# Eps

Machine learning for Ruby

- Build predictive models quickly and easily
- Serve models built in Ruby, Python, R, and more

Check out this post for more info on machine learning with Rails

build passing

## Installation

Add this line to your application's Gemfile:

```
1  gem "eps"
```

On Mac, also install OpenMP:

```
1  brew install libomp
```

## Getting Started

Create a model

```
1  data = [
2    {bedrooms: 1, bathrooms: 1, price: 100000},
3    {bedrooms: 2, bathrooms: 1, price: 125000},
4    {bedrooms: 2, bathrooms: 2, price: 135000},
5    {bedrooms: 3, bathrooms: 2, price: 162000}
6  ]
7  model = Eps::Model.new(data, target: :price)
8  puts model.summary
```

Make a prediction

```
1  model.predict(bedrooms: 2, bathrooms: 1)
```

Store the model

```
1  File.write("model.pmml", model.to_pmml)
```

Load the model

```
1  pmml = File.read("model.pmml")
2  model = Eps::Model.load_pmml(pmml)
```

A few notes:

- The target can be numeric (regression) or categorical (classification)
- Pass an array of hashes to `predict` to make multiple predictions at once
- Models are stored in PMML, a standard for model storage

## Building Models

### Goal

Often, the goal of building a model is to make good predictions on future data. To help achieve this, Eps splits the data into training and validation sets if you have 30+ data points. It uses the training set to build the model and the validation set to evaluate the performance.

If your data has a time associated with it, it's highly recommended to use that field for the split.

```
1  Eps::Model.new(data, target: :price, split: :listed_at)
```

Otherwise, the split is random. There are a number of other options as well.

Performance is reported in the summary.

- For regression, it reports validation RMSE (root mean squared error) - lower is better
- For classification, it reports validation accuracy - higher is better

Typically, the best way to improve performance is feature engineering.

### Feature Engineering

Features are extremely important for model performance. Features can be:

1. numeric
2. categorical
3. text

**Numeric**  For numeric features, use any numeric type.

```
1  {bedrooms: 4, bathrooms: 2.5}
```

**Categorical**    For categorical features, use strings or booleans.

```
1  {state: "CA", basement: true}
```

Convert any ids to strings so they're treated as categorical features.

```
1  {city_id: city_id.to_s}
```

For dates, create features like day of week and month.

```
1  {weekday: sold_on.strftime("%a"), month: sold_on.strftime("%b")}
```

For times, create features like day of week and hour of day.

```
1  {weekday: listed_at.strftime("%a"), hour: listed_at.hour.to_s}
```

**Text**    For text features, use strings with multiple words.

```
1  {description: "a beautiful house on top of a hill"}
```

This creates features based on word count.

You can specify text features explicitly with:

```
1  Eps::Model.new(data, target: :price, text_features: [:description])
```

You can set advanced options with:

```
 1  text_features: {
 2    description: {
 3      min_occurences: 5,          # min times a word must appear to be
            included in the model
 4      max_features: 1000,         # max number of words to include in the
            model
 5      min_length: 1,              # min length of words to be included
 6      case_sensitive: true,       # how to treat words with different
            case
 7      tokenizer: /\s+/,           # how to tokenize the text, defaults to
            whitespace
 8      stop_words: ["and", "the"]  # words to exclude from the model
 9    }
10  }
```

## Full Example

We recommend putting all the model code in a single file. This makes it easy to rebuild the model as needed.

In Rails, we recommend creating a app/ml_models directory. Be sure to restart Spring after creating the directory so files are autoloaded.

```
1  bin/spring stop
```

Here's what a complete model in app/ml_models/price_model.rb may look like:

```ruby
1  class PriceModel < Eps::Base
2    def build
3      houses = House.all
4
5      # train
6      data = houses.map { |v| features(v) }
7      model = Eps::Model.new(data, target: :price, split: :listed_at)
8      puts model.summary
9
10     # save to file
11     File.write(model_file, model.to_pmml)
12
13     # ensure reloads from file
14     @model = nil
15   end
16
17   def predict(house)
18     model.predict(features(house))
19   end
20
21   private
22
23   def features(house)
24     {
25       bedrooms: house.bedrooms,
26       city_id: house.city_id.to_s,
27       month: house.listed_at.strftime("%b"),
28       listed_at: house.listed_at,
29       price: house.price
30     }
31   end
32
33   def model
34     @model ||= Eps::Model.load_pmml(File.read(model_file))
35   end
36
37   def model_file
38     File.join(__dir__, "price_model.pmml")
39   end
40 end
```

Build the model with:

```
1  PriceModel.build
```

This saves the model to `price_model.pmml`. Check this into source control or use a tool like Trove to store it.

Predict with:

```
1  PriceModel.predict(house)
```

## Monitoring

We recommend monitoring how well your models perform over time. To do this, save your predictions to the database. Then, compare them with:

```
1  actual = houses.map(&:price)
2  predicted = houses.map(&:predicted_price)
3  Eps.metrics(actual, predicted)
```

For RMSE and MAE, alert if they rise above a certain threshold. For ME, alert if it moves too far away from 0. For accuracy, alert if it drops below a certain threshold.

## Other Languages

Eps makes it easy to serve models from other languages. You can build models in Python, R, and others and serve them in Ruby without having to worry about how to deploy or run another language.

Eps can serve LightGBM, linear regression, and naive Bayes models. Check out ONNX Runtime and Scoruby to serve other models.

### Python

To create a model in Python, install the sklearn2pmml package

```
1  pip install sklearn2pmml
```

And check out the examples:

- LightGBM Regression
- LightGBM Classification
- Linear Regression
- Naive Bayes

## R

To create a model in R, install the pmml package

```
1  install.packages("pmml")
```

And check out the examples:

- Linear Regression
- Naive Bayes

## Verifying

It's important for features to be implemented consistently when serving models created in other languages. We highly recommend verifying this programmatically. Create a CSV file with ids and predictions from the original model.

| house_id | prediction |
|----------|------------|
| 1        | 145000     |
| 2        | 123000     |
| 3        | 250000     |

Once the model is implemented in Ruby, confirm the predictions match.

```ruby
1   model = Eps::Model.load_pmml("model.pmml")
2
3   # preload houses to prevent n+1
4   houses = House.all.index_by(&:id)
5
6   CSV.foreach("predictions.csv", headers: true, converters: :numeric) do
       |row|
7     house = houses[row["house_id"]]
8     expected = row["prediction"]
9
10    actual = model.predict(bedrooms: house.bedrooms, bathrooms: house.
        bathrooms)
11
12    success = actual.is_a?(String) ? actual == expected : (actual -
        expected).abs < 0.001
13    raise "Bad prediction for house #{house.id} (exp: #{expected}, act:
        #{actual})" unless success
14
15    putc "√"
```

```
16  end
```

## Data

A number of data formats are supported. You can pass the target variable separately.

```
1  x = [{x: 1}, {x: 2}, {x: 3}]
2  y = [1, 2, 3]
3  Eps::Model.new(x, y)
```

Data can be an array of arrays

```
1  x = [[1, 2], [2, 0], [3, 1]]
2  y = [1, 2, 3]
3  Eps::Model.new(x, y)
```

Or Numo arrays

```
1  x = Numo::NArray.cast([[1, 2], [2, 0], [3, 1]])
2  y = Numo::NArray.cast([1, 2, 3])
3  Eps::Model.new(x, y)
```

Or a Rover data frame

```
1  df = Rover.read_csv("houses.csv")
2  Eps::Model.new(df, target: "price")
```

Or a Daru data frame

```
1  df = Daru::DataFrame.from_csv("houses.csv")
2  Eps::Model.new(df, target: "price")
```

When reading CSV files directly, be sure to convert numeric fields. The `table` method does this automatically.

```
1  CSV.table("data.csv").map { |row| row.to_h }
```

## Algorithms

Pass an algorithm with:

```
1  Eps::Model.new(data, algorithm: :linear_regression)
```

Eps supports:

- LightGBM (default)

- Linear Regression
- Naive Bayes

**LightGBM**

Pass the learning rate with:

```
1  Eps::Model.new(data, learning_rate: 0.01)
```

**Linear Regression**

By default, an intercept is included. Disable this with:

```
1  Eps::Model.new(data, intercept: false)
```

To speed up training on large datasets with linear regression, install GSL. With Homebrew, you can use:

```
1  brew install gsl
```

Then, add this line to your application's Gemfile:

```
1  gem "gslr", group: :development
```

It only needs to be available in environments used to build the model.

**Probability**

To get the probability of each category for predictions with classification, use:

```
1  model.predict_probability(data)
```

Naive Bayes is known to produce poor probability estimates, so stick with LightGBM if you need this.

**Validation Options**

Pass your own validation set with:

```
1  Eps::Model.new(data, validation_set: validation_set)
```

Split on a specific value

```
1  Eps::Model.new(data, split: {column: :listed_at, value: Date.parse("
       2019-01-01")})
```

Specify the validation set size (the default is `0.25`, which is 25%)

```
1  Eps::Model.new(data, split: {validation_size: 0.2})
```

Disable the validation set completely with:

```
1  Eps::Model.new(data, split: false)
```

### Database Storage

The database is another place you can store models. It's good if you retrain models automatically.

> We recommend adding monitoring and guardrails as well if you retrain automatically

Create an Active Record model to store the predictive model.

```
1  rails generate model Model key:string:uniq data:text
```

Store the model with:

```
1  store = Model.where(key: "price").first_or_initialize
2  store.update(data: model.to_pmml)
```

Load the model with:

```
1  data = Model.find_by!(key: "price").data
2  model = Eps::Model.load_pmml(data)
```

### Jupyter & IRuby

You can use IRuby to run Eps in Jupyter notebooks. Here's how to get IRuby working with Rails.

### Weights

Specify a weight for each data point

```
1  Eps::Model.new(data, weight: :weight)
```

You can also pass an array

```
1  Eps::Model.new(data, weight: [1, 2, 3])
```

Weights are supported for metrics as well

```
1  Eps.metrics(actual, predicted, weight: weight)
```

Reweighing is one method to mitigate bias in training data

## Upgrading

### 0.3.0

Eps 0.3.0 brings a number of improvements, including support for LightGBM and cross-validation. There are a number of breaking changes to be aware of:

- LightGBM is now the default for new models. On Mac, run:

  ```
  1  brew install libomp
  ```

  Pass the `algorithm` option to use linear regression or naive Bayes.

  ```
  1  Eps::Model.new(data, algorithm: :linear_regression) # or :
        naive_bayes
  ```

- Cross-validation happens automatically by default. You no longer need to create training and test sets manually. If you were splitting on a time, use:

  ```
  1  Eps::Model.new(data, split: {column: :listed_at, value: Date.parse
        ("2019-01-01")})
  ```

  Or randomly, use:

  ```
  1  Eps::Model.new(data, split: {validation_size: 0.3})
  ```

  To continue splitting manually, use:

  ```
  1  Eps::Model.new(data, validation_set: test_set)
  ```

- It's no longer possible to load models in JSON or PFA formats. Retrain models and save them as PMML.

### 0.2.0

Eps 0.2.0 brings a number of improvements, including support for classification.

We recommend:

1. Changing `Eps::Regressor` to `Eps::Model`
2. Converting models from JSON to PMML

```ruby
model = Eps::Model.load_json("model.json")
File.write("model.pmml", model.to_pmml)
```

3. Renaming `app/stats_models` to `app/ml_models`

## History

View the changelog

## Contributing

Everyone is encouraged to help improve this project. Here are a few ways you can help:

- Report bugs
- Fix bugs and submit pull requests
- Write, clarify, or fix documentation
- Suggest or add new features

To get started with development:

```sh
git clone https://github.com/ankane/eps.git
cd eps
bundle install
bundle exec rake test
```