

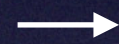
Chapter 2

End-to-End Machine Learning Project

Latitude



Longitude



District Housing Price

Per Block Group(i.e. population of 600 to 3000 people)



District	Population	Median Income	Medium Price						
----------	------------	---------------	--------------	--	--	--	--	--	--

How will they use the model?

Model output will be fed into another system to determine whether it is worth investing in certain area.



How does company benefit from model?

It guides where we make investments, which directly impacts revenue at this company.

Features

Target

Univariate regression (predict single value)

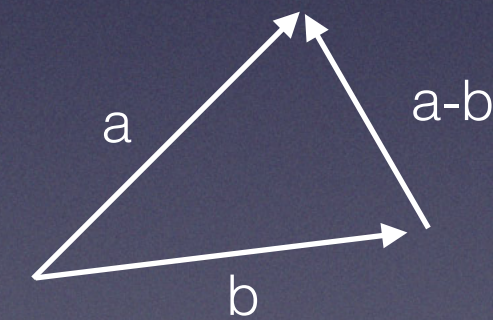
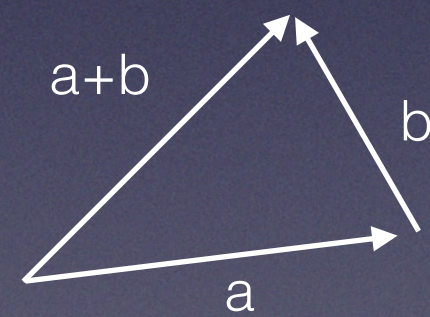
$$MAE(X, h) = \sqrt{\frac{1}{m} \sum_{i=1}^m |h(x^{(i)}) - y^{(i)}|}$$

$$RSME(X, h) = \sqrt{\frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2}$$

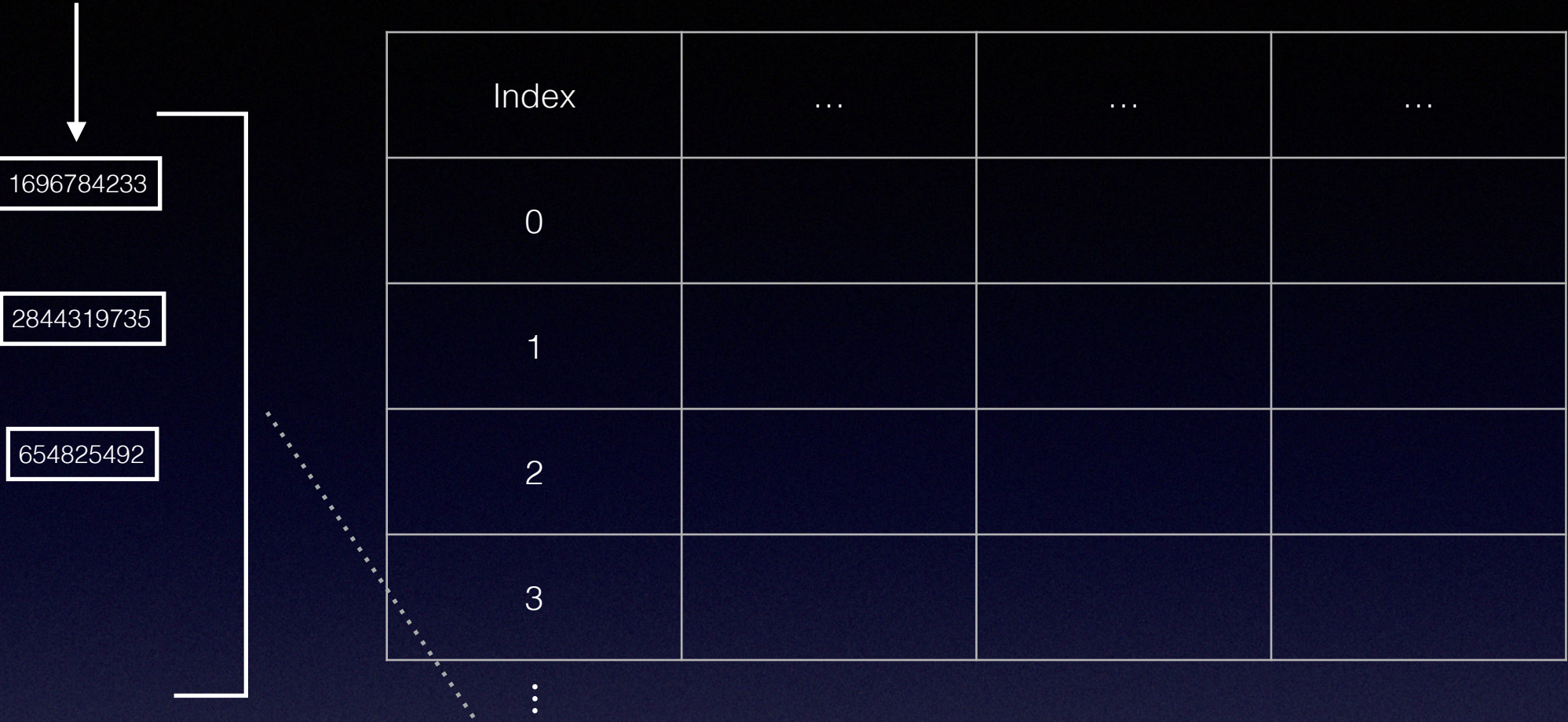
Cost functions measure distance of prediction vectors to target value vectors

RMSE - sensitive to outliers , use when outliers are negligible

MAE - not sensitive to outliers existing in dataset

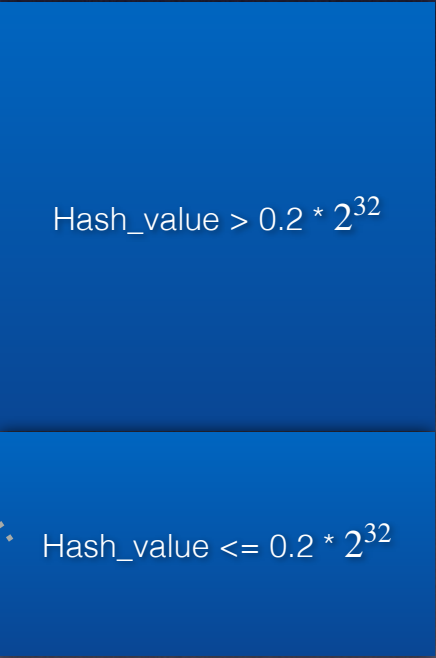


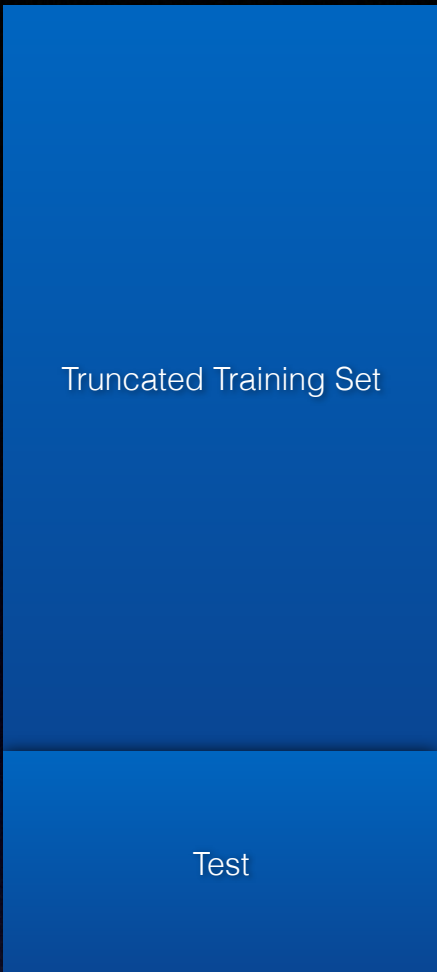
crc32 redundancy per index



crc32 redundancy per index

Create custom ID using features if dataset gets deleted





Stratified sampling - sampling which guarantees the test is representative of overall population

Design Architect reveals the median income is very important or sensitive to predict median housing prices



Capture indices which represent the median income feature(i.e. attribute) when creating Test set.

Analyze Correlation to gain insights on data

Fill missing values in datasets using imputer

Encode Categorical Attributes

Ordinal Encoder

One Hot Encoder

Embedding

Low dimensional vector which is updated during training to support model performance

Custom Encoding

Replace country code with courtly population

Feature Scaling

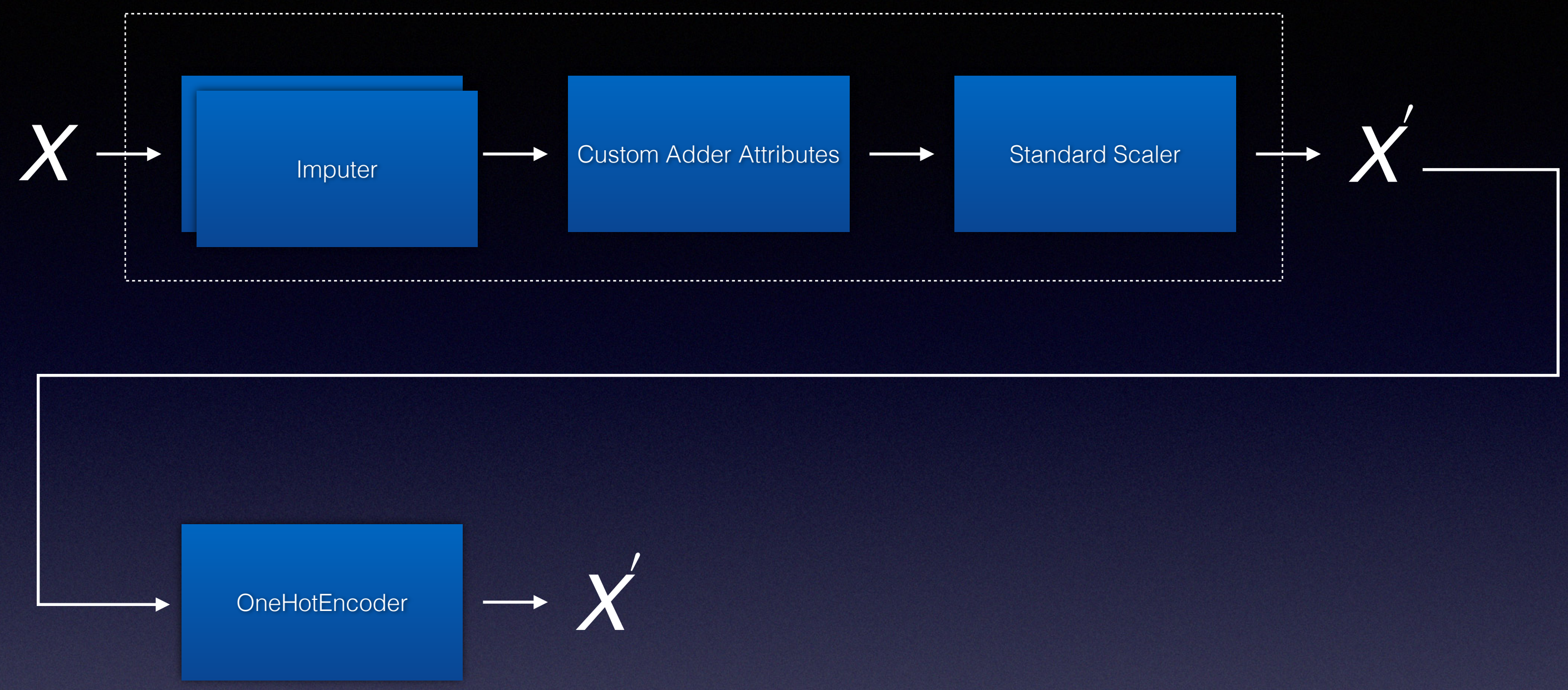
min-max

Values are rescaled to range from 0 to 1

standardization

Does not bound values to range like min-max.. This scaling option is not affected by outliers

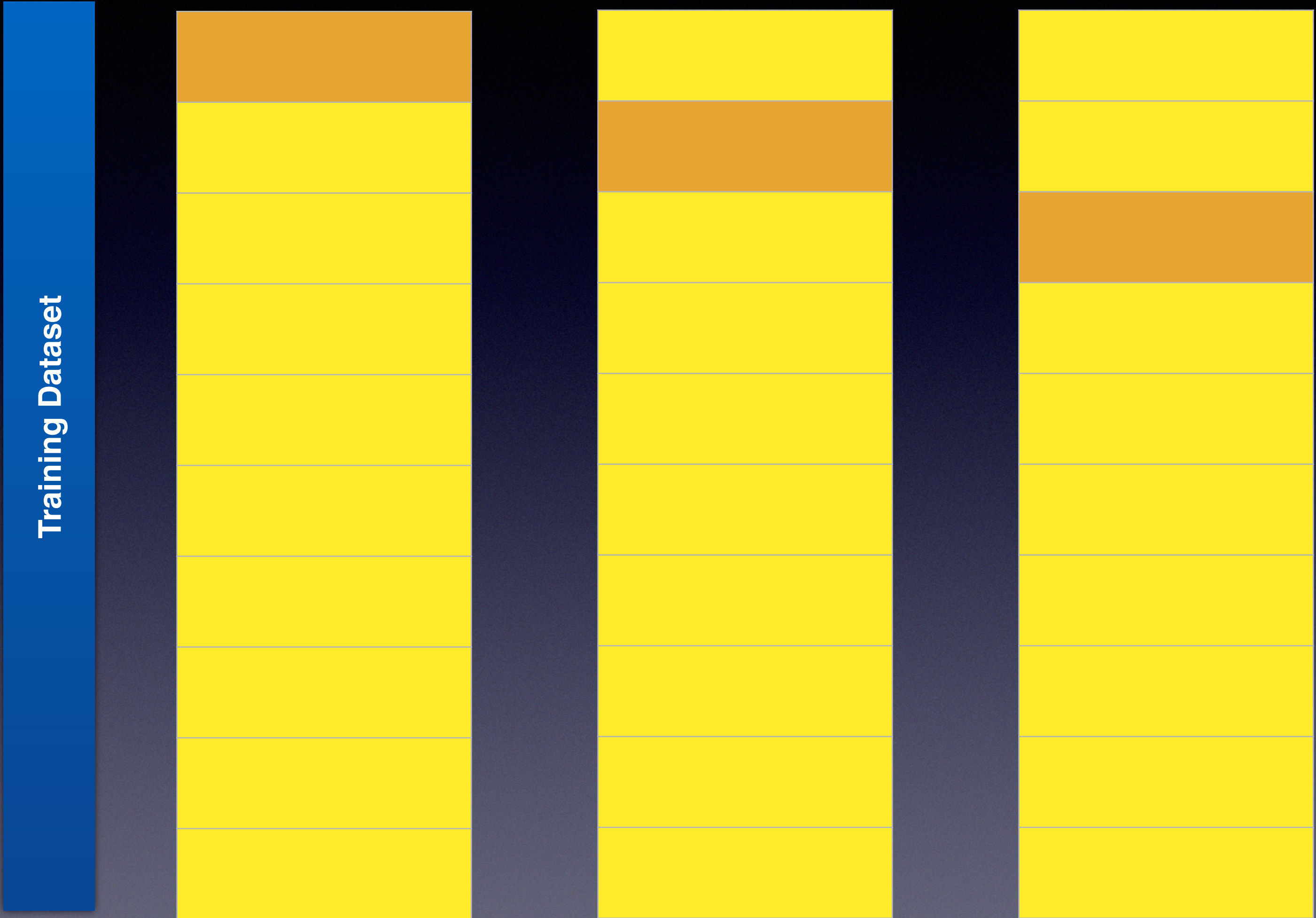
Simple Transformation Pipelines



Train Set
Evaluation Set

CV = 10

$X \rightarrow \text{[blue box]} \rightarrow \textit{Model}$



...

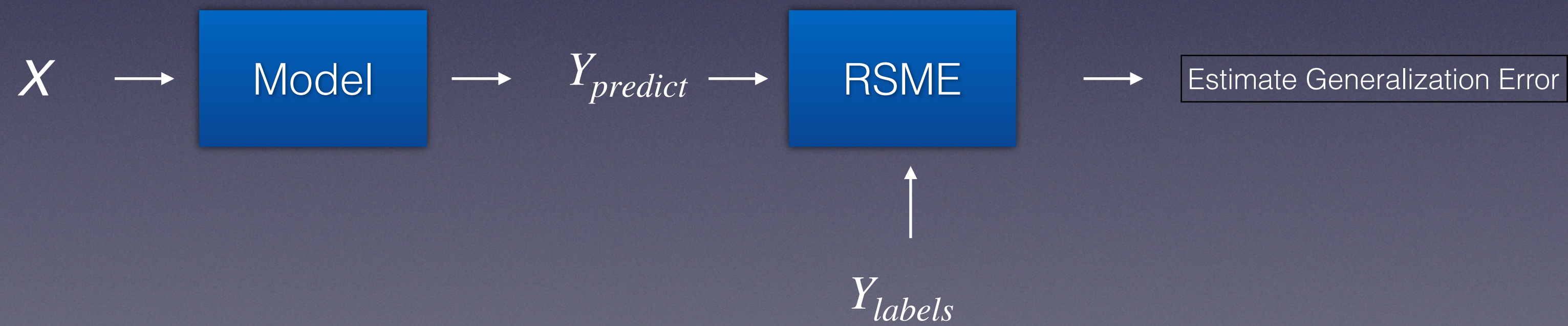
Validation Score (RMS)

Validation Score (RMS)

Validation Score (RMS)

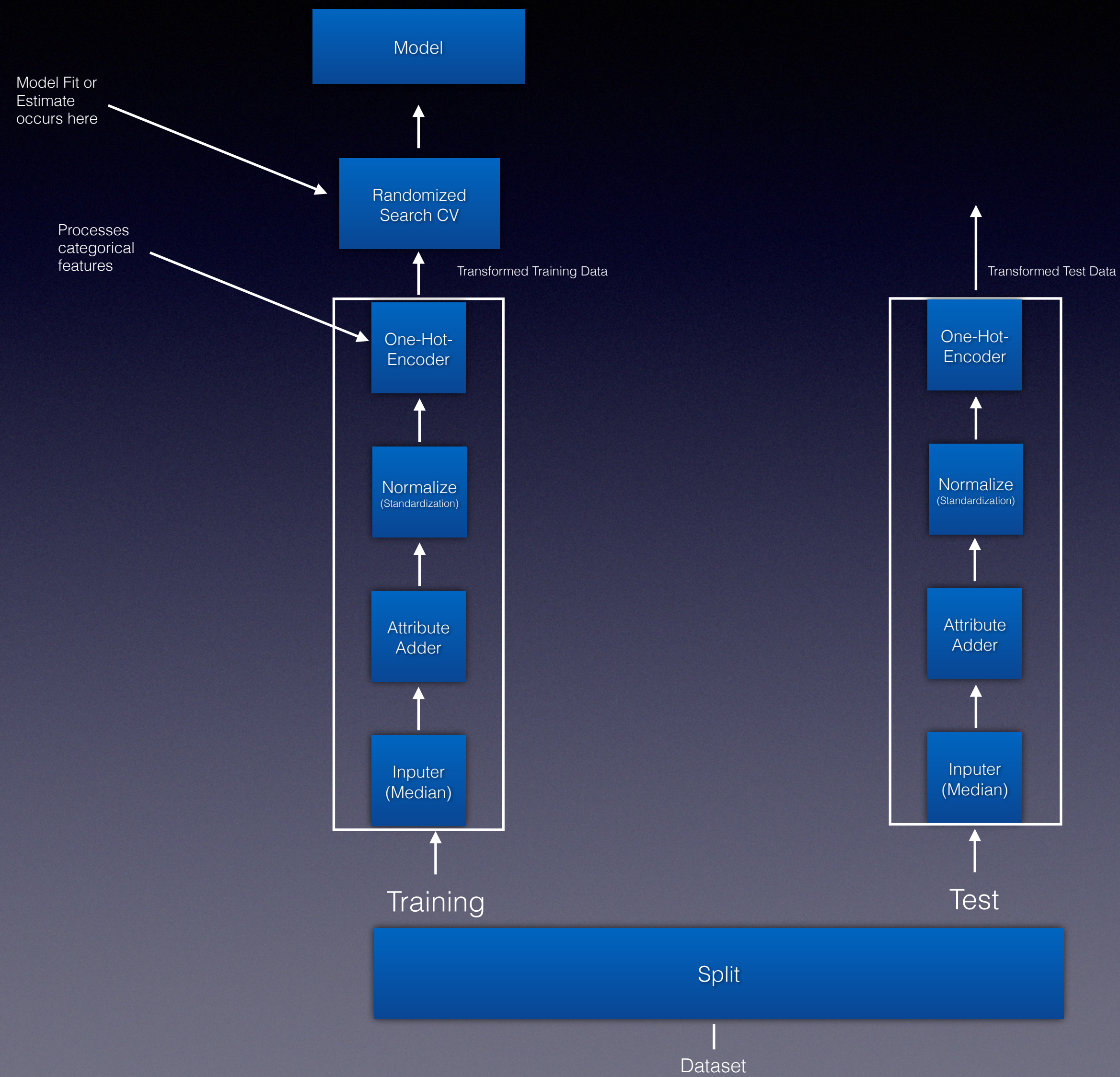
Model Score (RMS)

Training Model Score < Validation Score	Overfitting Model
Training Score = 0	Likely overfitting, run cross validation
Note	Cross validation estimates the performance of model on dataset 'folds', and estimates how precise the model is. Precision being its relation to other validation runs on the model



Estimate generalization error is estimated to lie between A(lower bound) and B(upper bound) with confidence level

Lets researcher know how precise the estimate is.



Test Dataset



Model



Predictions



Statistics



Test Set Labels

Exercise Chapter 2

