

fn make_expr_strptime

Michael Shoemate

This proof resides in “**contrib**” because it has not completed the vetting process.

Proves soundness of `make_expr_strptime` in `mod.rs` at commit `f5bb719` (outdated¹).

`make_expr_strptime` returns a Transformation that parses strings in a dataframe column into temporal types.

1 Hoare Triple

Precondition

Compiler-verified

- Argument `input_domain` of type `WildExprDomain`
- Argument `input_metric` of type `M`
- Generic `M` implements `OuterMetric`
 - `OuterMetric` defines an associated type `InnerMetric` that must implement `UnboundedMetric`
- `(ExprDomain, M)` implements `MetricSpace`,
- `Expr` implements `StableExpr<M, M>`

Human-verified

None

Pseudocode

```
1 def make_expr_strptime(
2     input_domain: WildExprDomain, input_metric: M, expr: Expr
3 ) -> Transformation:
4     match expr: #
5         case Expr.Function(
6             inputs,
7             function=FunctionExpr.StringExpr(
8                 StringFunction.Strptime(to_type, strftime_options)
9             ),
10        ):
11            pass
12        case _:
13            raise ValueError("expected str.strptime expression")
14    input, ambiguous = inputs
```

¹See new changes with `git diff f5bb719..f99639a rust/src/transformations/make_stable_expr/namespace_str/expr_strptime/mod.rs`

```

15     t_prior = input.make_stable(input_domain, input_metric) #
16     middle_domain, middle_metric = t_prior.output_space()
17
18     if strptime_options.format is None: #
19         raise ValueError("format must be specified")
20
21     if to_type == DataType.Time and not strptime_options.exact:
22         raise ValueError("non-exact not implemented for Time data type")
23
24     # never raise on error
25     strptime_options.strict = False #
26
27     try: #
28         ambig_value = literal_value_of(ambiguous)
29     except Exception:
30         ambig_value = None
31     ambiguous = lit(ambig_value if ambig_value in {"earliest", "latest"} else "null")
32
33     output_domain = middle_domain.clone()
34     series_domain = output_domain.column
35
36     # check input and output types
37     if series_domain.dtype() != DataType.String: #
38         raise ValueError("str.strptime input dtype must be String")
39
40     if to_type not in {DataType.Time, DataType.Datetime, DataType.Date}: #
41         raise ValueError("str.strptime output dtype must be Time, Datetime or Date")
42
43     # in Rust, this assigns to the series domain in output_domain
44     series_domain.set_dtype(to_type) #
45     series_domain.nullable = True
46
47     def function(expr: Expr) -> Expr:
48         return expr.str.strptime(to_type, strptime_options, ambiguous)
49
50     return t_prior >> Transformation.new( #
51         middle_domain,
52         output_domain,
53         Function.then_expr(function),
54         middle_metric,
55         middle_metric,
56         StabilityMap.new(lambda d_in: d_in),
57     )
58

```

Postcondition

Theorem 1.1. For every setting of the input parameters (`input_domain`, `input_metric`, `expr`, `M`) to `make_expr_strptime` such that the given preconditions hold, `make_expr_strptime` raises an error (at compile time or run time) or returns a valid transformation. A valid transformation has the following properties:

1. (Data-independent runtime errors). For every pair of members x and x' in `input_domain`, `invoke(x)` and `invoke(x')` either both return the same error or neither return an error.
2. (Appropriate output domain). For every member x in `input_domain`, `function(x)` is in `output_domain` or raises a data-independent runtime error.
3. (Stability guarantee). For every pair of members x and x' in `input_domain` and for every pair (d_{in}, d_{out}) , where d_{in} has the associated type for `input_metric` and d_{out} has the associated type for `output_metric`, if x, x' are d_{in} -close under `input_metric`, `stability_map(d_{in})` does not raise an

error, and `stability_map(d_in) = d_out`, then `function(x), function(x')` are `d_out`-close under `output_metric`.

2 Proof

Starting from line 4, `expr` is matched as if it were a `strptime` expression, or otherwise rejects the expression.

All preconditions for `make_stable` on line 16 are compiler-verified, therefore by the postcondition `t_prior` is a valid transformation.

To prove that the output is a valid transformation, we must first prove that the transformation on line 51 is a valid transformation.

Proof. Data-Independent Errors Line 38 ensures that the input data is always a string type and line 41 ensures that the output data is always a temporal type, which are necessary preconditions for the `strptime` function in Polars. However, even if these checks were ignored when they would fail, then all possible choices of input dataset would fail, resulting in data-independent errors.

Line 26 disables raising an error when conversion fails. Conversion cannot be ambiguous because line 19 requires the format to be specified, but even if parsing were to be ambiguous, line 28 configures the parsing of ambiguous values to not raise.

Therefore, all errors that may be encountered while parsing result in null values, so no error can be raised by the `strptime` transformation. \square

Proof. Appropriate Output Domain Line 45 constructs a series domain whose elements are of type `to_type` or null. Each string in the input data may either parse into `to_type` or fail and by 26 and 28 evaluate to null.

Therefore the output domain is the same as the input domain, but where the active column domain's data type is nullable values of type `to_type`. \square

Proof. Stability Guarantee `strptime` can be considered a 1-stable row-by-row transformation, as the date parsing is applied independently to each and every row. See `make_row_by_row` for a proof of the stability of row-by-row functions. \square

Since it has been shown that both `t_prior` and the `strptime` transformation are valid transformations, then the preconditions for `make_chain_tt` are met (invoked via the right-shift operator shorthand).