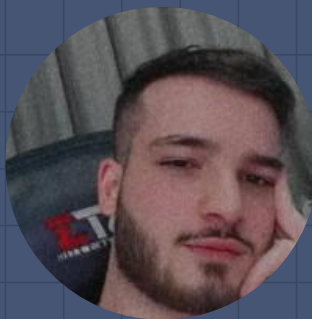


# Treinamento e validação de modelos preditivos com AM

# GRUPO



**Rafael Petry**  
00314280



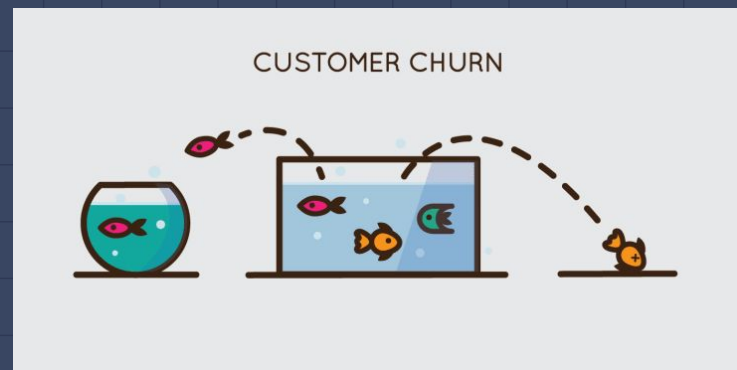
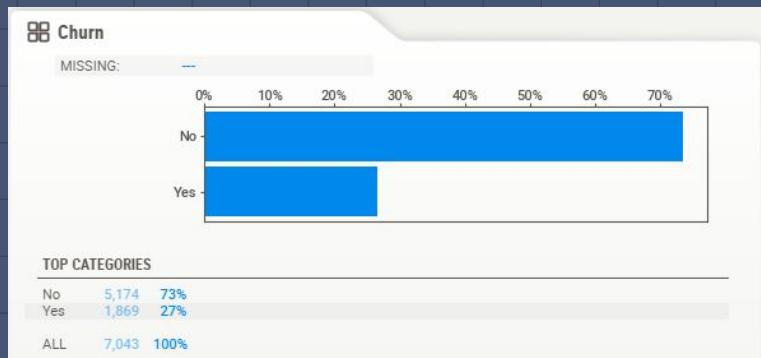
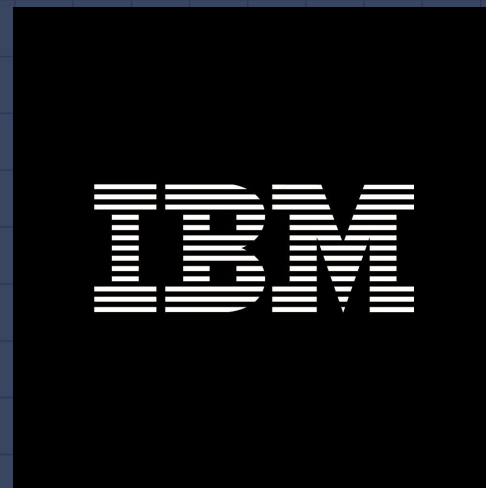
**Lucas Melo**  
00315747



**Jean Andrade**  
00252846

# DATASET

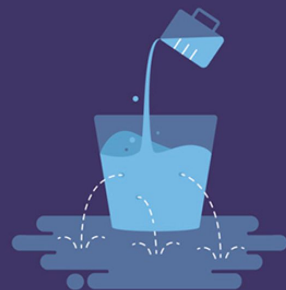
Optamos pelo dataset  
"Telco Customer Churn" da IBM  
que contém dados sobre os  
clientes e serviços contratados



Coluna	Intervalo
gender	[Male, Female]
SeniorCitizen	[0, 1]
Partner	[Yes, No]
Dependents	[Yes, No]
Tenure	[0, ..., 72]
PhotoService	[Yes, No]
MultipleLines	[Yes, No, No service]
InternetService	[DSL, Fiber, No service]
OnlineSecurity	[ Yes, No, No service]
OnlineBackup	[ Yes, No, No service]

Coluna	Intervalo
DeviceProtection	[Yes, No, No service]
TechSupport	[Yes, No, No service]
StreamingTV	[Yes, No, No service]
StreamingMovies	[Yes, No, No service]
Contract	[Monthly, One Year, Two Years]
PaperlessBilling	[Yes, No]
PaymentMethod	[E-check, Mail-Check, Transfer, Credit Card]
MonthlyCharges	[18.30, ..., 119.00]
*TotalCharges	[18.80, 8684.80]
Churn	[Yes, No]

**Qual a probabilidade de  
um cliente cancelar o  
serviço?**



# ATRIBUTOS

## Irrelevantes / Removidos

- customerID - Somente ID do cliente não é importante.
- StreamingTV-Baseado na análise de correlação.

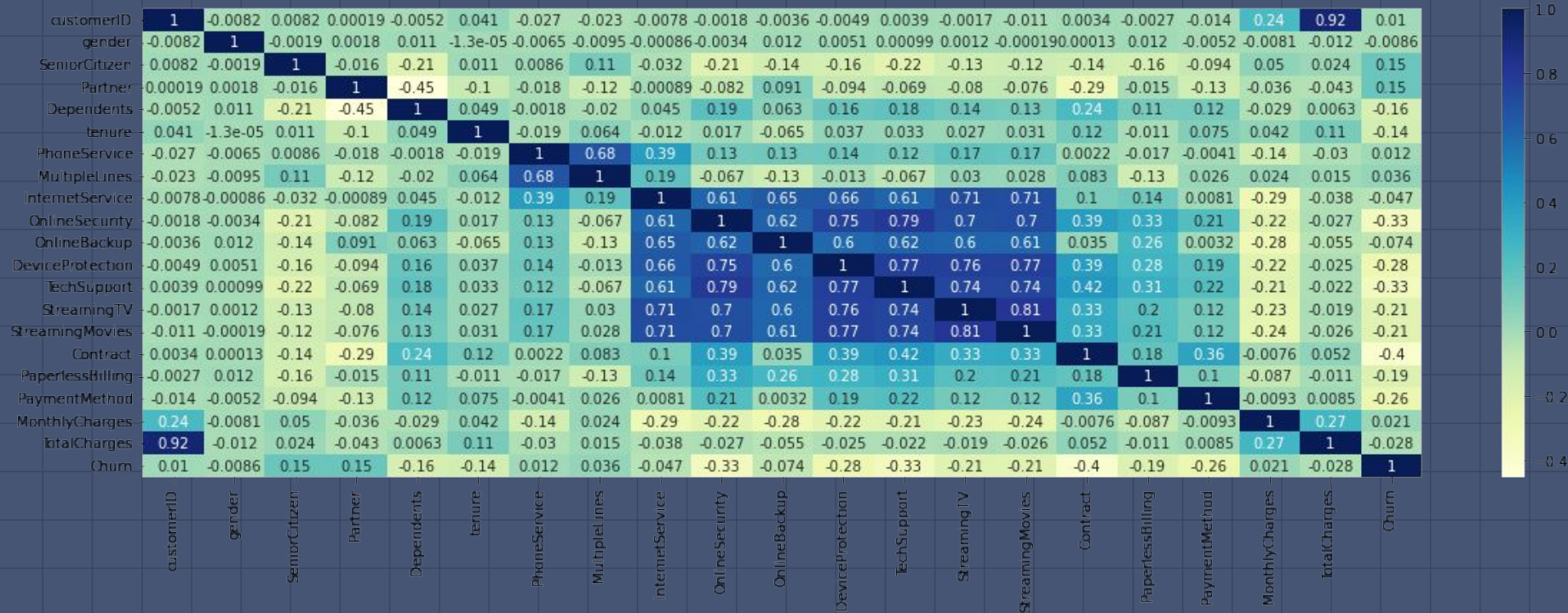
## Principais

- TotalCharges - Total pago no quarter
- MonthlyCharges - Mensalidade total paga
- Tenure - Total de meses que a pessoa é cliente
- Churn (target - se cliente cancelou)



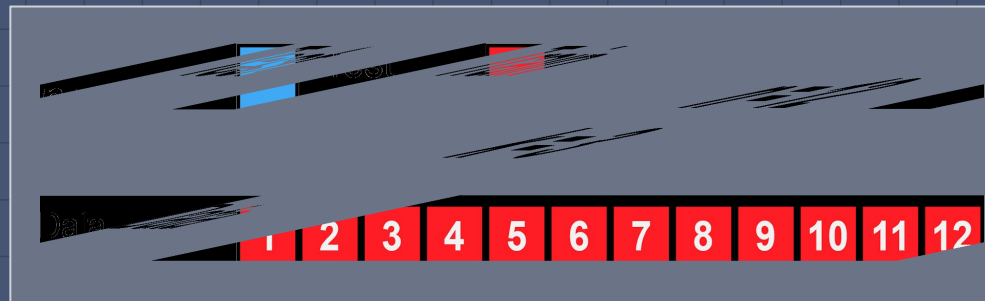
# Correlação

7



# K-fold cross validation

1. Divida aleatoriamente os dados em  $k$  subconjuntos (folds)
2. Use a parte dos dados como treino
3. Use o fold  $j$  para teste
4. Repita o processo  $k$  vezes





# K-fold cross validation

```
def random_k_folds(instance_list: List[Instance], k: int, merge_remainders=True):
    instances = instance_list.copy()
    random.Random(0).shuffle(instances)
    max_fold_size = len(instances) // k
    k_folds = [instances[i * max_fold_size:(i + 1) * max_fold_size] for i in range(k)
    ]
    remainders = instances[k * max_fold_size:]
    return remainders_merge_policy(k_folds, remainders) if merge_remainders else (
        k_folds, remainders)
```

```
def stratified_k_folds(instance_list: List[Instance], k: int, merge_remainders=True):
    instances = instance_list.copy()
    labels = {i.label() for i in instances}
    labeled_instances = [l: [i for i in instances if i.label() == l] for l in labels]
    label_folds = [random_k_folds(labeled_instances[l], k, merge_remainders=False)
    for l in labels]
    labeled_k_folds = [[e for l in f for e in l] for f, r in label_folds]
    k_folds = [[e for l in labeled_k_folds for e in l[i * len(l) // k:(i + 1) * len(l)
    ] // k]] for i in range(k)]
    k_folds = [sorted(k, key=operator.attrgetter('_instance_data')) for k in k_folds]
    remainders = [e for f, r in label_folds for e in r]
    return remainders_merge_policy(k_folds, remainders) if merge_remainders else (
        k_folds, remainders)
```

```
class Instance:
    def __init__(self, instance_data: Any, label: Union[str, int]):
        self._instance_data = instance_data
        self._label = label

    def data(self):
        return self._instance_data

    def label(self):
        return self._label

    def __repr__(self):
        return f"{self._instance_data}-{self._label}"
```

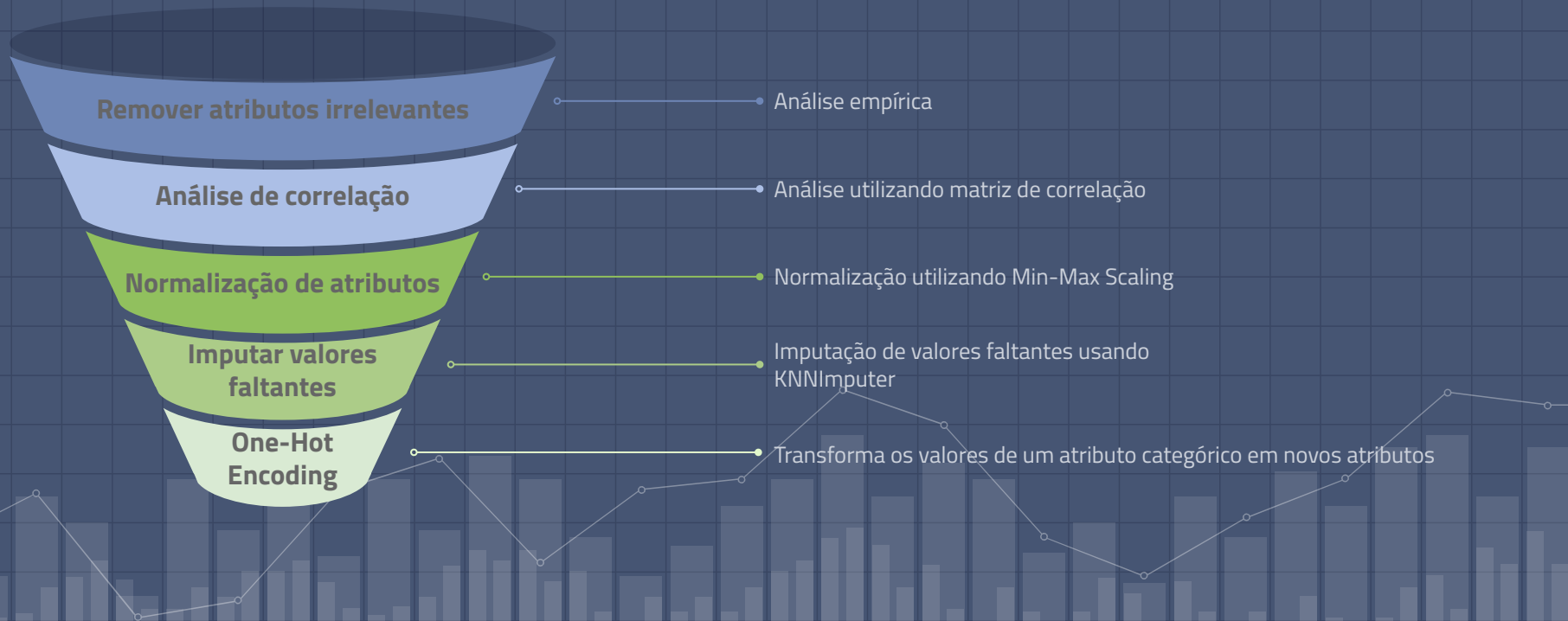
```
def return_dfs(X:pd.DataFrame, y:pd.DataFrame, k:int):
    listt = []
    for i, yy in y.items():
        """
        cria uma instancia da classe para cada instancia do dataset
        cada instancia possui o index no dataset e a label correspondente
        """
        listt.append(Instance(i, yy))
    s_folds = stratified_k_folds(listt, k)

    folds = []
    data = []
    labels = []
    for i in range(k):
        """
        para cada i em k(quantidade de folds) pega o indices que estao no fold
        e utiliza esses indices para montar os fold de dados e de labels
        e em seguida adiciona esse fold criado na lista de todos folds
        """
        indx = [a._instance_data for a in s_folds[i]]

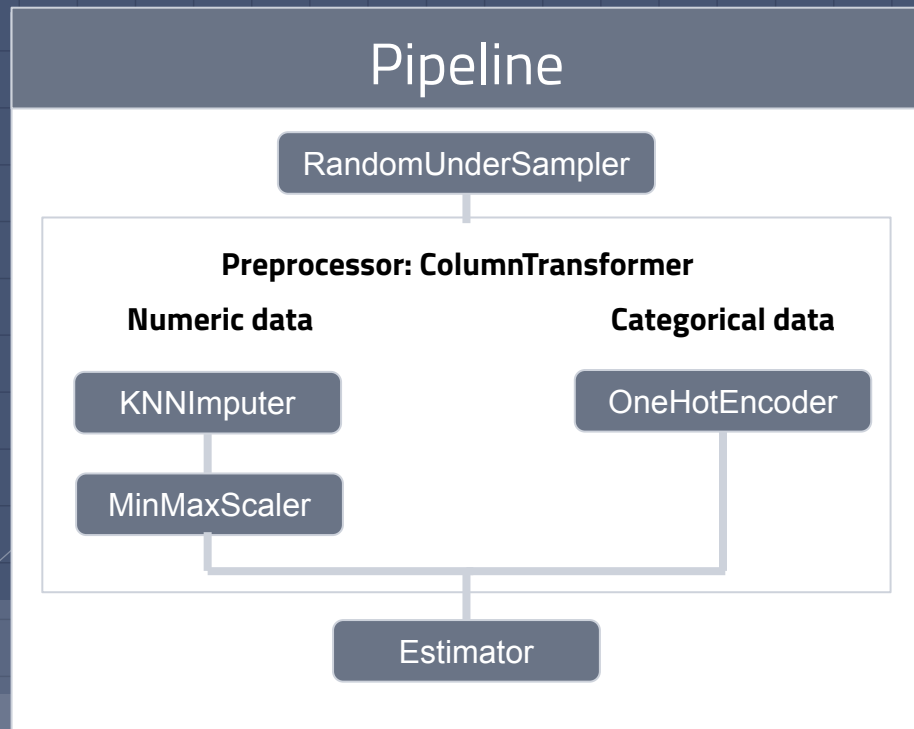
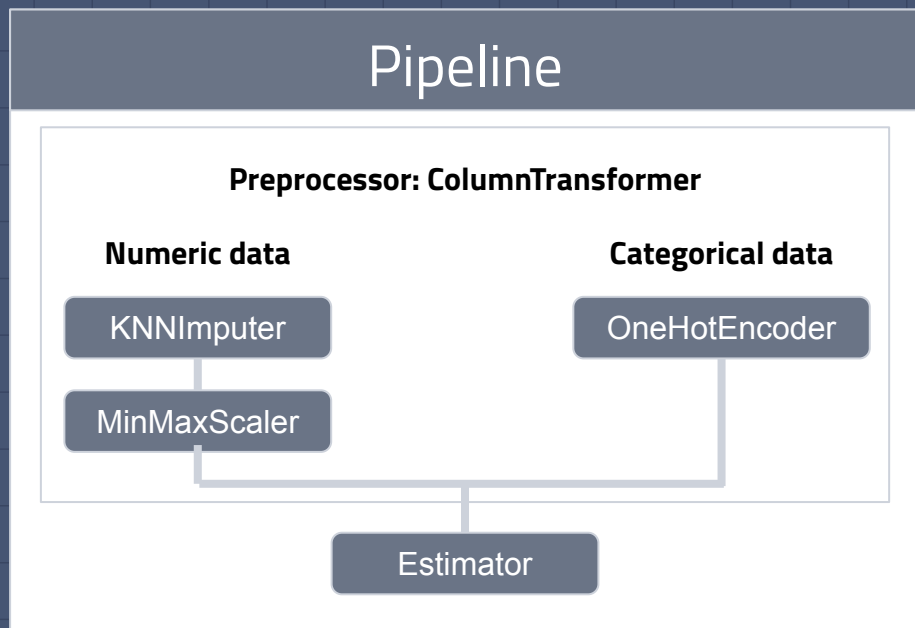
        folds.append(indx)
        data.append(X.iloc[indx])
        labels.append(y.iloc[indx])

    return data, labels
```

# PRÉ PROCESSAMENTO



# Pipeline



# Pipeline

```

numeric_features = ["MonthlyCharges", "TotalCharges", "tenure", "gender", "
Partner", "Dependents",
                    "PhoneService", "PaperlessBilling"]
numeric_transformer = Pipeline(
    steps=[("knnImp", KNNImputer(n_neighbors=3)), ("scaler", MinMaxScaler())])
categorical_features = ["InternetService", "PaymentMethod", 'MultipleLines',
'OnlineSecurity', 'OnlineBackup',
                        'DeviceProtection', 'TechSupport', 'StreamingMovies',
'Contract']
categorical_transformer = Pipeline(
    steps=[("OneHot", OneHotEncoder(handle_unknown="ignore"))])
preprocessor = ColumnTransformer(
    transformers=[
        ("num", numeric_transformer, numeric_features),
        ("cat", categorical_transformer, categorical_features)
    ]
)
if rus:
    pipe = Pipeline(
        steps=[("rus", RandomUnderSampler(random_state=0)), ("pp",
preprocessor), ("clf", clf)])
else:
    pipe = Pipeline(
        steps=[("pp", preprocessor), ("clf", clf)])

search = RandomizedSearchCV(pipe, param_grid[i], verbose=True, n_iter=20,
refit=True, random_state=0)
search.fit(X_train, y_train)

y_pred = search.predict(X_test)

y_test = np.array(y_test)

for metric in ['accuracy', 'precision', 'recall', 'f1measure']:
    if metric not in metadata[run][type(clf).__name__]:
        metadata[run][type(clf).__name__][metric] = {}
    metadata[run][type(clf).__name__]['accuracy'] = accuracy(y_test, y_pred)
    metadata[run][type(clf).__name__]['precision'] = precision(y_test, y_pred)
    metadata[run][type(clf).__name__]['recall'] = recall(y_test, y_pred)
    metadata[run][type(clf).__name__]['f1measure'] = f1Measure(y_test, y_pred)
return metadata

```

# PseudoCódigo

---

**Algorithm 1** K-fold Cross Validation

---

```
1:  $X, y \leftarrow \text{return\_dfs}(\text{data}_X, \text{data}_y, k)$   
2:  $\text{all\_metadata} \leftarrow \emptyset$   
3: for  $X_{\text{test}}, y_{\text{test}} \in X, y$  do  
4:    $X_{\text{train}} \leftarrow X - X_{\text{test}}$   
5:    $y_{\text{train}} \leftarrow y - y_{\text{test}}$   
6:    $\text{metadata} \leftarrow \text{Pipeline}(X_{\text{train}}, y_{\text{train}}, X_{\text{test}}, y_{\text{test}})$   
7:    $\text{all\_metadata} \leftarrow \text{all\_metadata} \cup \text{metadata}$   
8: end for  
9: return  $\text{all\_metadata}$ 
```

---

# ALGORITMOS ESCOLHIDOS

## **SVM - Support Vector Machine**

Algoritmo classificador linear binário não probabilístico. Procura o hiperplano entre dados de duas classes

## **XGBoost**

Algoritmo baseado em árvores de decisão com GradientBoosting.

## **GaussianNB**

Naive Bayes Gaussian.

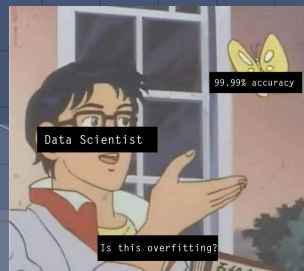
# MÉTRICAS

	Desempenho médio	Desvio padrão
Algoritmo 1	10	20
Algoritmo 2	30	15
Algoritmo 3	5	24

# MÉTRICAS (Acurácia e Precisão) - UnderSampling

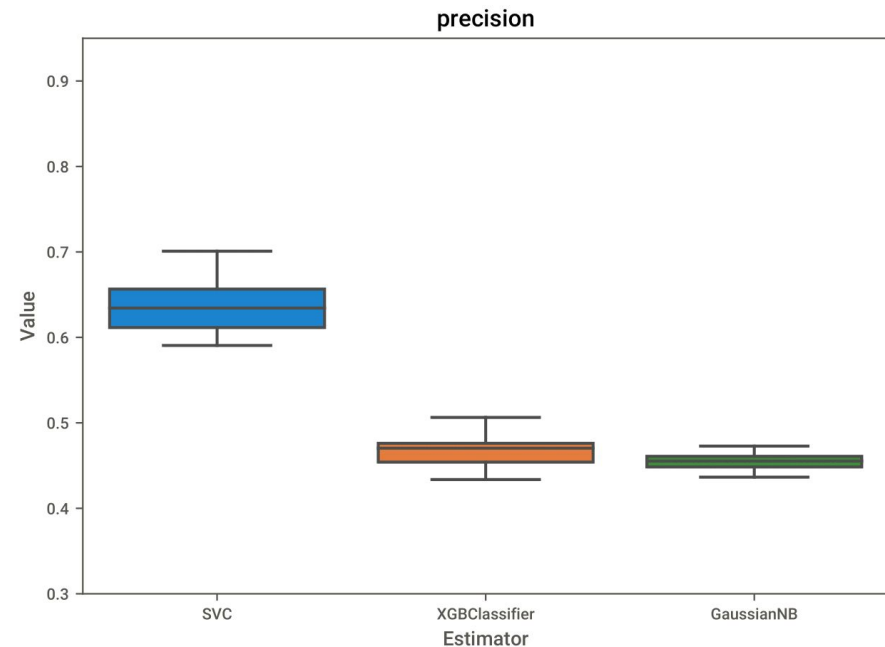
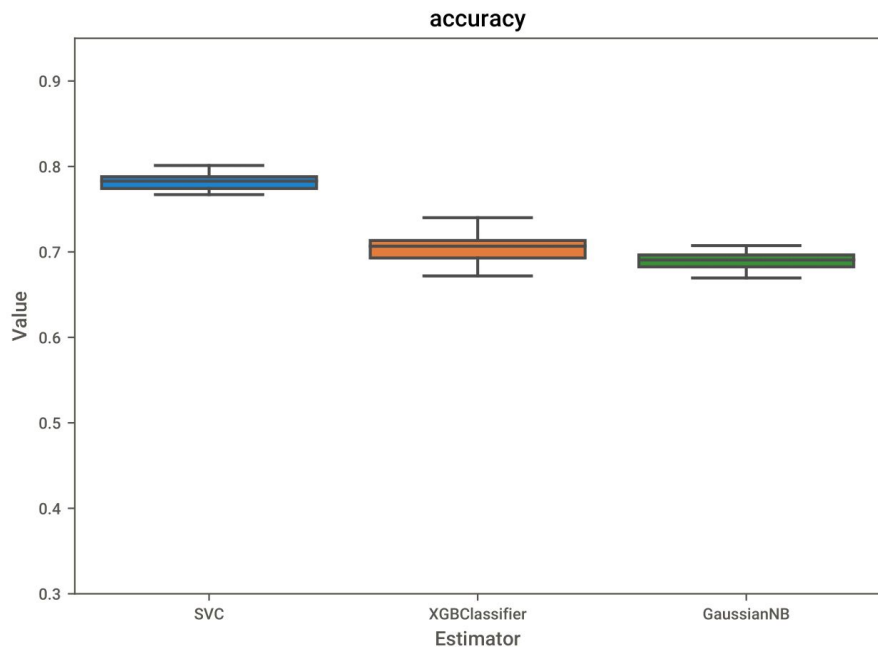
	Desempenho médio	Desvio padrão
SVC	0.782765	0.011507
XGBOOST	0.704531	0.019574
GaussianNB	0.689197	0.012587

	Desempenho médio	Desvio padrão
SVC	0.638734	0.035831
XGBOOST	0.467399	0.020973
GaussianNB	0.454311	0.011943





# MÉTRICAS (Acurácia e Precisão) - UnderSampling

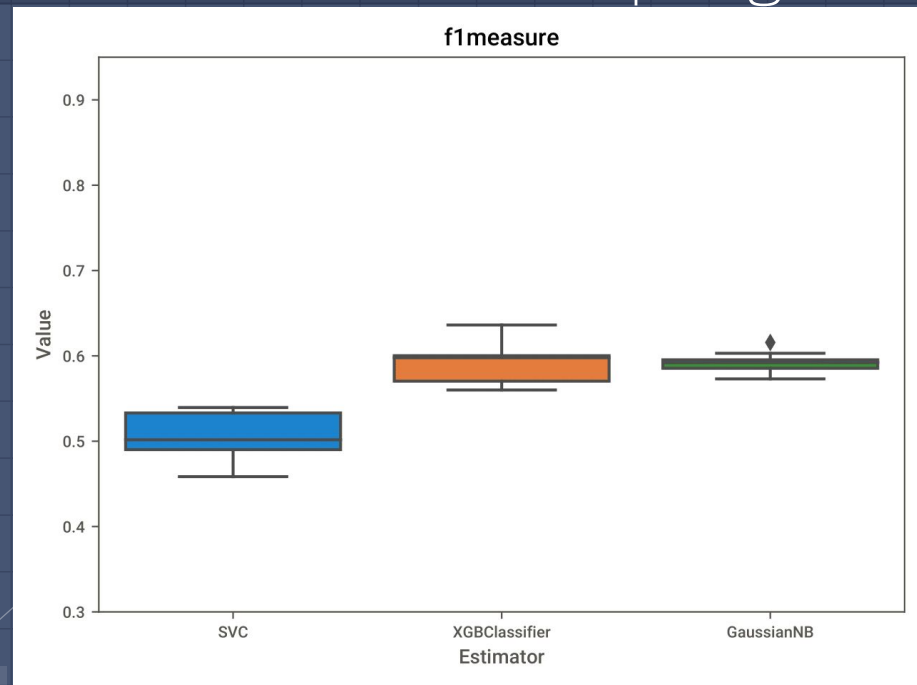
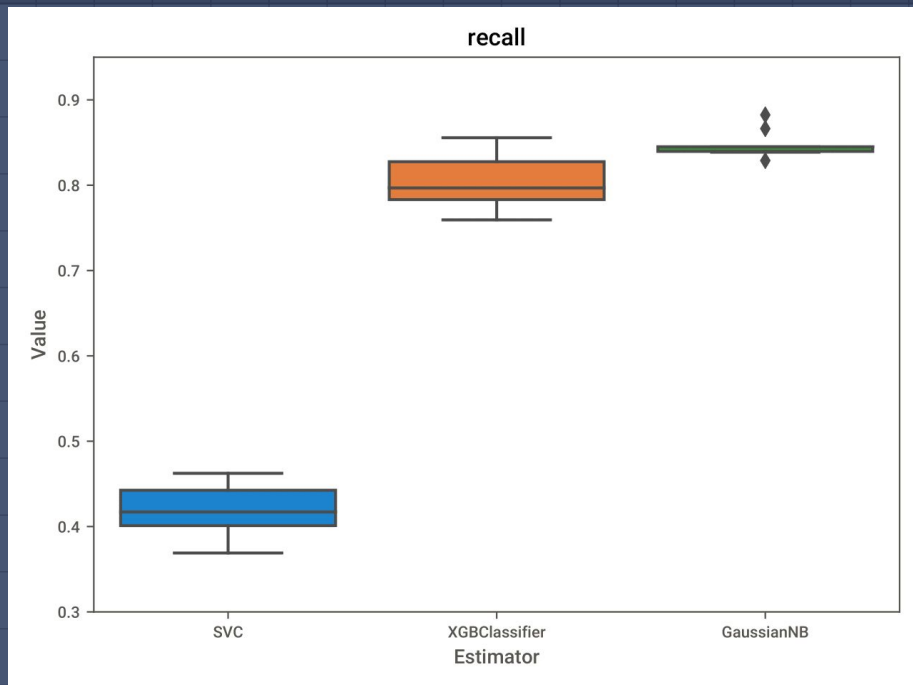


# MÉTRICAS (Recall e F1Measure) - UnderSampling

	Desempenho médio	Desvio padrão
SVC	0.419499	0.029269
XGBOOST	0.803631	0.02978
GaussianNB	0.847507	0.015438

	Desempenho médio	Desvio padrão
SVC	0.505901	0.028015
XGBOOST	0.590902	0.023053
GaussianNB	0.591464	0.01246

# MÉTRICAS (Recall e F1Measure) - UnderSampling



# MÉTRICAS (Acurácia e Precisão)

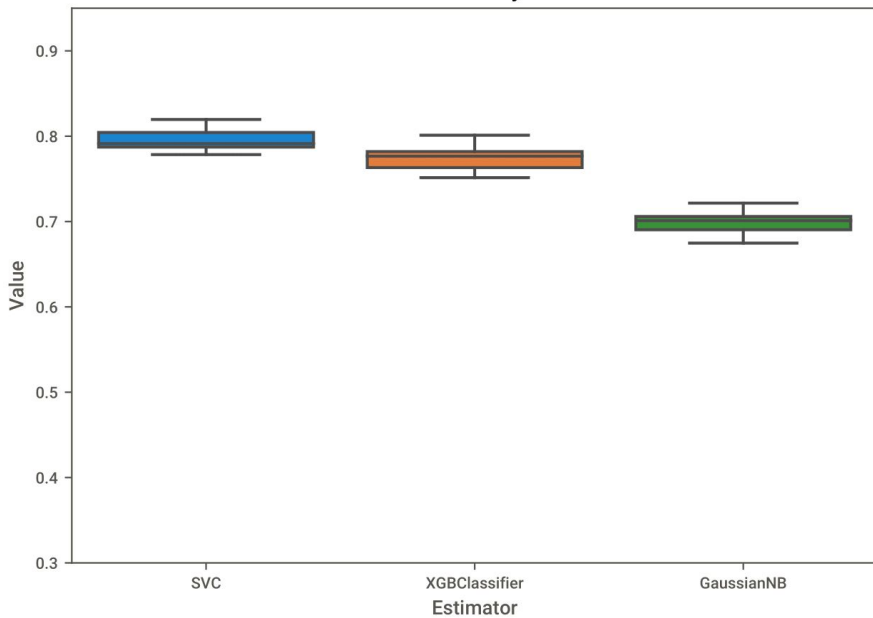
	Desempenho médio	Desvio padrão
SVC	0.795544	0.013279
XGBOOST	0.773534	0.015294
GaussianNB	0.698849	0.013008

	Desempenho médio	Desvio padrão
SVC	0.643444	0.029137
XGBOOST	0.563995	0.02542
GaussianNB	0.46294	0.013204

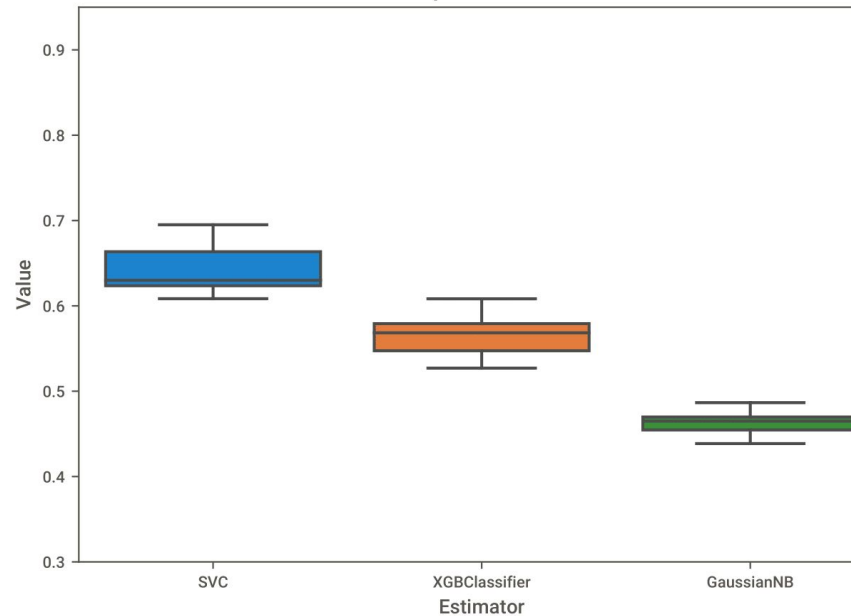
# MÉTRICAS (Acurácia e Precisão)

21

accuracy



precision



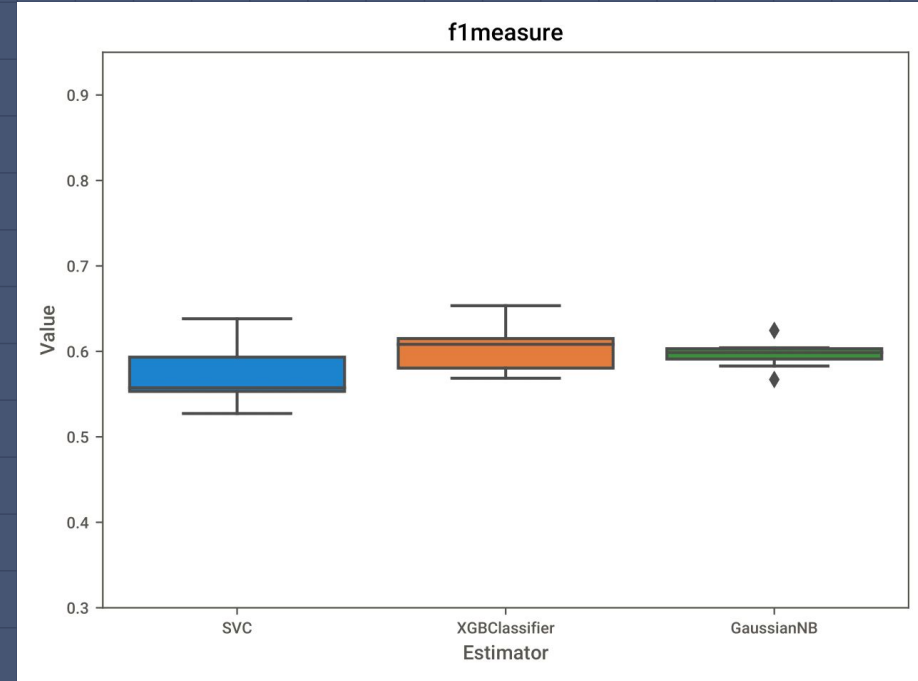
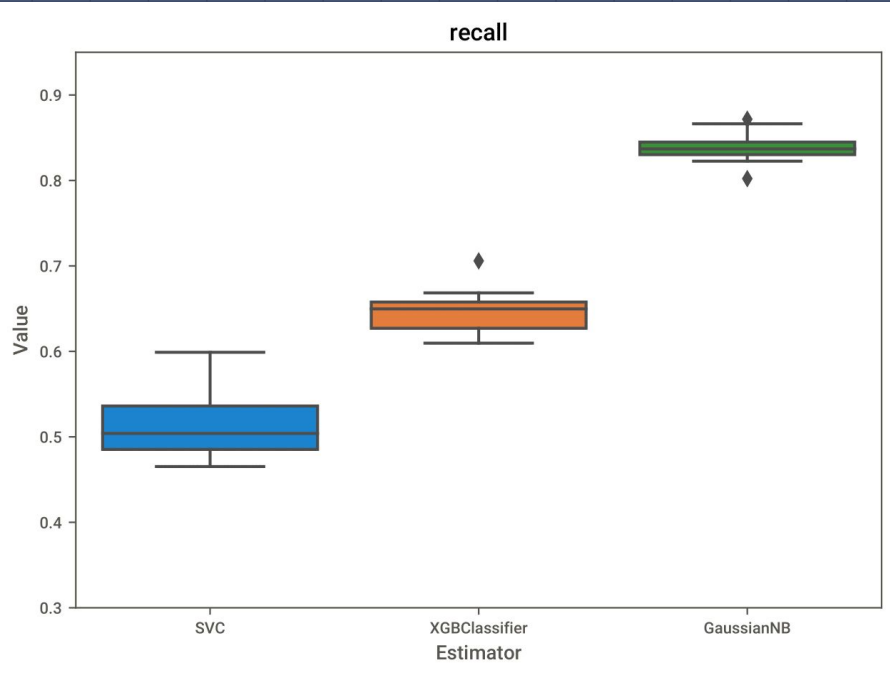
# MÉTRICAS (Recall e F1Measure)

	Desempenho médio	Desvio padrão
SVC	0.514174	0.040452
XGBOOST	0.646852	0.028786
GaussianNB	0.838943	0.020164

	Desempenho médio	Desvio padrão
SVC	0.571172	0.033598
XGBOOST	0.60254	0.02644
GaussianNB	0.596593	0.015013

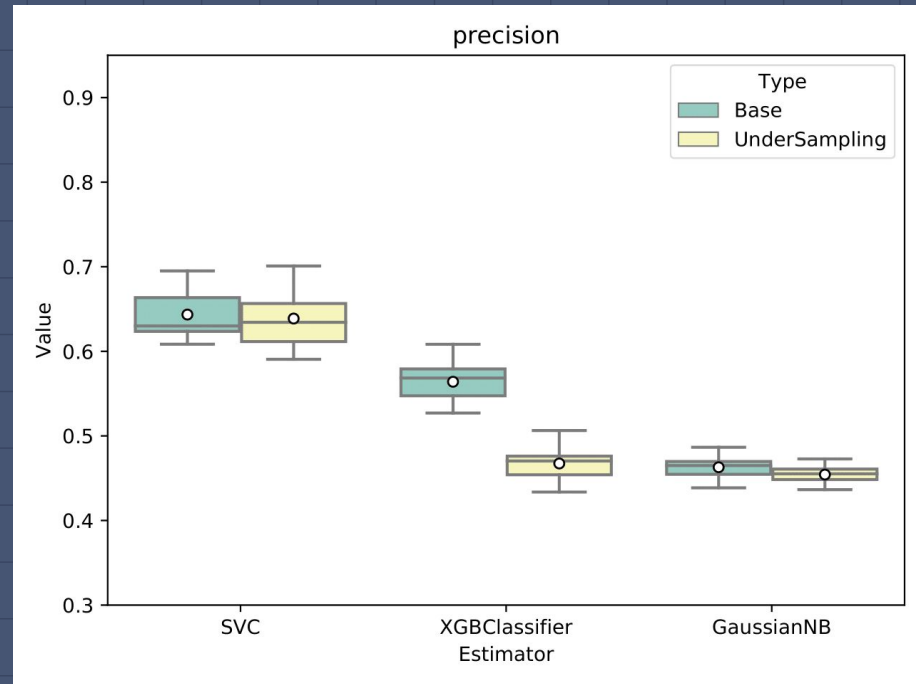
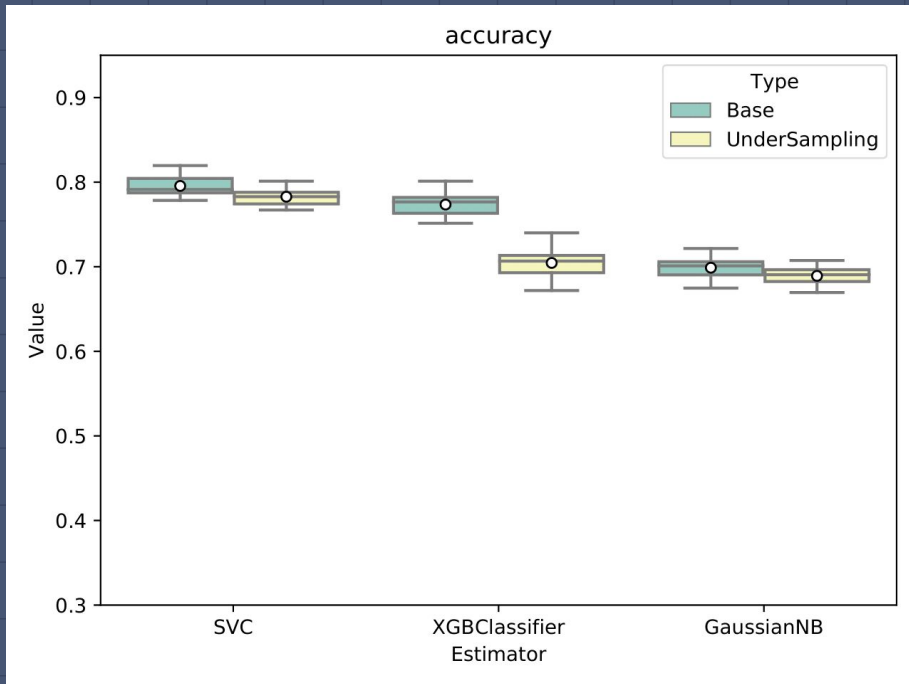
# MÉTRICAS (Recall e F1 Measure)

23



# MÉTRICAS (Acurácia e Precisão) - Comparação

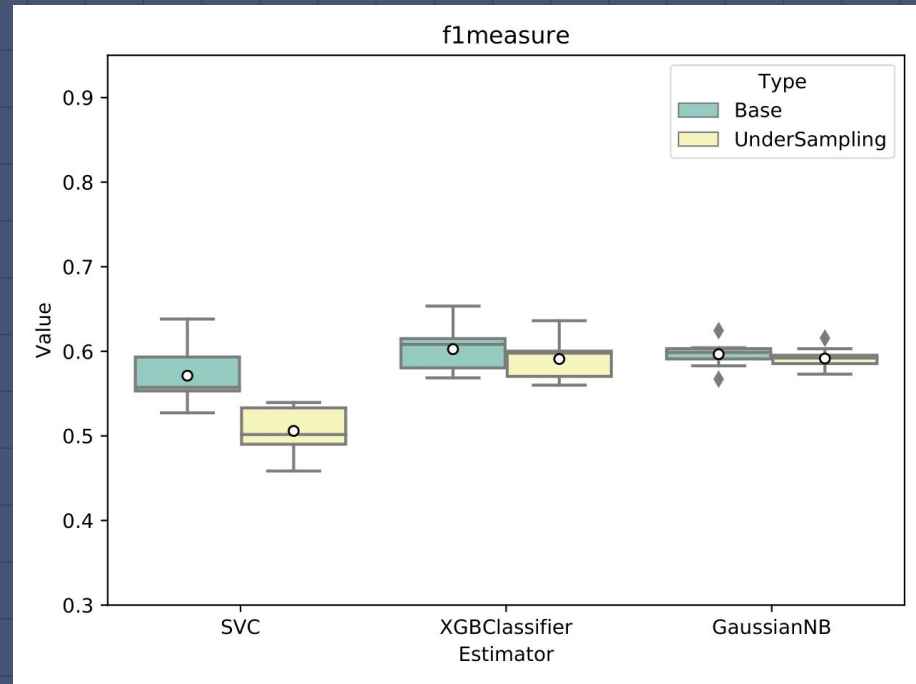
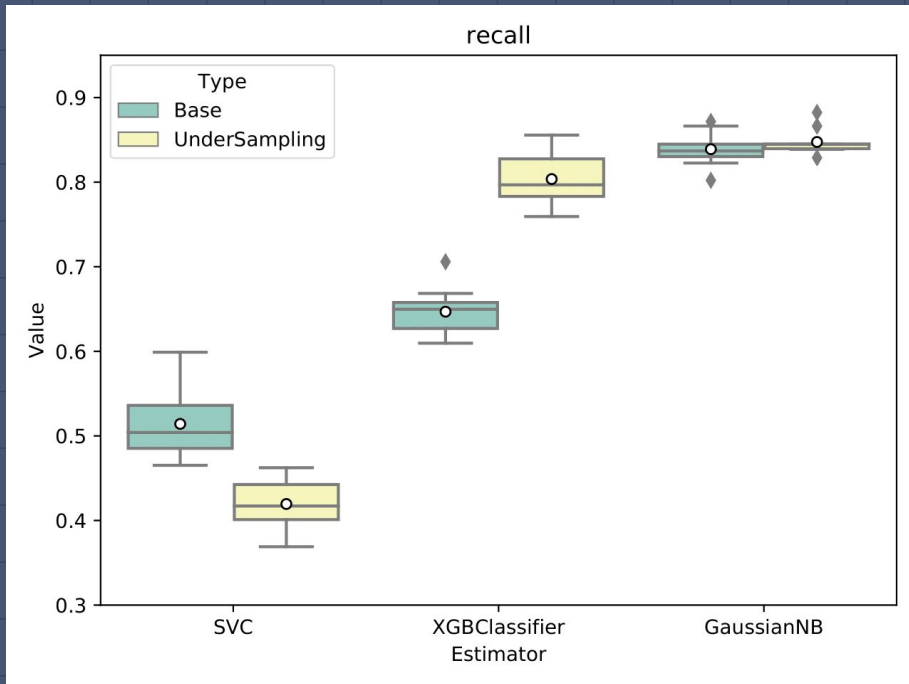
24





# MÉTRICAS (Recall e F1Measure) – Comparação

25



# Conclusão

Como o *churn* gera impacto de caráter econômico, acreditamos que a melhor opção de modelo, para este caso, seria o Gaussian Naïve Bayes, pois apresentou o melhor desempenho para **F-measure** e **Recall**.

Acreditamos que a priorização do **Recall** é interessante por priorizar manter os clientes da empresa. E a da **F-measure** por permitir ajuste de parâmetros baseado em características do mercado como custo/retorno de prevenção de *churn*, pois o quanto é interessante poupar caixa da empresa (**reduzir FP**) ou aumentar a arrecadação mantendo clientes (**reduzir FN**) pode variar de acordo com as condições de mercado e custos de aquisição ou manutenção de cliente.

Me: \*uses machine learning\*  
Machine: \*learns\*  
Me:



# Obrigado!

**Perguntas?**

