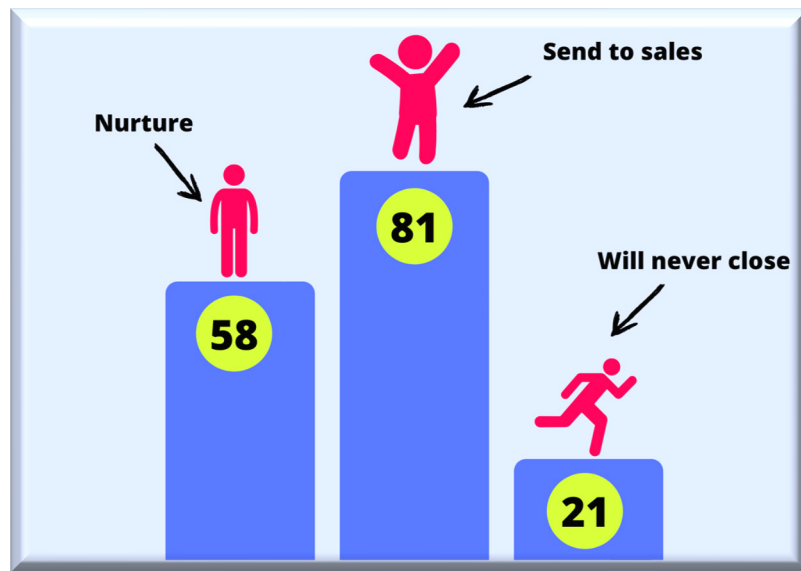


$$\frac{\sqrt{2.8}}{3+2^+}$$

# Lead Scoring with Cliently



$$\frac{-\sqrt{2}}{(\frac{1}{2})^2}$$



# Introduction

## LET'S DEFINE WHAT LEAD SCORING REALLY IS.

- It is the process of “scoring” company’s potential hot leads based on data collected from their website and marketing sources.
- This will prioritize your sales efforts starting with your most promising opportunities. In other words, it will show you which leads are worth following up on so that they turn into paying customers.
- Lead scoring is a quantitative approach that helps you show your sales reps which leads they should focus their efforts on.



What is lead scoring and why does it matter?

**In this guide, we'll discuss how you can make lead scoring work for your business and improve the effectiveness of your marketing campaigns.**

## THIS GUIDE WILL BE HELPFUL IF YOU NEED TO:

- Identify ways to prioritize prospects based on data collected from your website, email list, forms and other sources
- Establish a formula that works best for your company
- Calculate real-time lead scores with automated tools
- Identify high-quality leads and weed out unqualified ones
- Get answer to the most frequently asked questions on lead scoring

# Table of Contents

1.0

Why lead scoring matters

2.0

First things first – Am I collecting the proper data?

3.0

Point-based scoring – traditional methodology of lead scoring

3.1 Pardot scores

3.2 Why point-based scoring is primitive

4.0

AI-based lead scoring

4.1 Standard practice to calculate predictive scoring

4.1.1 Building various machine learning models

4.1.1.1 XGBoost

4.1.1.2 LightGBM

4.1.1.3 CatBoost

4.1.1.4 Random Forest

4.1.1.5 Neural Networks

4.1.1.6 Ensembles

4.1.2 Hyperparameter tuning

4.1.2.1 Types of hyperparameter tuning

4.1.3 Comparing machine learning models

4.1.4 Operationalize or MLOps

5.0

How Cliently can help you

5.1 How our algorithms are unique

5.1.1 A preview into Cliently's data cleaning

5.1.2 A preview into Cliently's feature engineering

5.1.2.1 Feature construction

5.2 Comparing our accuracies to other autoML solutions

5.3 Causal Modelling – Why? What to do?

6.0

Limitations of predictive scoring

7.0

Return on investment

○

8.0

Contact us



$$\frac{\sqrt{2.8}}{3+2^+}$$



1.0

Why lead  
scoring  
matters



# 1.1 Why lead scoring matters

## Part 1

Consider this - your outbound, inbound channels are bringing in 1000's of leads. How do you send all those leads to sales?

Would you just ask your Sales Reps to email or call all of them? There are so many steps that sometimes it takes months to close a deal. On the journey, sales reps find out that most of the leads are garbage but only after spending a considerable amount of time on each of them. When they are not able to meet their quotas, the finger-pointing begins.

If you ask any executive today in how they qualify leads, 80% will say BANT. To be considered for the next round of sales and marketing opportunity, it is important that a lead has a budget, authority, a need, and a timeline. This is usually produced after multiple conversations with the lead. A survey, a decade ago, concluded it takes 65 minutes on average per lead per sales rep to know BANT manually. Do you have 65 minutes to spend on each lead to know BANT?



# 1.2 Why lead scoring matters

## Part 2

The buying process has changed drastically as well. Information gathering for buyers is no longer a one-time event. It's an ongoing process that starts long before they've finalized their budget and timeline. They rely heavily on the internet to gather information too. For instance, they are watching your webinars, downloading your white papers, interacting with your social media pages, and browsing various pages on your website. Even if you want to connect with them over a sales call, it is unlikely you will get more information than their digital behaviors.

All these data points, when collected and used in lead scoring models, can give a true picture of the lead and determine if they are sales-ready or not. Just relying on BANT alone is a bad practice. Lead scoring steers companies to prioritize their time and efforts when reaching out to all their contacts -- using less of their valuable resources for those who are unlikely to convert and dedicating more of your resources to those who are more likely to convert



$$\frac{\sqrt{2.8}}{3+2^+}$$



“B2B marketers who emphasize lead volume over lead quality reduces sales efficiencies, increases campaign costs, and fuels the gap between sales and marketing. To generate a qualified demand, marketers need technology and processes that capture lead quality information; validate, score, and classify leads; develop programs to nurture leads that don’t yet warrant sales attention; and define metrics that directly identify marketing’s contribution to the sales pipeline and closed deals.”

— Laura Ramos

Forrester Research  
Improving B2B Lead Management

$$\frac{4+6+(2\sqrt{3})}{\sqrt{276}}$$

$$\frac{\sqrt{2.8}}{3+2^+}$$

# 2.0

## First Things First

Am I collecting the proper  
data?





## 2.0 First Things First

### Am I collecting the proper data?

We need lots of data points to differentiate serious leads from non-serious leads. The more accurate data we have for learning, the better the predictions. Overall, this is a two-way learning process. The machine will learn from the salesperson (from data that the salesperson enters in the system) and the salesperson will learn from the machine (on predictions it makes).

Companies need to consider two different kinds of information in their lead scoring: explicit and implicit.

Implicit data is based on information that you observe or infer about the prospect such as their online behaviors. It consists of tracking your lead's online behaviors. Examples include the number of clicks, time spent on the website, how many pages visited, forms filled, number of blogs visited, white papers downloaded, the location of their IP, etc.

Explicit data is based on information the prospect tells you or otherwise directly identifiable information. This is often collected through an online form or registration process. You can identify the prospect's demographical and firmographic variables like job title, company size, and annual revenue to see how well they compare against your ideal buyer profile. Many use data appending services to add more information. Third party companies sell data of individuals such as their past purchases, what applications they are using, etc. By enriching data, lead scoring becomes more accurate.



$$C = \frac{B^3 + C^2 + A}{3BA}$$



01 Watched a demo? For how many minutes?

02 Total time spent on the website?

03 How many pages visited?

04 If visited the Pricing page?

05 Downloaded white papers? How many?

06 Registered and attended a webinar?

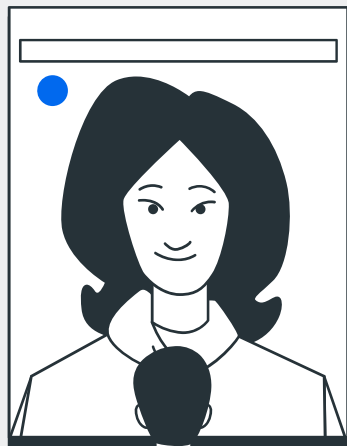
07 Engagement on social media

08 Number of clicks on website

A table to indicate some of the fields you should be collecting in your CRM -



$$\frac{\sqrt{2.8}}{3+2^+}$$



3.0

# Point-Based Scoring

Traditional Methodology of  
Lead Scoring



# 3.0 Point-Based Scoring

## Traditional Methodology of Lead Scoring

Remember that we said 80% of the executives use BANT to figure out qualified leads. That means that the other 19% use point-based scoring.

Some companies create internal ratings (based on experience and intuition) on behaviors (variables). For example, a person who watches a webinar is scored 10 points; a person who calls the company directly is scored 20 points, and so on.

Many companies use tools like Salesforce Pardot to automatically score their prospects.

Marketing Channel	Behavior	Score
Website	Requested a Demo	Most Important Factor <b>(+17)</b>
	Visited Careers Page	Negative Factor <b>(-10)</b>
	Visited Multiple Pages	Influencing Factor <b>(+8)</b>
	Visited Pricing Page	Important Factor <b>(+10)</b>
Event	Attended Event	Influencing Factor <b>(+7)</b>
	Had an Amazing Conversation	Important Factor <b>(+10)</b>
	Had a Good Conversation	Influencing Factor <b>(+30)</b>
Email	Opened Email	Influencing Factor <b>(+7)</b>
	Click within Email	Most Important Factor <b>(+5)</b>
	Forwarded Email	Influencing Factor <b>(+8)</b>
	Unsubscribed	Negative Factor <b>(-14)</b>
Content	Downloaded White Paper	Influencing Factor <b>(+7)</b>
	Downloaded a Specific White Paper	Important Factor <b>(+5)</b>
	Completed a Piece of Interactive Content	Influencing Factor <b>(+10)</b>
Webinar	Attended Webinar	Important Factor <b>(+30)</b>
	Registered for Webinar	Influencing Factor <b>(+20)</b>

*An example of point-based scoring*

## 3.1 Pardot Scores

Many companies who use point-based scoring use Salesforce Pardot. Pardot scores are Pardot's point system based on how many interactions a prospect has had with your Pardot elements. The higher the Pardot score, the more interactions they have had with your site or emails and thus become more aware of your product or service.

The default prospect scoring system comes with Pardot's pre-set point allocation for specific actions such as clicking on an email, page views, filling out a landing page form, etc. If a prospect performs one of these actions, the default number associated with that action is added to their Pardot score.

These scores can be customized as well (for every action, you can decide the number of points to be added or subtracted).



**If Ravi fills out a form on the website**



**His Pardot score increases 50 points**



**If Mary downloads a white paper**



**Her Pardot score increases 10 points**



**If Anna fills out a form on your website and later downloads the white paper from your email**



**Her Pardot score would increase to 60 points.  
(50+10=60)**

$$\frac{\sqrt{2.8}}{3+2^+}$$



4.0

# AI Based Lead Scoring

# 4.0 AI-Based Lead Scoring Saves the Day



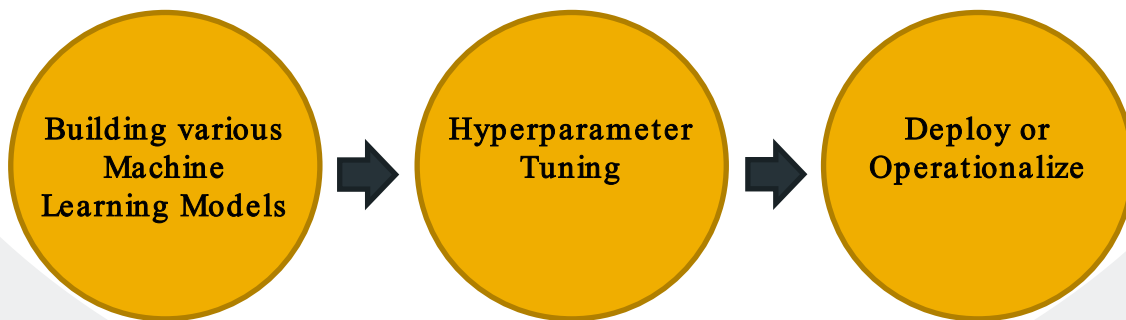
Instead of defining attributes and their corresponding weight factors, predictive lead scoring uses existing sales data, along with data mining and analytics techniques to build the 'ideal' lead. No faulty intuition or guesses needed! Subsequent leads are compared with present leads and correspondingly labelled as qualified or not.

In predictive lead scoring, the algorithms create a formula for automatically ranking leads. The model is continuously fed with data on leads that successfully converted to customers or those that failed thus eliminating the need for 'run and check' processes.



## 4.1 Standard Practices to Calculate Lead Scoring

For AI-based scoring, we have to combine explicit and implicit together. Data is often presented in Salesforce, Hubspot, Intercom, Outreach, Salesloft, Zoominfo, etc. And we must bring it all together. Once this is done, there are three steps:

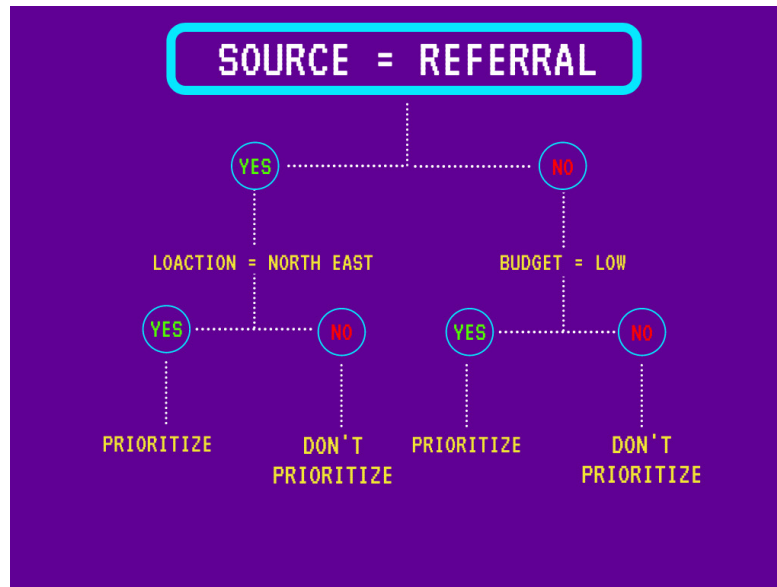




## 4.1.1 Building Various Machine Learning Models

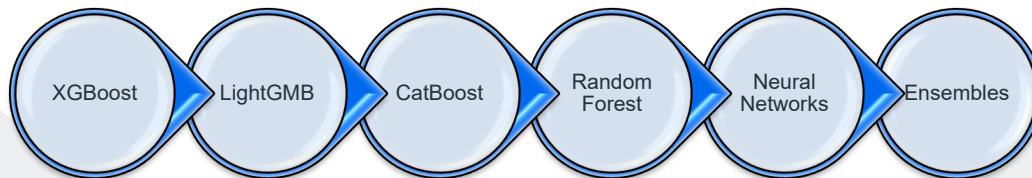
- Most of the time tree-based models, like LightGBM and XGBoost are winners as stand-alone models.

Neural Networks are popular too, but they don't do as well. When various models are combined, we get Ensemble Models, and they usually outperform single models.



*An example of a simple Decision Tree (tree-based model).*

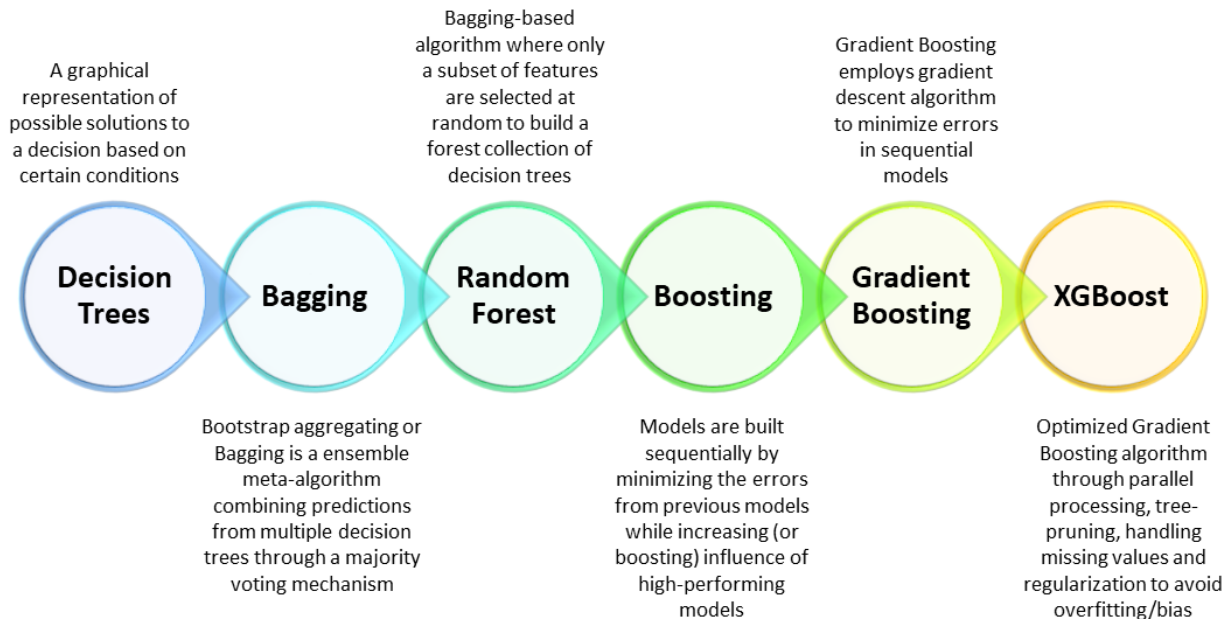
### Some of the Machine Learning Models



# XGBoost

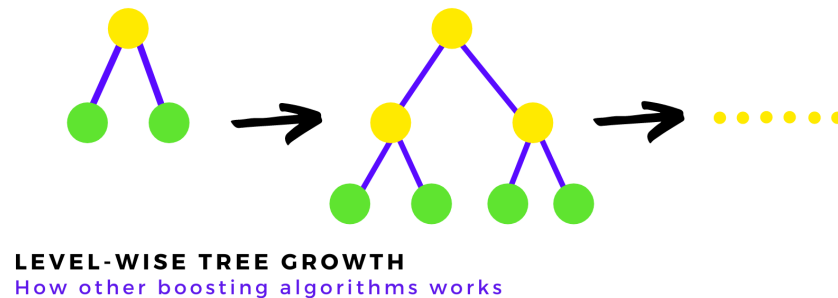
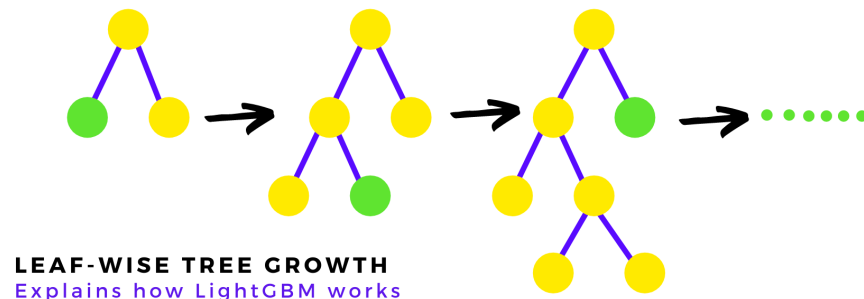
XGBoost stands for “Extreme Gradient Boosting”. XGBoost is an optimized distributed gradient boosting library designed to be highly efficient, flexible, and portable. It implements Machine Learning algorithms under the Gradient Boosting framework. It provides parallel tree boosting to solve many data science problems in a fast and accurate way.

XGBoost algorithm was developed as a research project at the University of Washington. Tianqi Chen and Carlos Guestrin presented their paper at the SIGKDD Conference in 2016 and caught the Machine Learning world on fire. Since its introduction, this algorithm has not only been credited with winning numerous Kaggle competitions, but it also has been the driving force under the hood for several cutting-edge industry applications.



# LightGBM

LightGBM trees grow vertically while other algorithms trees grow horizontally meaning that LightGBM grows trees leaf-wise while other algorithms grow level-wise. It will choose the leaf with max delta loss to grow. When growing the same leaf, a Leaf-wise algorithm can reduce more loss than a level-wise algorithm. LightGBM's are highly efficient on large datasets.



CatBoost is an algorithm for gradient boosting on decision trees. Developed by Yandex researchers and engineers, it is the successor of the MatrixNet algorithm.

## CatBoost & Random Forest

- ◇ CatBoost has gained much popularity compared to other gradient boosting algorithms primarily due to the following features:

- Ordered Boosting to overcome overfitting.
- Native handling for categorical features.
- Using Oblivious Trees or Symmetric Trees for faster execution.

Random forest is a supervised learning algorithm. The "forest" it builds is an ensemble of decision trees, usually trained with the "bagging" method. The general idea of the bagging method is that a combination of learning models increases the overall result. Put simply, random forests build multiple decision trees and merge them together to get a more accurate and stable prediction.



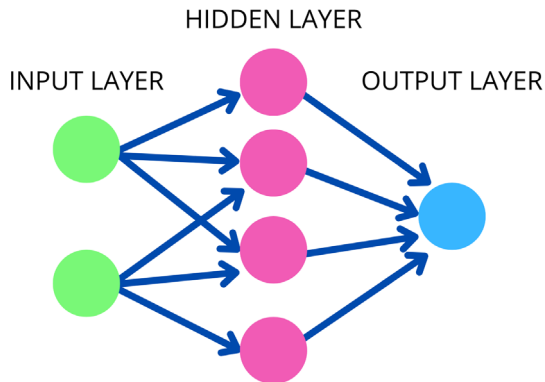
# Neural Networks

Their name and structure are inspired by the human brain, mimicking the way that biological neurons signal to one another. Artificial neural networks (ANN's) are comprised of node layers containing an input layer, one or more hidden layers, and an output layer. Each node, or artificial neuron, connects to another and has an associated weights and thresholds. If the output of any individual node is above the specified threshold value, that node is activated, sending data to the next layer of the network. Otherwise, no data is passed along to the next layer.

Neural networks reflect the behavior of the human brain allowing computer programs to recognize patterns and solve common problems in the fields of AI, machine learning, and deep learning.

Neural networks rely on training data to learn and improve their accuracy over time. However, once these learning algorithms are fine-tuned for accuracy, they are powerful tools, allowing us to classify and cluster data at a high velocity.

## A SIMPLE NEURAL NETWORK

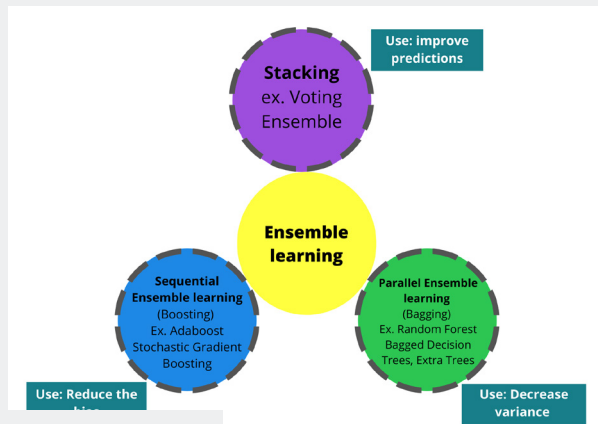


# Ensembles

○

Ensemble methods consists of a machine learning technique that combines several base models in order to produce one optimal predictive model.

The most popular ensemble methods are boosting, bagging, and stacking.



## 4.1.2 Hyperparameter Tuning

The same kind of machine learning model can require different constraints, weights or learning rates to generalize different data patterns. These measures are called hyperparameters and must be tuned so that the model can optimally solve the machine learning problem. We need to see, at which given hyperparameters, a given model gives us the highest accuracy.

The best way to think about hyperparameters is like the settings of an algorithm that can be adjusted to optimize performance, just as you might turn the knobs of an AM radio to get a clear signal.

Hyperparameter tuning relies more on experimental results than theory. Thus, the best method to determine the optimal settings are to try many different combinations to evaluate the performance of each model.

With so much uncertainty it is understandable that most Data Scientists forego hyperparameter tuning and stick with the default values provided by the model. This works in most cases, but when it comes to sales data with low conversion, even a 1-3% accuracy difference can make or break sales quotas.

What should the value be for the maximum depth of the Decision Tree?

How many trees should I select in a Random Forest model?

Should use a single layer or multiple layer Neural Network?

If multiple layers, then how many layers should there be?

How many neurons should I include in the Neural Network?

What should be the minimum sample split value be for the Decision Tree?

What value should I select for the minimum sample leaf for my Decision Tree?

How many iterations should I select for the Neural Network?

What should the value of the learning rate for gradient descent be?

Which solver method is best suited for my Neural Network?

What should the value be for C and sigma in the Support Vector Machine?

**Hyperparameter tuning helps us answer questions like these.**



## 4.1.2.1 Types of Hyperparameter Tuning

### Grid Search

Grid search is arguably the most basic hyperparameter tuning method. With this technique, we simply build a model for each possible combination for all of the hyperparameter values, evaluate each model, and select the architecture which produces the best results.

### Random Search

Random search differs from grid search in that we no longer provide a discrete set of values to explore for each hyperparameter. Instead, we provide a statistical distribution for each hyperparameter from which values may be randomly sampled.

### Bayesian Optimization

The previous two methods performed individual experiments, building models with various hyperparameter values, and recording the model performance for each. Because each experiment was performed in isolation, it's very easy to parallelize this process. However, because each experiment was performed in isolation, we're not able to use the information from one experiment to improve the next experiment. Bayesian optimization belongs to a class of sequential model-based optimization (SMBO) algorithms that allow for one to use the results of a previous iteration to improve sampling methods of the next experiment. Simply put, Bayesian Optimization addresses the pitfalls of grid and random search by incorporating a “belief” of what the solution space looks like and by “learning” from the configurations it evaluates.





$$\frac{A}{3B}$$



### 4.1.3 Comparing Machine Learning Models

Remember, since data is often imbalanced, we need to consider metrics like Weighted F1, MCC, KS Statistic, and Cohen's Kappa to measure model performance.

These metrics put almost equal importance in classes into target (what you are predicting, lead scoring in this case).

The higher the value of these metrics, the better the performance of the models.

We also check the accuracy on test datasets. We split the data into training (70%-80%) and testing datasets (20%-30%). The training dataset is used to train, and the testing dataset is a holdout dataset, where we apply models to see performance.

$$\frac{5 \pm \sqrt{3-4}}{2}$$

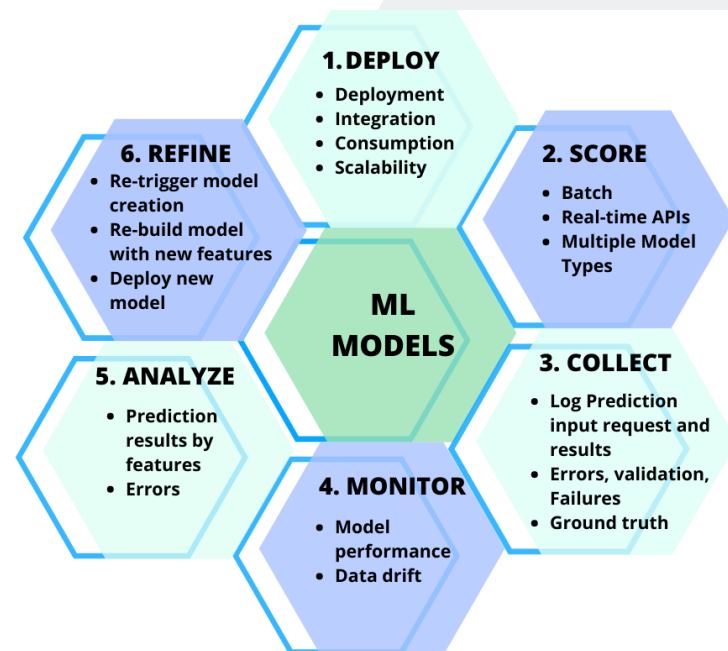


$$\frac{10+17}{3.45}$$

## 4.1.4 Operationalize or MLOps

One of the most critical, but often overlooked aspects of the Machine Learning process is the operationalization of Machine Learning models. Operationalization refers to the deployment of models to be consumed by business applications in order to predict the target class or target value of a classification/regression problem. Machine Learning models have no tangible business benefits until they're operationalized.

In other words, you need consumable, actionable insights and predictions for your Sales Reps in your desired CRM. And these predictions must be in real-time. A new lead at the moment it enters CRM system should have a conversion score. If that lead watches your webinar tomorrow, that should be reflected in lead scoring in real-time (the score will increase).



$$\frac{\sqrt{2.8}}{3+2^+}$$



5.0

# How Cliently Can Help You

## 5.0 How Cliently Can Help You

All the work we have discussed requires a team of data scientists and engineers to make it possible.

We have a suitable platform to give lead scoring predictions in real-time, no matter the number of variables you have or the amount of data. In fact, the more data the more accurate the output.

The best part is that we give you real-time predictions.

As your data updates, so does the lead scoring.

We will also let you know what can increase your scores and where you should focus to get more and bigger opportunities using our causal modelling reports.

If the junk leads take actions like downloading a paper, at very later stages (after many months of cadence), the scores are then adjusted and are displayed as hot leads.

On top of all of this, we not only tell you lead scores, but also give you the best possible action to engage with them to have the highest chances of conversion.



Real-time View > **Contacts**

**Contact Recipes**

My Contacts	1192
All Contacts	1192
Active Contacts	544
Hot Leads	19
<b>Most Engaged</b>	<b>698</b>
Glengarry	365
Back from the Dead	501
It's Been a While	198
Might Churn	67

**Most Engaged**

contact owner is Spencer Farber and

and phone number is not empty and

**698 contacts**  
@ 575 companies Engage More

☐ CONTACT INFO

- Dontae Little**  
Technical Writer @ Slack, Inc.
- Sammy Lawson**  
Financial Analyst @ Eare
- Nuria Pelayo**  
Help Desk Technician @ Midel
- Saga Lindén**  
Computer Systems ... @ Jaxworks
- Maria Paula Morterero**  
Budget/Accounting ... @ Linkline

**Contact GPS Roadmap**

**Dontae Little**  
Technical Writer  
Slack, Inc. **88%**

**Send Email today @ 2pm**  
Zengo Team Demo Follow-up

**Create** **Ignore**

**July 12 @ 2pm**  
New Articles on the Blog

**July 19 @ 2pm**  
Join Bright Side Now!

**July 22 @ 2pm**  
Call Script 7

**Global Lead Scoring GPS**

Lead Score > 75%  
Company Revenue > \$5M

**Dontae Little**  
{Company} Team  
Demo Follow-up **88%**

**Send this email today @ 2pm**  
88% likelihood of engaging

**Miranda Kalkofen**  
Need some  
feedback on this **75%**

**Send this email on Jul 12 @ 2pm**  
46% likelihood of engaging

**Sim Machen**  
Reminder about  
your Appointment **79%**

**Send this video on Jul 12 @ 2pm**  
65% likelihood of engaging

$$C = \frac{B^3 + C^2 + A}{3BA}$$

A screenshot of what the view looks like inside our app -

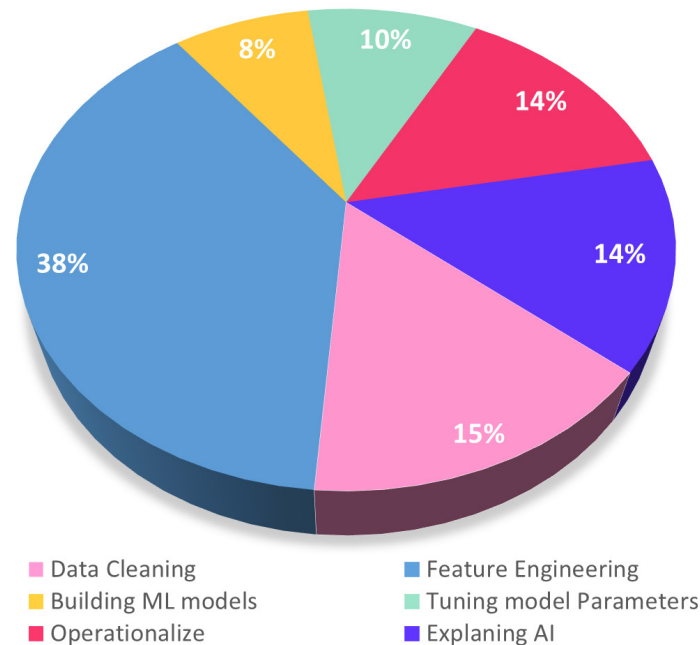


## 5.1 How Our Algorithms are Unique

Our accuracies are very high because of our algorithms.

From day 1, our focus has been on data cleaning and feature engineering. The kind of information we create from current information is unparalleled. We have 50,000 lines of code to clean all sorts of data, take various permutations, combinations into account, and create more variables from current ones.

We have learned best practices from Data Science competitions in Kaggle and implemented them into ours. Not to mention, our auttml product performs in top 1% of rankings in every competition.



## 5.1.1 A Preview Into Cliently's Data Cleaning

$$C = \frac{B^3 + C^2 + A}{3BA}$$

Modelling

Preparing the Data

“Quality Data Beats Fancy Algorithms”

Not everyone has the same goals and data. For some, the most important thing is booking demos that can be converted into opportunities. For others it can be sales. Our team of Data Scientists understands your requirements and will set up an automated infrastructure specific to your data and needs.

Let's take a peek into some of our data cleaning practices....



# Data Cleaning Practices

## Step 1

Remove duplicate or irrelevant observations - Irrelevant observations are when you notice observations that do not fit into the specific problem you are trying to analyze. For example, you should be looking at leads you acquired within the last few years and not include leads acquired few decades ago because their buying behaviors have changed so much in recent times. In the future you would be acquiring leads similarly to what have done recently. This can make analysis more efficient and minimize distraction from your primary target—as well as creating a more manageable and more performant dataset. Similarly, when we join datasets from all the places, there can be a lot of duplicates, we use fuzzy match to find duplicates and remove them.

## Step 2

Fix structural errors - Structural errors occur when is measured or transferred. It appears as strange naming conventions, typos, or incorrect capitalizations. These inconsistencies can cause mislabeled categories or classes. For example, we may find “N/A” and “Not Applicable” both appear, but we analyze them as the same category. The misspellings are corrected here using our automated data cleaning.

## Step 3

Filter unwanted outliers - Often, there will be one-off observations where, at a glance, do not appear to fit within the data we are analyzing. If we have a legitimate reason to remove an outlier, like with data-entry mistakes, doing so will help the performance of the data you are working with. However, sometimes it is the appearance of an outlier that will prove a theory you are working on. Remember, just because an outlier exists, doesn't mean it is incorrect. This step is needed to determine the validity of that number. If an outlier proves to be irrelevant for analysis or is a mistake, we consider removing it.





## Step 4

Handle missing data - You can't ignore missing data because most machine learning algorithms will not accept missing values. We have a couple of ways to deal with missing data. Neither is optimal, but both can be considered depending on the situation. As a first option, we can drop observations that have missing values, but doing this will drop or lose information, so we are very mindful of this before we remove it. As a second option, we can impute missing values based on other observations. We use mean imputation methods and tree-based imputations. A third option can be to give some a number or value like "NA" or "9999" and treat them as individual values.

## Step 5

Standardize + Normalize - Standardization and normalization are crucial to the effectiveness of the data cleaning process as they make data ripe for statistical analysis and easy to compare and analyze. Standardization is a process during which you're ensuring that all of your values adhere to a specific standard, such as deciding whether to go with kilos or grams, upper- or lower-case letters, etc. Normalization is the process of adjusting the values to a common scale. For example, you can rescale values into the 0-1 range. This action is necessary also because many models require normalized data.



## Step 6

Handle an Imbalanced Dataset - Modelling imbalanced data is a major challenge that we face when we train a model. When dealing with classification problems the class balance of the target class label plays an important role in modelling. For imbalanced class problems, such as with the presence of minority class in the dataset, the models try to learn only the majority class which results in a biased prediction.

In lead scoring, we see this a lot. Since opportunity is 2%-10%, we need to balance the data. At Cliently, we use undersampling, where we make the numbers of major classes (non-opportunities) equal to minor classes (opportunities). So, if there were 5% opportunities and 95% non at the start, we bring both classes to be 50:50.

At the time of scoring, we desire that estimated classes should be a true sampling of original training data, hence we calibrate the probabilities. When we are scoring or deploying the models, we will be predicting 50% of leads will convert into opportunities (due to oversampling). This would be wrong as only 5% were opportunities. Hence, we bring them to the original ratio of 5:95 using Platt Scaling. We desire that the estimated class probabilities are reflective of the true underlying probability of the sample. That is, the predicted class probability (or probability-like value) needs to be well-calibrated.

## Step 7

Validate and QA - At the end of the data cleaning process, we answer these questions as a part of basic validation:

Does the data make sense?

Does the data follow the appropriate rules for its field?

Does it prove or disprove our working theory or bring any insight to light?

Can we find trends in the data to help us form our next theory?

If not, is that because of a data quality issue?



○

## 5.1.2 A Preview Into Cliently's Feature Engineering ○

$$C = \frac{B^3 + C^2 + A}{3BA}$$

When our goal is to get the best possible results from a predictive model, we need to get the most from what we have. This includes getting the best results from the algorithms we are using. It also involves getting the most out of the data for our algorithms to work with.

Feature engineering is the process of transforming raw data into features that better represent the underlying problem to the predictive models, resulting in improved model accuracy on unseen data.

“You have to turn your inputs into things the algorithm can understand!”

◇

○



## 5.1.2.1 Feature Construction

Much of the raw data must be converted into numeric variables that Machine Learning models can understand. 30% of the code focuses on this part alone. There are just so many permutations and combinations possible. If a data science projects duration is 1 month, we are able to bring it down to 2 to 3 days only because of our automated feature engineering and data cleaning!

70% of the variables contain non-numeric information. Examples include date, time, address, text, alpha-numeric characters, phone numbers (ex- (972) 672-6434), emails, urls, speical characters etc.

**Let us show you some of the examples in our Feature Construction**



## Decomposing a Date-Time

A date-time contains a lot of information that can be difficult for a model to take advantage of in its native form, such as ISO 8601 (i.e., 2014-09-20T20:45:40Z). So we break it down into individual components.

Here is an example of 2 data columns (with different formats) and how we transform them. Our smart algorithms not just only convert date into same formats to do subtractions but also break down date into smaller numeric components. This data can't be used as it is in machine learning.

Date the person replied first: 10NOV2020:03:49:19	
Date when person booked a demo: 09/21/2021T08:03:00 AM	
↓	
Day of the month when person replied	10
Month when the person replied	11
Year when the person replied	2020
Hour when person replied	3
Minute of that hour when person replied	49
Difference in days between nearest holiday and date of reply	16*
Difference in days between from the day person replied to day person booked a demo	315
Day of the month when person booked a demo	21
Month when the person booked a demo	09
Year when the person booked a demo	2021
Hour when person booked a demo	08
Minute of that hour when person booked a demo	03
Difference in days between nearest holiday and date of demo booking	15**

\* Difference calculated from Thanksgiving day 2020

\*\* Difference calculated from Labour day 2021





## Decomposing Text Columns:

We break down text (emails, phone conversations, chat support etc.) into 3 components:

- Topic analysis (also called topic detection, topic modeling, or topic extraction) is a machine learning technique that organizes and understands large collections of text data by assigning “tags” or categories according to each individual text’s topic or theme. Topic analysis uses natural language processing (NLP) to break down human language so that we can find patterns and unlock semantic structures within texts to extract insights and help make data-driven decisions.
- We calculate sentiments of a text, whether positive, neutral, or negative and gauge by how much.
- Subjectivity analysis recognizes the contextual polarity of opinions, attitudes, emotions, feelings, etc. regarding products, services, topics, or issues. Subjectivity classification categorizes the given text as subjective or objective. While an objective text contains one or more facts about a product or an issue, a subjective text expresses a person's opinions.



## Decomposing an Address Column

📍 Address of the lead	
416, 1400 20 <sup>th</sup> St NW, Washington DC, 20036	
Apartment number	416
Building Number	1400
Street Number	20
Neighbourhood/district	NW
City	Washington
State	DC
Zip code	20036
Country	USA



## 5.2 Comparing Our Accuracies to Other autoML Solutions



A few automl solutions like google automl, Microsoft azure, aws sagemaker and datarobot could be good solutions if you have data scientists and engineers to set up infrastructure. They may be more cost-effective as well.

**We have compared 100+ datasets and 85% of the time, we performed better than these 4 automl platforms in accuracies (weighted f1, mcc, rmse) !!!**

How do we achieve higher accuracies? It stems from the intensive data cleaning and feature engineering we discussed in the above sections. If we were just to build models on raw data, our accuracies would be similar or even lower.

## 5.3 Causal Modeling

### Why? What to do?



At Cliently, we give you action plans and recommendations to improve lead conversion. We call this Explainable AI or Causal modelling. It focuses on determining which factors effectively influence outcomes. This is significant in that it helps determine which kinds of sales strategies are most effective.

We give you answers to questions like -

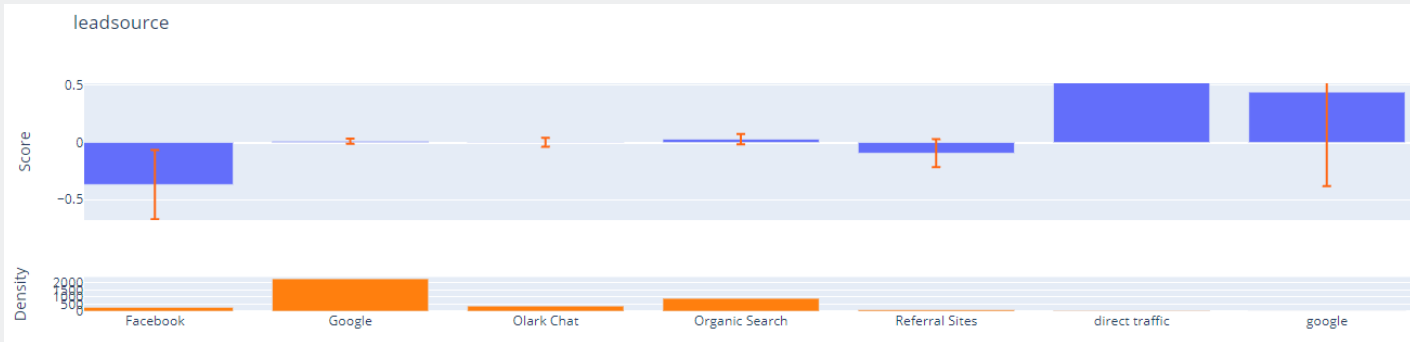
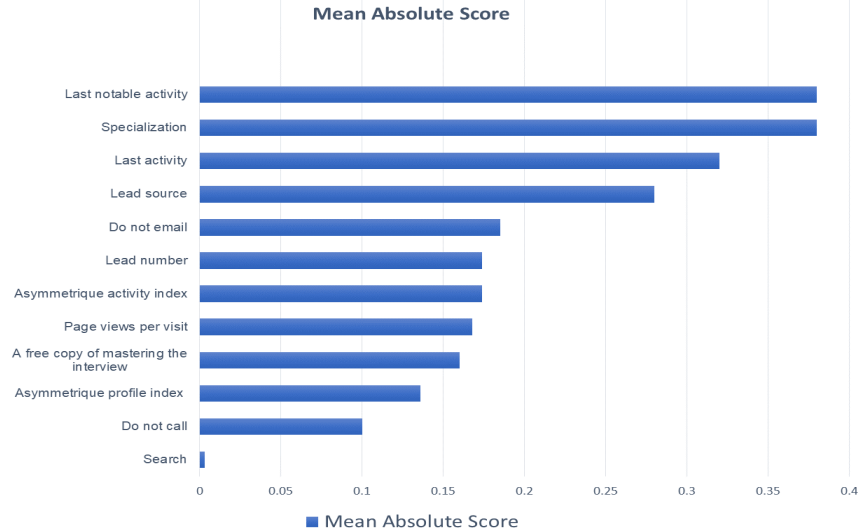
- ❖ What variables or fields improve the chances of conversion?
- ❖ By changing which variables by how much or by how, will the conversion be affected?
- ❖ What are the ideal prospects and how to acquire more of such prospects?
- ❖ ROI analysis
- ❖ What should be the cadence for outreach to maximize lead scoring?





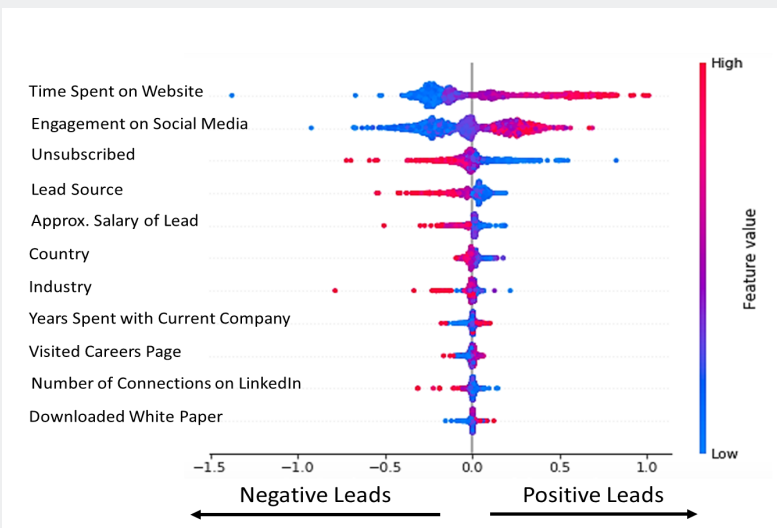
## Some snippets from our action plans -

A graph telling which variables are important in lead conversion.



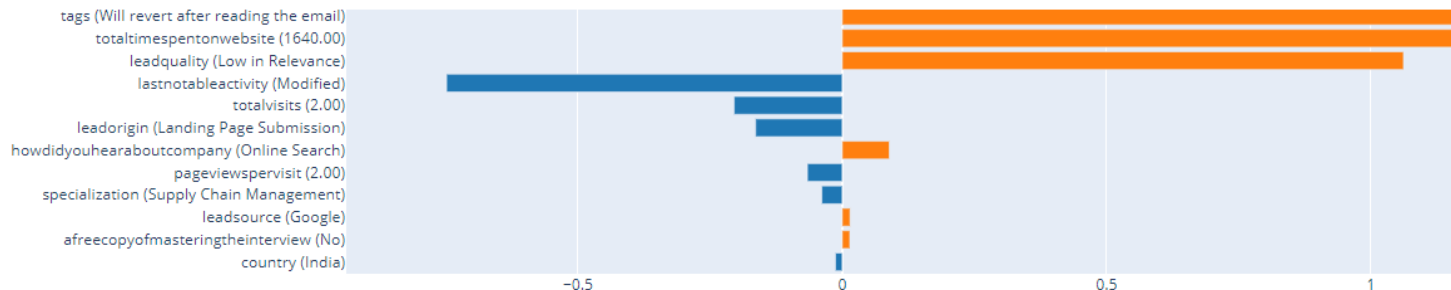
A breakdown of important individual variables can help us understand how changing values affects conversion.

For instance, looking at this, we can understand, which channels lead to higher conversion.

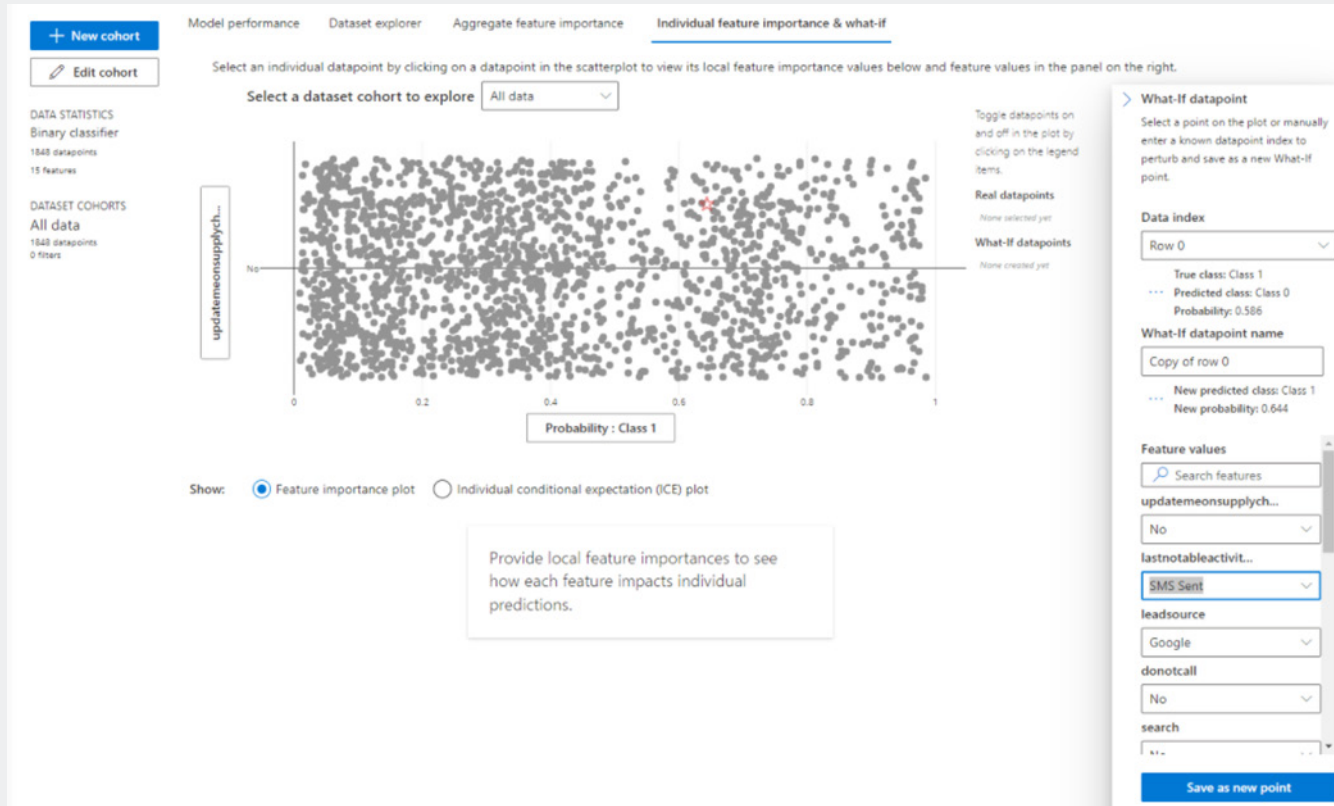


Shap plots can help us understand how increasing or decreasing variables affects conversion. For instance, in this example, as time spent on website increases, there is higher conversion. A counterintuitive point here - higher the salary of the lead, lesser the chances of conversion

Predicted (1.0): 0.972 | Actual (1.0): 0.972



This graph depicts a random lead which has 97% chance of converting. The variables with orange bars are helping this person to convert, the ones in blue are pushing him to non conversion. Since weightage of orange bars is more, this person has positive outcome.



## What if Scenarios

Here 1 means conversion. and 0 non-conversion

On the right-hand side, as you change the values of variables, you get new probabilities of conversion.

This example indicates, if we sent this person's sms, that his/her chances of conversion would have gone up by 6 points.

$$\frac{\sqrt{2.8}}{3+2^+}$$

6.0

# Limitations of Predictive Scoring



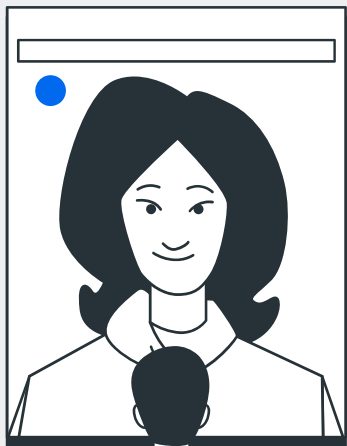
# 6.0 Limitations of Predictive Scoring

**Predictive scoring is not for all companies. Here are the shortcomings —**

- Historical data: Since models depend on historical data, you need to have at least a few months (if not longer) of records.
- The more data the better. Collect as much as possible. For instance, take web activity and social media activity and combine it with third party data. For example, credit scores, annual spending, etc.
- The explicit data we collect can, at times, be incorrect. For instance, when we ask BANT information over forms, calls, or emails, many people answer incorrectly or won't answer. Some reasons could be - They work with changing timeline and budget priorities; - Too early into the buying process where decision-makers have yet been decided on; Or they might prefer that Salespeople avoid contacting them!
- Most predictive scoring systems are black boxes: The simple models like logistic regression and decision trees can easily be explained, but more complex algorithms like ensemble models and lightGBM's are harder to explain. We use more complex AI models because they are more accurate. However, at Cliently, we have come up with innovative ways to explain black-box models, which we covered in causal modelling.
- Your target market and product should be stable. If your target market changes, your sales data loses its value for analytics. However, if you know the critical junction after which your company decided to focus on a certain market, you could take data from that point on for training data for the models.



$$\frac{\sqrt{2.8}}{3+2^+}$$



7.0

Return on  
Investment



## 3.2 Why Point-Based Scoring is Primitive

- ◇ In point-based lead scoring, you have predefined criteria that are multiplied by weighting factors to deduce the lead's total score. The problem is three-fold.

### Problem #1

You are using your past knowledge for weighting factors. For instance, you could assign more weight to the type of email address (company email or personal email) than it deserves.

### Problem #2

The criteria that you are using may not be predictive of lead conversions. You are using your intuition to help make predictions.

### Problem #3

Every time your intuition goes wrong, you must go back to the old rules and adjust. For example, earlier, you were rewarding leads from California. Now, the leads from CA are not so good so you go back, adjust, and reduce points for leads living in California.

# 7.0 Return on Investment

Engaging with leads is an intensive process. Focusing on the right accounts and contacts, engaging only with warm leads versus reaching out to all, can translate into a 40%+ increase in productivity.

Creating common “lead” definitions simplifies follow-up processes and drives alignment.

The most significant expense is people. Senior Reps can focus on higher-ranked leads. The average salary of a Senior Sales Rep is \$55,000. Every hour of their time is a real cost to the company.

Let's take an example of a team receiving 3000 leads a month comparing point based and predictive scored leads.

	No Lead Scoring	Point-based Lead Scoring	AI-based Lead Scoring
Total time spent on leads	2000 hours	1300 hours	900 hours
Number of SDRs needed *	13	8	5
Percentage conversion	8%	12%	15%
Opportunities converted	240	360	450

This is an example, based on industry averages, but can also be heavily dictated by variables around team size, target, data accuracy, and lead source.

In a study of 10 B2B organizations using lead scoring systems, [Eloqua](#) found that, on average, deal close rates increased by 30%, company revenue increased by 18% and the revenue per deal increased by 17%.

According to [Aberdeen Research](#), companies that focus on lead scoring the correct way have a 192% higher average lead qualification rate than those that do not.





$$\frac{\sqrt{2.8}}{3+2^+}$$



8.0

Contact  
Us



# 8.0 Contact Us

$$\frac{\sqrt{2.8}}{3+2^+}$$

Cliently is the first truly AI-based Revenue Intelligence Application.

Create custom Recipes (Views) to understand the entire Sales Journey from all of your sources in one place.

Automatically generate customized Real-Time AI predictions that tells reps which accounts and contacts to engage with and which action to take in order to maximize sales and save countless hours. Reps get insights, recommendations, and predictions from data in a very consumable way.

Create engaging automated outreach playbooks for your Sales team with an omnichannel approach using everything from email, to videos, to gifts. Reps can take action directly from Cliently's UI.

If you are interested in AI-based scoring automations, contact us.



[www.cliently.com](http://www.cliently.com)



[info@cliently.com](mailto:info@cliently.com)



$$\frac{\sqrt{2}}{(\frac{1}{2})^2}$$

