

Package ‘rollama’

March 21, 2024

Title Communicate with 'Ollama'

Version 0.0.3

Description Wraps the 'Ollama' <<https://ollama.com>> API, which can be used to communicate with generative large language models locally.

License GPL (>= 3)

Encoding UTF-8

RoxygenNote 7.3.1

Imports callr, cli, dplyr, httr2, jsonlite, methods, prettyunits, purrr, rlang, tibble

Suggests base64enc, knitr, rmarkdown, spelling, testthat (>= 3.0.0)

VignetteBuilder knitr

Config/testthat/edition 3

Language en-US

URL <https://jbgruber.github.io/rollama/>,
<https://github.com/JBGruber/rollama>

NeedsCompilation no

Author Johannes B. Gruber [aut, cre] (<<https://orcid.org/0000-0001-9177-1772>>),
Maximilian Weber [aut, ctb] (<<https://orcid.org/0000-0002-1174-449X>>)

Maintainer Johannes B. Gruber <JohannesB.Grubergmail.com>

Repository CRAN

Date/Publication 2024-03-21 10:10:02 UTC

R topics documented:

chat_history	2
create_model	2
embed_text	3
list_models	4
ping_ollama	4
pull_model	5
query	6

Index**9**

chat_history	<i>Handle conversations</i>
--------------	-----------------------------

Description

Shows and deletes (`new_chat`) the local prompt and response history to start a new conversation.

Usage

```
chat_history()
```

```
new_chat()
```

Value

`chat_history`: tibble with chat history

`new_chat`: Does not return a value

create_model	<i>Create a model from a Modelfile</i>
--------------	--

Description

Create a model from a Modelfile

Usage

```
create_model(model, modelfile, server = NULL)
```

Arguments

`model` name of the model to create

`modelfile` either a path to a model file to be read or the contents of the model file as a character vector.

`server` URL to an Ollama server (not the API). Defaults to "http://localhost:11434".

Details

Custom models are the way to save your system message and model parameters in a dedicated shareable way. If you use `show_model()`, you can look at the configuration of a model in the column `modelfile`. To get more information and a list of valid parameters, check out <https://github.com/ollama/ollama/blob/main/docs/modelfile.md>. Most options are also available through the `query` and `chat` functions, yet are not persistent over sessions.

Value

Nothing. Called to create a model on the Ollama server.

Examples

```
modelfile <- system.file("extdata", "modelfile.txt", package = "rollama")
## Not run: create_model("mario", modelfile)
modelfile <- "FROM llama2\nSYSTEM You are mario from Super Mario Bros."
## Not run: create_model("mario", modelfile)
```

embed_text	<i>Generate Embeddings</i>
------------	----------------------------

Description

Generate Embeddings

Usage

```
embed_text(
  text,
  model = NULL,
  server = NULL,
  model_params = NULL,
  verbose = getOption("rollama_verbose", default = interactive())
)
```

Arguments

text	text vector to generate embeddings for.
model	which model to use. See https://ollama.com/library for options. Default is "llama2". Set option(rollama_model = "modelname") to change default for the current session. See pull_model for more details.
server	URL to an Ollama server (not the API). Defaults to "http://localhost:11434".
model_params	a named list of additional model parameters listed in the documentation for the Modelfile such as temperature.
verbose	Whether to print status messages to the Console (TRUE/FALSE). The default is to have status messages in interactive sessions. Can be changed with options(rollama_verbose = FALSE).

Value

a tibble with embeddings.

Examples

```
## Not run:
embed_text(c("Here is an article about llamas...",
            "R is a language and environment for statistical computing and graphics."))

## End(Not run)
```

list_models	<i>List models that are available locally.</i>
-------------	--

Description

List models that are available locally.

Usage

```
list_models(server = NULL)
```

Arguments

server URL to an Ollama server (not the API). Defaults to "http://localhost:11434".

Value

a tibble of installed models

ping_ollama	<i>Ping server to see if Ollama is reachable</i>
-------------	--

Description

Ping server to see if Ollama is reachable

Usage

```
ping_ollama(server = NULL, silent = FALSE)
```

Arguments

server URL to an Ollama server (not the API). Defaults to "http://localhost:11434".
 silent suppress warnings and status (only return TRUE/FALSE).

Value

TRUE if server is running

pull_model	<i>Pull, show and delete models</i>
------------	-------------------------------------

Description

Pull, show and delete models

Usage

```
pull_model(model = NULL, server = NULL, insecure = FALSE)
```

```
show_model(model = NULL, server = NULL)
```

```
delete_model(model, server = NULL)
```

```
copy_model(model, destination = paste0(model, "--copy"), server = NULL)
```

Arguments

model	name of the model. Defaults to "llama2" when NULL (except in delete_model).
server	URL to an Ollama server (not the API). Defaults to "http://localhost:11434".
insecure	allow insecure connections to the library. Only use this if you are pulling from your own library during development. description
destination	name of the copied model.

Details

- pull_model(): downloads model
- show_model(): displays information about a local model
- copy_model(): creates a model with another name from an existing model
- delete_model(): deletes local model

Model names: Model names follow a model:tag format, where model can have an optional namespace such as example/model. Some examples are orca-mini:3b-q4_1 and llama2:70b. The tag is optional and, if not provided, will default to latest. The tag is used to identify a specific version.

Value

(invisible) a tibble with information about the model (except in delete_model)

Examples

```
## Not run:
model_info <- pull_model("mixtral")
# after you pull, you can get the same information with:
model_info <- show_model("mixtral")

## End(Not run)
```

query	<i>Chat with a LLM through Ollama</i>
-------	---------------------------------------

Description

Chat with a LLM through Ollama

Usage

```

query(
  q,
  model = NULL,
  screen = TRUE,
  server = NULL,
  images = NULL,
  model_params = NULL,
  template = NULL
)

chat(
  q,
  model = NULL,
  screen = TRUE,
  server = NULL,
  images = NULL,
  model_params = NULL,
  template = NULL
)

```

Arguments

q	the question as a character string or a conversation object.
model	which model(s) to use. See https://ollama.com/library for options. Default is "llama2". Set option(rollama_model = "modelname") to change default for the current session. See pull_model for more details.
screen	Logical. Should the answer be printed to the screen.
server	URL to an Ollama server (not the API). Defaults to "http://localhost:11434".
images	path(s) to images (for multimodal models such as llava).
model_params	a named list of additional model parameters listed in the documentation for the Modelfile such as temperature.
template	the prompt template to use (overrides what is defined in the Modelfile).

Details

query sends a single question to the API, without knowledge about previous questions (only the config message is relevant). chat treats new messages as part of the same conversation until [new_chat](#) is called.

Value

an httr2 response.

Examples

```
## Not run:
# ask a single question
query("why is the sky blue?")

# hold a conversation
chat("why is the sky blue?")
chat("and how do you know that?")

# save the response to an object and extract the answer
resp <- query(q = "why is the sky blue?")
answer <- resp$message$content

# ask question about images (to a multimodal model)
images <- c("https://avatars.githubusercontent.com/u/23524101?v=4", # remote
            "/path/to/your/image.jpg") # or local images supported
query(q = "describe these images",
      model = "llava",
      images = images)

# set custom options for the model at runtime (rather than in create_model())
query("why is the sky blue?",
      model_params = list(
        num_keep = 5,
        seed = 42,
        num_predict = 100,
        top_k = 20,
        top_p = 0.9,
        tfs_z = 0.5,
        typical_p = 0.7,
        repeat_last_n = 33,
        temperature = 0.8,
        repeat_penalty = 1.2,
        presence_penalty = 1.5,
        frequency_penalty = 1.0,
        mirostat = 1,
        mirostat_tau = 0.8,
        mirostat_eta = 0.6,
        penalize_newline = TRUE,
        stop = c("\n", "user:"),
        numa = FALSE,
        num_ctx = 1024,
        num_batch = 2,
        num_gqa = 1,
        num_gpu = 1,
        main_gpu = 0,
        low_vram = FALSE,
        f16_kv = TRUE,
```

```
        vocab_only = FALSE,
        use_mmap = TRUE,
        use_mlock = FALSE,
        embedding_only = FALSE,
        rope_frequency_base = 1.1,
        rope_frequency_scale = 0.8,
        num_thread = 8
    ))

# this might be interesting if you want to turn off the GPU and load the
# model into the system memory (slower, but most people have more RAM than
# VRAM, which might be interesting for larger models)
query("why is the sky blue?",
      model_params = list(num_gpu = 0))

# You can use a custom prompt to override what prompt the model receives
query("why is the sky blue?",
      template = "Just say I'm a llama!")

# Asking the same question to multiple models is also supported
query("why is the sky blue?", model = c("llama2", "orca-mini"))

## End(Not run)
```


Index

chat (query), 6
chat_history, 2
copy_model (pull_model), 5
create_model, 2

delete_model (pull_model), 5

embed_text, 3

list_models, 4

new_chat, 6
new_chat (chat_history), 2

ping_ollama, 4
pull_model, 3, 5, 6

query, 6

show_model (pull_model), 5