improve` is the minimum relative improvement in OOB error for the search to go on. - `trace` is a boolean that indicates where to print the search progress. - `plot` is a boolean that indicates whether to plot the search results.

Below, we conduct a search from `mtryStart = 1` with a `stepFactor = 2`. The search result indicates that 2 variables per tree are optimal.

```
#Search for optimal number of input features
fit.tune <- tuneRF(x = train[,c("avg50", "max50", "min50", "sd50")],
                     y = factor(train$activity),
                     ntreeTry = 500, mtryStart = 1,
                     stepFactor = 2, improve = 0.001,
                     plot = FALSE, trace = FALSE)

## 0.3214629 0.001
## -0.1390071 0.001

#Extract best parameter
tune.param <- fit.tune[fit.tune[, 2] == min(fit.tune[, 2]), 1]
```

Normally, we can plug the tuned parameter back into the `randomForest()` method and re-train the algorithm, but it unnecessary in this case as the default model already uses the same parameters. When applied to the test set, we see that the mean F1-statistic is much improved – or a whole 10-percentage point increase.

```
#Predict classes in test
yhat <- predict(fit.rf, test)
```

```
#Calculate mean F1  
meanF1(test$activity, yhat)  
  
## [1] 0.8516852
```

This result does not mean that Random Forests will always turn better results, but rather multiple techniques should be tested when tackling prediction problems. Also, remember the policy goal: is the objective to predict or to explain? If a little of both, then it is worth understanding the value of increased accuracy at the cost of interpretability.