

Lessons Learned Using GPT-3 to Build Klaviyo's SMS Assistant

Author: Charlie Natoli

Claps: 80

Date: Jan 24

Over the last year, Klaviyo launched tools that write creative [email subject lines](#) and [SMS messages](#) for our users. The way these work is simple â€” the user provides a description of their message (e.g. â€œweâ€™re having a sale on our scented candle collectionâ€”) and we return multiple marketing copy ideas. Under the hood, we use OpenAI's GPT-3 model.

The image shows two side-by-side screenshots of the Klaviyo SMS assistant interface. The left screenshot shows the input form with the following fields: 'Product name' (filled with 'Natoli Premium Cat Foods') and 'Campaign details' (filled with 'Introducing our classic tuna flavor, now reformulated for senior cats over the age of 7.'). Below the text area is a 'Get ideas' button. The right screenshot shows the output of the assistant, displaying three generated SMS ideas, each with a character count (e.g., '1 SMS', '2 SMS') and a 'Get more ideas' button at the bottom.

Input Form:

- Product name: Natoli Premium Cat Foods
- Campaign details: Introducing our classic tuna flavor, now reformulated for senior cats over the age of 7.
- Get ideas

Generated SMS Ideas:

- Introducing Natoli Premium Classic flavors - now reformulated for senior cats over the age of 7. (1 SMS)
- Natoli has reformulated our classic flavors for senior cats over the age of 7. Learn more here: (1 SMS)
- Introducing our classic flavors, now reformulated to help meet the nutritional needs of senior cats over the age of 7. Made with a limited amount of ingredients to help prevent allergies and stomach issues. Learn more here: (2 SMS)
- Get more ideas

Putting our SMS assistant to use for my *totally real* cat food business.

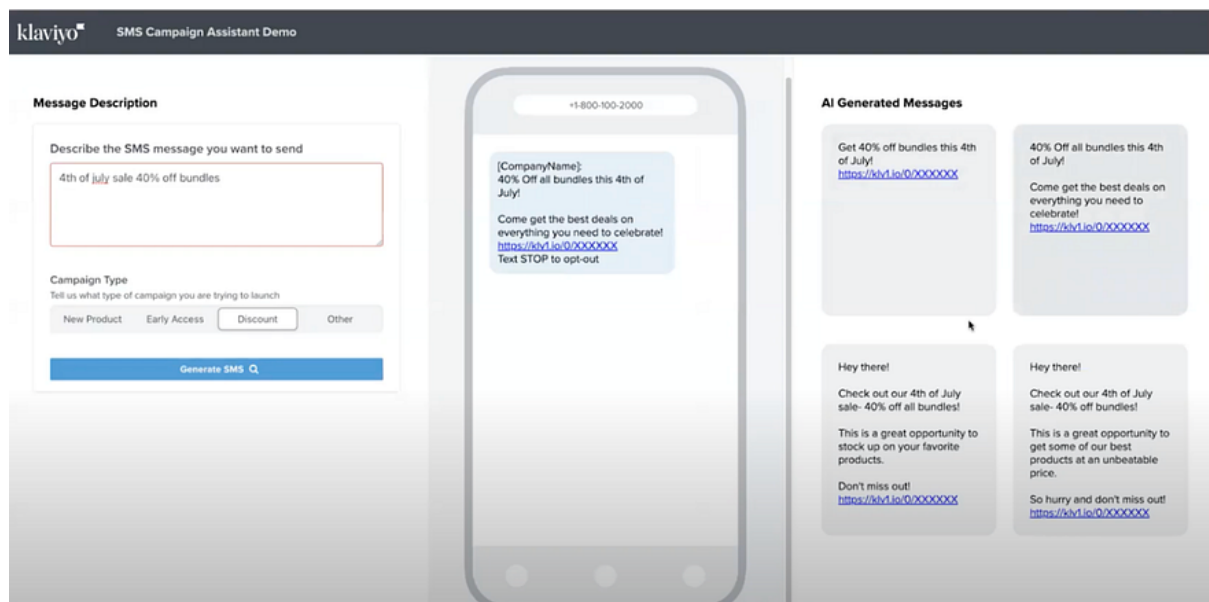
I was the data science lead on our SMS assistant. The project was different from other data science problems I've tackled in the past. First, working with generative AI models (whether generating text or images) is different than using models for other data science applications, such as classification or prediction. Second, while customer understanding is always key for a data scientist, I felt it was especially important here due to the subjective and open-ended nature of copywriting.

Here are lessons I learned.

When building a product using generative AI, actively participate in user research.

In many data science projects, it's often somewhat clearer what the goal of a model might be, for example classifying fraud, recommending products, or forecasting churn. However, for something as open ended and subjective as creative text generation, it was less clear what exactly the model should be doing. For example, are there certain marketing best practices to follow? Do users want you to closely follow their prompt? Or do they want you to show them new, creative ideas?

To get a sense for this, it was critical to work with the product manager and designer on the project in interviewing customers. We started by asking customers about their creative process and what makes it hard. However, just asking in the abstract often wasn't enough; the best way to get specific and actionable feedback was to show users real examples and let them play around with a prototype. To do this, we built a lightweight app that allowed users to play with the model, and used it in our conversations.



A screenshot of the tool we built to demo the SMS Assistant in interviews. It proved to be a great conversation guide.

From these interviews, we were able to get at more granular requirements for how the model should behave:

- **Be creative.** Many users, particularly small business owners wearing multiple hats, were constantly struggling to come up with new and creative ways of introducing things like product releases or sales. This suggested to us that taking creative leaps from the user's prompt was a key goal.
- **Avoid made up details.** Our model sometimes made up details that weren't true (such as discount codes that the user didn't specify in their prompt). While these seem like realistic details that might go in a promotional text message, they are unhelpful and (understandably) look weird to users. (There seems to be increasing public awareness of this problem, sometimes referred to as "hallucination," especially given the examples of this happening with ChatGPT.)

- **Limit message length:** Users also cared a lot about having a mix of longer and shorter options. On the one hand, longer SMS messages cost more to send. On the other, it can be hard to convey your full point in a creative way within a tight character limit.

Performance of your generative AI model is hard to quantify before launching.

Before launching, quantitative metrics are tempting but often not worth it. As data scientists, we're always thinking about how to quantify things like model performance. In this project, I was tempted to do this here as well. Many academic papers test models against existing NLP benchmarks, and we wondered if we could do something similar. Doing so would help more formally compare different model versions, and be more reproducible than saying "yeah, we looked at some examples and we think it's good."

However, was there a benchmark that closely captured what we were trying to do? The answer for us was "not really." As described above, just figuring out what we wanted to optimize for was hard. Many existing benchmarks measure something much more narrow like sentiment, formality, or politeness, whereas our goal was a lot harder to define. In the end, meticulously looking through examples and validating that they were good enough based on customer learnings gave us confidence to move forward.

Test things live, running experiments where you can "there's no better way to learn than that."

Even with deep customer research and testing, it's still hard to say exactly how well your model will fare with customers, for a few reasons:

- Users may react differently in reality compared to how they react in interviews. For example, customer interviewees may shy away from giving you negative feedback out of politeness, or they can't tell how useful the tool is until they're actually sitting down to write a promotional SMS.
- There may be important edge cases (in our case, types of prompts people give to the SMS Assistant) that you didn't think of going in. Your model may not be robust to all of these.

In some cases (for example, a medical chatbot), there may be severe consequences to a bad message being generated. However, for a creative/marketing use case, this wasn't a problem as long as there were enough good answers that the user got value from the tool. So our aim became to get the model into the hands of real users. But, to test and tweak how our model fared, we ran a longer than normal beta period. We also considered (and continue to look at):

- **Message selection rate.** Of everyone who used the tool, what percent clicked on at least one of the messages generated? What factors were correlated with higher selection rates?
- **Use retention rate.** What percent of users came back and used the tool again a second time? What factors were correlated with higher retention?
- **Message quality.** We also looked at a large sample of both prompts and messages, to get a sense for what types of things people do and don't put in. Here, we saw a much wider

range of inputs than we expected, from some users putting in only a word or two, to others pasting in large amounts of text from other marketing copy.

- **Message filtering rate.** While we only need to return three messages to the user, we generate more than we need, and then filter some out that don't meet various criteria we have set. For example, messages that are too long, contain toxic content, or contain made up discount codes all get filtered and not shown to the user. Sometimes, we aren't able to generate three good ones. To ensure long-term success, we studied if certain types of prompts were more likely to struggle generating enough acceptable responses.

Wrapping up

Overall, working with generative AI has been a fantastic challenge, and a very different type of challenge than traditional data science work. These technologies will likely get integrated into lots of types of software. This seems even more clear after the recent release of ChatGPT which has raised awareness of generative AI and is changing user expectations for what's possible.

For all data scientists interested, I encourage you to read up on how the models work, and what training data was used. Read how others online have tested GPT-3 (and increasingly, ChatGPT). Push your teams to consider both the range of possible opportunities, as well as risks, in using these models. And most importantly, I encourage you to play around with the models yourself.

Links

NLP benchmarks: [This paper](#) gives a good review of some of these benchmarks, as well as their pros and cons compared to human evaluation.

Exploring what GPT-3 can do:

- See our recent [blog post](#) exploring the range of GPT-3's cousin, ChatGPT can do.
- Another great and really thorough account I've found of what GPT-3 can do are blog posts from Gwern Branwen on [nonfiction](#) and [fiction](#).

Possible harms from large language models: I found [this paper](#) to give a really helpful overview of potential future harms from large language models overall. Some of these apply directly to our case of generating marketing copy.