# Decision Trees



How can a dataset be represented as a tree?

| Age | Marital-Status | Educational | Sex | High-Income |
|-----|----------------|-------------|-----|-------------|
| 28 | Never-Married | Bachelors | Female | <=50k |
| 46 | Never-Married | Assoc-acdm | Female | <=50k |
| 35 | Married-civ-spouse | Some-college | Male | <=50k |
| 27 | Married-civ-spouse | Bachelors | Male | >50k |
| 59 | Divorced | Some-college | Female | <=50k |

# Decision Tree (classification)

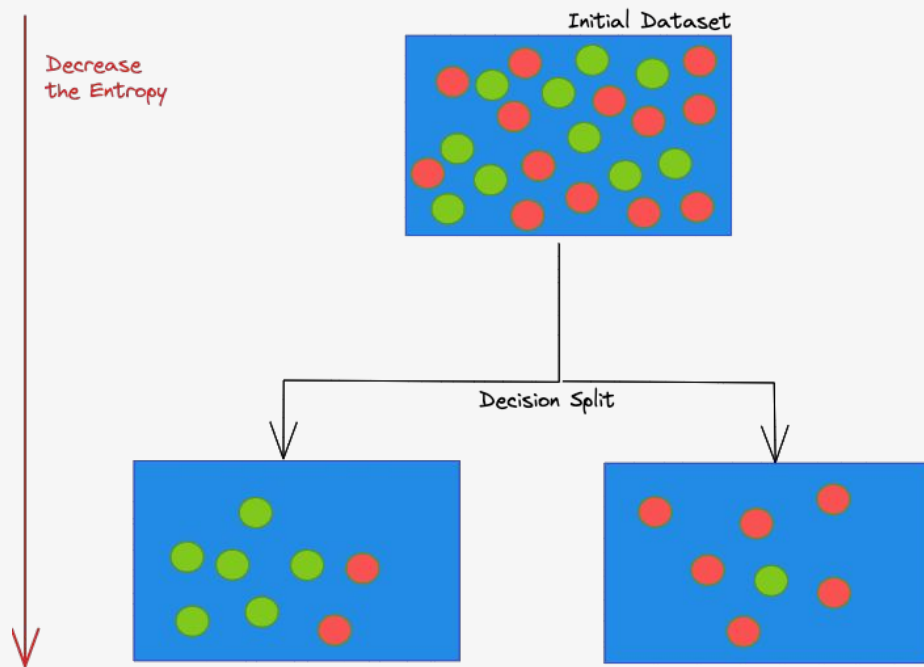# How can we split the tree?

Age

Algorithm used in Decision Trees

1. ID3 (Entropy)
2. Gini Index
3. Chi-Square
4. Reduction in Variance
   a. C4.5, pruning
5. ...

Entropy is an indicator of how messy your data is.

# Why Entropy in Decision Trees?



Decrease the Entropy

Initial Dataset

Decision Split

- The goal is to tidy the data.
- You try to separate your data and group the samples together in the classes they belong to.
- You maximize the purity of the groups as much as possible each time you create a new node of the tree
- Of course, at the end of the tree, you want to have a clear answer.

# Mathematical Definition of Entropy



Suppose a set of N items, these items fall into two categories:

+ gain > 50k (k)
- gain <= 50k (m)

$$p = \frac{k}{N}, q = \frac{m}{N}$$

$$Entropy = -p \log p - q \log q$$

# Generalization

Feature x

$$E(x) = - \sum_{i=1}^{c} P(x_i) \log_b P(x)$$

$P(x_i)$ is the fraction of examples in a given class i

| | |
|---|---|
| <= 50k. | 17288 |
| > 50k. | 5487 |

```
from scipy.stats import entropy
entropy(df_train.high_income.value_counts(), base=2)
0.7965702796015677
```

# Entropy using the frequency table of two attributes

| High Income | | |
|---|---|---|
| | **<= 50k** | **> 50k** |
| **<=37** | 7206 | 3883 | 11089 (48%) |
| **>37** | 10082 | 1604 | 11686 (52%) |

**Age**

$$E(T \mid X) = \sum_{c \in X} \frac{|X_c|}{|X|} E(T \mid X_c)$$

```
cross = pd.crosstab(
        df_train.age <= df_train.age.median(),
        df_train.high_income)
```

```
0.486894 * entropy(cross.iloc[0], base=2) \
+ 0.513106 * entropy(cross.iloc[1], base=2)
0.7509335429830957
```

# Information Gain

$$IG (T,X) = E(T) - E(T|X)$$

Information Gain from X on T

The information gain is based on the **decrease in entropy after a dataset is split** on an attribute.

Constructing a decision tree is all about finding attribute that returns the **highest information gain** (i.e., the most homogeneous branches).
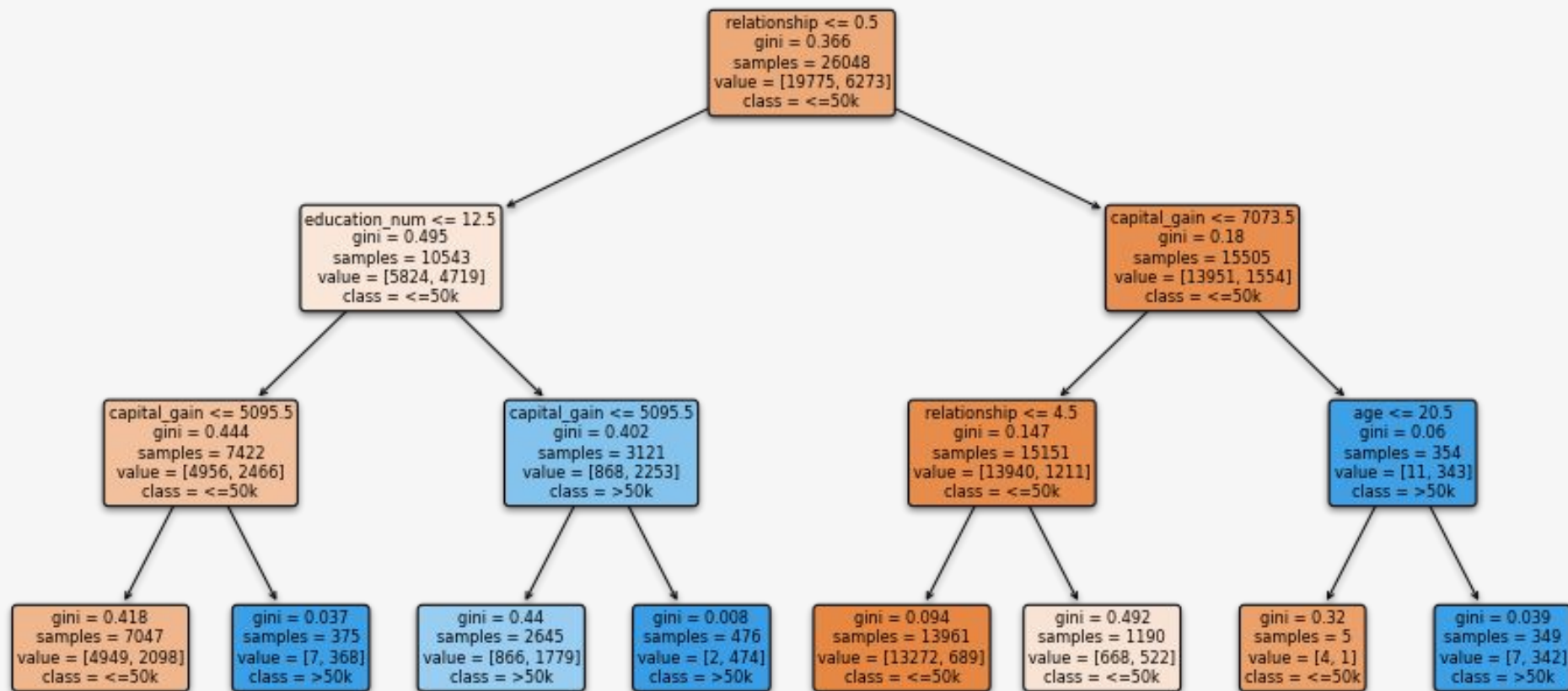
$$Gini(x) = 1 - \sum_{i=1}^{c} P(x_i)^2$$

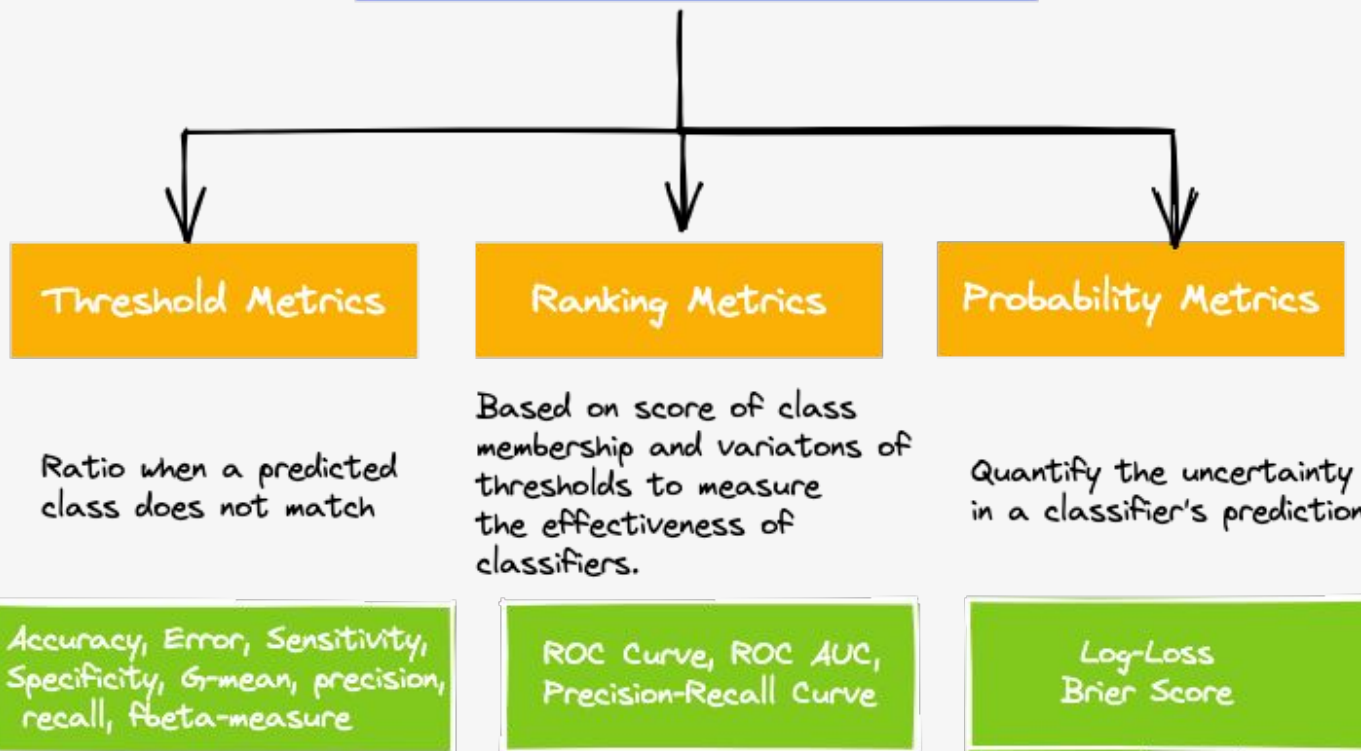$$Entropy(x) = -\sum_{i=1}^{c} P(x_i) \log_b P(x_i)$$

Gini index or Entropy is the criterion for calculating **Information Gain**. Both of them are measures of impurity of a node.

```
from sklearn.tree import plot_tree
```



relationship <= 0.5
gini = 0.366
samples = 26048
value = [19775, 6273]
class = <=50k

education_num <= 12.5
gini = 0.495
samples = 10543
value = [5824, 4719]
class = <=50k

capital_gain <= 7073.5
gini = 0.18
samples = 15505
value = [13951, 1554]
class = <=50k

capital_gain <= 5095.5
gini = 0.444
samples = 7422
value = [4956, 2466]
class = <=50k

capital_gain <= 5095.5
gini = 0.402
samples = 3121
value = [868, 2253]
class = >50k

relationship <= 4.5
gini = 0.147
samples = 15151
value = [13940, 1211]
class = <=50k

age <= 20.5
gini = 0.06
samples = 354
value = [11, 343]
class = >50k

gini = 0.418
samples = 7047
value = [4949, 2098]
class = <=50k

gini = 0.037
samples = 375
value = [7, 368]
class = >50k

gini = 0.44
samples = 2645
value = [866, 1779]
class = >50k

gini = 0.008
samples = 476
value = [2, 474]
class = >50k

gini = 0.094
samples = 13961
value = [13272, 689]
class = <=50k

gini = 0.492
samples = 1190
value = [668, 522]
class = <=50k

gini = 0.32
samples = 5
value = [4, 1]
class = <=50k

gini = 0.039
samples = 349
value = [7, 342]
class = >50k

**Taxonomy of Classifier Evaluation Metrics**

**Threshold Metrics**

Ratio when a predicted class does not match

Accuracy, Error, Sensitivity, Specificity, G-mean, precision, recall, fbeta-measure

**Ranking Metrics**

Based on score of class membership and variatons of thresholds to measure the effectiveness of classifiers.

ROC Curve, ROC AUC, Precision-Recall Curve

**Probability Metrics**

Quantify the uncertainty in a classifier's prediction

Log-Loss
Brier Score

# How to choose an evaluation metric?

What do you want to predict?

- Class Labels
  - Are both classes equally important?
    - Do < 80-90% of examples belong to the majority class?
      - Accuracy
    - G-Mean
  - Is the positive class more imortant?
    - Are false negative and false positives equally costly?
      - F1 Score
    - Are false positives more costly?
      - F0.5 Score
    - Are false negatives more costly?
      - F2 Score
- Probabilities
  - Do you need probabilities?
    - Brier Score
  - Do you need class labels?
    - Is the positive class more important?
      - PR AUC
    - Are both classes equally important?
      - ROC AUC

@Brownlee, Jason. Imbalanced classification with python.

# Confusion Matrix

## Expected



Predicted

|  | Positive Class (1) | Negative Class (0) |
|---|---|---|
| Positive class (1) | True Positive (TP) — Predicted / Expected | False Positive (FP) — Predicted / Expected |
| Negative class (0) | False Negative (FN) — Predicted / Expected | True Negative (TN) — Predicted / Expected |

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

$$Error = 1 - Accuracy$$

# Confusion Matrix



Expected

Positive Class (1)     Negative Class (0)

Predicted / Positive class (1) / Negative class (0)

True Positive (TP)
Predicted   Expected

False Positive (FP)
Predicted   Expected

False Negative (FN)
Predicted   Expected

True Negative (TN)
Predicted   Expected

$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

$$\text{Specificity} = \frac{TN}{FP + TN}$$

$$\text{G-mean} = \sqrt{\text{Sensitivity} \times \text{Specificity}}$$

# Confusion Matrix

## Expected

|  | Positive Class (1) | Negative Class (0) |
|---|---|---|
| **Positive class (1)** | True Positive (TP)<br>Predicted    Expected | False Positive (FP)<br>Predicted    Expected |
| **Negative class (0)** | False Negative (FN)<br>Predicted    Expected | True Negative (TN)<br>Predicted    Expected |

**Predicted**

$$\text{Precision} = \frac{TP}{TP + FP}$$
(positive predicte value - PPV)

$$\text{Precision} = \frac{TN}{TN + FN}$$
(negative predicte value - NPV)

$$\text{Recall} = \frac{TP}{TP + FN}$$

# Confusion Matrix

## Expected

**Predicted**

|  | Positive Class (1) | Negative Class (0) |
|---|---|---|
| **Positive Class (1)** | True Positive (TP)<br>Predicted 😀 Expected 😀 | False Positive (FP)<br>Predicted 😀 Expected 😀 |
| **Negative Class (0)** | False Negative (FN)<br>Predicted 😀 Expected 😀 | True Negative (TN)<br>Predicted 😀 Expected 😀 |

$$\text{Fbeta-measure} = \frac{(1 + \beta^2) \times \text{Precision} \times \text{Recall}}{\beta^2 \times \text{Precision} + \text{Recall}}$$
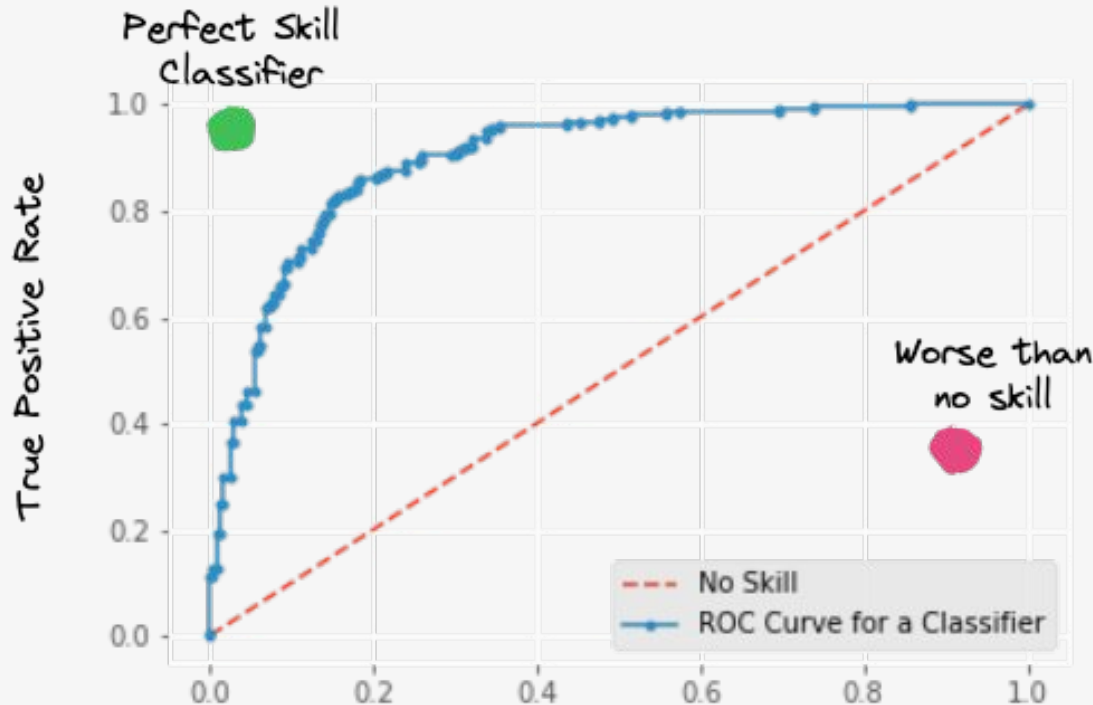
$$\beta == \begin{cases} 0.5, & \text{more weight on precision} \\ 1.0, & \text{balance on weight PR and RE} \\ 2.0, & \text{less weight on precision} \end{cases}$$

Rank metrics are more concerned with evaluating classifiers based on **how effective** they are at separating classes.

These metrics require that a **classifier predicts a score** or a probability of class membership. From this score, **different thresholds** can be applied to **test the effectiveness of classifiers**. Those models that maintain a good score across a range of thresholds will have good class separation and will be ranked higher.

Ranking Metrics

Perfect Skill Classifier

**True Positive Rate**

$$TPR = \frac{TP}{TP + FN}$$

Worse than no skill

1.0

0.8

0.6

0.4

0.2

0.0

0.0    0.2    0.4    0.6    0.8    1.0

- - - No Skill
—•— ROC Curve for a Classifier

**False Positive Rate**

$$FPR = \frac{FP}{FP + TN}$$

Expected

Positive Class (1)        Negative Class (0)

Predicted

Positive Class (1)

True Positive (TP)        False Positive (FP)

Predicted    Expected        Predicted    Expected

Negative Class (0)

False Negative (FN)        True Negative (TN)
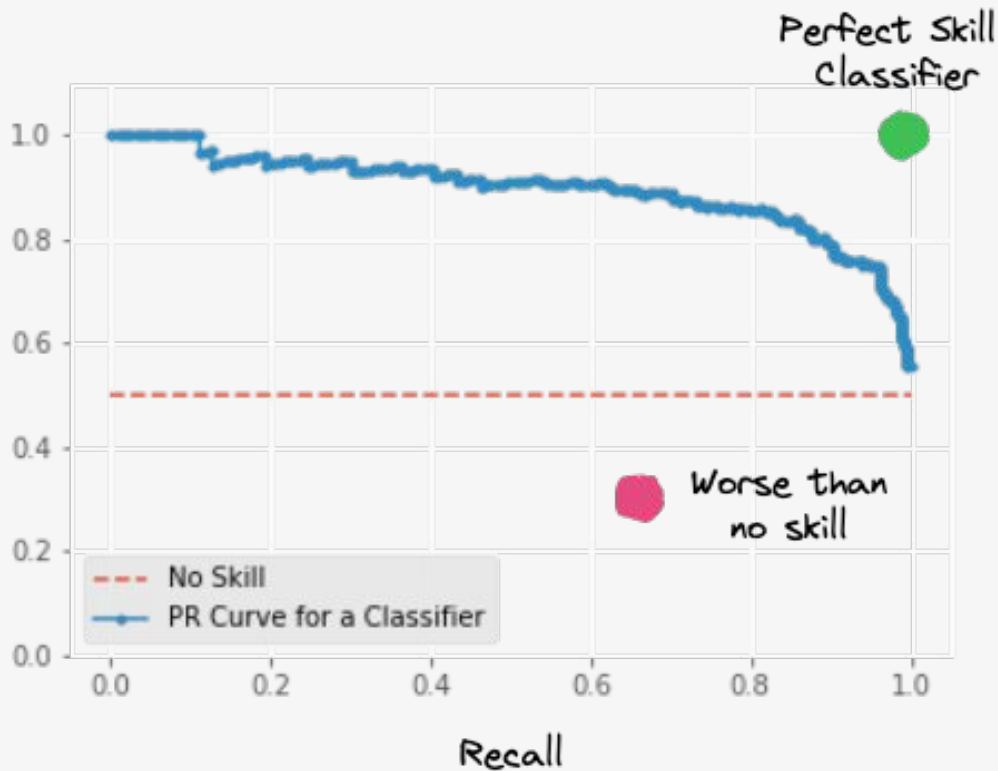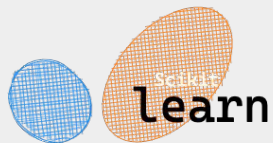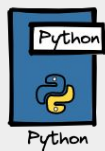
Predicted    Expected        Predicted    Expected

# Precision-Recall (PR) Curve

Perfect Skill Classifier

Worse than no skill

Legend:
- No Skill
- PR Curve for a Classifier

$$Precision = \frac{TP}{TP + FP}$$

Axis labels: Precision (y-axis), Recall (x-axis)

$$Recall = \frac{TP}{TP + FN}$$

Expected

| | Positive Class (1) | Negative Class (0) |
|---|---|---|
| Positive Class (1) | True Positive (TP) | False Positive (FP) |
| Negative Class (0) | False Negative (FN) | True Negative (TN) |

Predicted

Within each cell: Predicted / Expected

Hands ON

**UCI**

## Machine Learning Repository

# Adult Data Set

*Download*: **Data Folder**, **Data Set Description**

**Abstract**: Predict whether income exceeds $50K/yr based on census data. Also known as "Census Income" dataset.
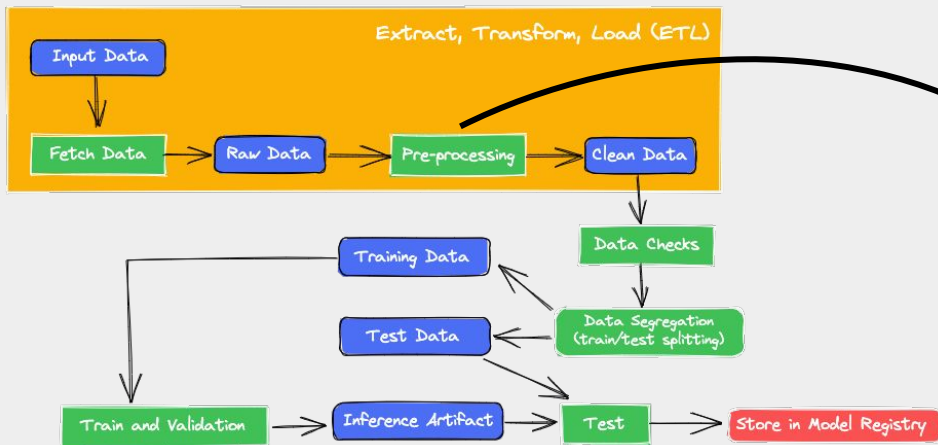
| Data Set Characteristics: | Multivariate | Number of Instances: | 48842 | Area: | Social |
|---|---|---|---|---|---|
| Attribute Characteristics: | Categorical, Integer | Number of Attributes: | 14 | Date Donated | 1996-05-01 |
| Associated Tasks: | Classification | Missing Values? | Yes | Number of Web Hits: | 2437279 |

**Source:**

Donor:

Ronny Kohavi and Barry Becker
Data Mining and Visualization
Silicon Graphics.
e-mail: ronnyk '@' live.com for questions.

Exploratory Data Analysis

# Extract, Transform, Load (ETL)

```
Input Data
   │
   ▼
Fetch Data → Raw Data → Pre-processing → Clean Data
                                              │
                                              ▼
                                         Data Checks
                                              │
                                              ▼
Training Data ◄─── Data Segregation
Test Data     ◄─── (train/test splitting)

Train and Validation → Inference Artifact → Test → Store in Model Registry
```

**Raw Data**

**Data Segregation**

Train Set

Test Set

**Holdout**

Train Set

Validation Set

**cross-validation**



Split

Data Transforms

Numerical Type

Categorical Type

Change Scale | Change Distribution | Engineer

Nominal type | Ordinal type

Normalize | Power | Polynomial

Standardize | Quantile

Robust | Discretize

One Hot Encode | Label Encode

Dummy Enconde

Extract, Transform, Load (ETL)

Input Data → Fetch Data → Raw Data → Pre-processing → Clean Data

Clean Data → Data Checks → Data Segregation (train/test splitting)

Data Segregation → Training Data, Test Data

Training Data → Train and Validation → Inference Artifact → Test → Store in Model Registry