

A Model of ChatGPT

Tom Cunningham

2025-06-14

This paper develops a simple model of human and AI ability to answer questions. Each question \mathbf{q} is a high-dimensional vector, with a true scalar answer a . An agent's estimate of the answer is an interpolation based on previously-seen questions and answers $(\mathbf{q}^i, a^i)_{i=1, \dots, n}$. This framework extends an [earlier model](#) developed for a different purpose.

The model yields several implications:

1. **The quality of an answer to a new question depends on its distance from the training set.** For a new question \mathbf{q} , the expected error is a function of the distance between \mathbf{q} and the training set \mathbf{Q} .
2. **The quality of answers increases with the size of the training set.** The expected error decreases linearly with the number of linearly-independent examples in the training set.
3. **The value of advice from another agent depends on the distance between their training sets.**

This framework can be interpreted as a model of an agent, "the user," who must provide an estimate for the answer to a question \mathbf{q} and can choose whether to consult an AI model like ChatGPT. The key components of the model are:

1. **The dimensionality of the question (p).** A higher-dimensional problem may be more costly to enter into the AI, but it also increases the potential benefit.
2. **The public information set.** These are the training questions that the AI has observed, which we can conceptualize as the corpus of public knowledge (e.g., the internet).
3. **The private information set.** These are the questions that the user has personally encountered and for which they have observed the true answer.

A user will consult the AI if and only if the expected improvement in their answer exceeds the associated cost. The model predicts that an AI will be most useful for questions with components that are novel to the user but contained within the AI's public training data.

This leads to several corollaries:

1. An AI will not be used for questions the user has encountered before.
2. An AI is more likely to be used for domains with higher *latent* dimensionality (p).
3. An AI is more likely to be used for domains with lower *surface* dimensionality, as this reduces the cost of specifying the question.
4. An AI is more likely to be used by humans with less experience in a domain (i.e., smaller n_{private}).

We can make some conjectures about adoption by occupation:

We can make some conjectures about adoption by task:

Additional things to add:

1. **High-dimensional answers.** Our model assumes *scalar* answers. In fact ChatGPT gives high-dimensional outputs. I think we can say some nice things here.
2. **Tacit knowledge.** ChatGPT will be more likely to be used for domains where humans have tacit knowledge.

Table 1: Conjectures about adoption by occupation

Occupation	Predicted ChatGPT use	Reason
software engineer	high	many novel discrete problems, similar to those on
software engineer - idiosyncratic language	low	many novel discrete problems, not similar to those on
physician	high	many novel discrete problems, similar to those on
contact center worker	low	novel problems, but not similar to those on the i
architect	low	novel problems, not discrete, not text-based
manual worker	low	not not text-based

Table 2: Conjectures about adoption by task

Task	Predicted ChatGPT use	Reason
Intellectual curiosity	high	novel discrete problem, similar to those on
Diagnosing medical problems	high	novel discrete problem, similar to those on
Problems with widely-adopted systems (car, house, computer)	high	novel discrete problem, similar to those on
Problems with idiosyncratic systems (custom setups)	low	novel discrete problem, <i>not</i> similar to those on

1 Model

The State of the World and Questions. The state of the world is defined by a vector of p unobserved parameters, $\mathbf{w} \in \mathbb{R}^p$. A question is a vector of p binary features, $\mathbf{q} \in \{-1, 1\}^p$. The true answer to a question \mathbf{q} is a scalar a determined by the linear relationship:

$$a = \mathbf{q}'\mathbf{w} = \sum_{k=1}^p q_k w_k$$

Agents and Information. There is a set of agents, indexed by $i \in \mathcal{I}$. Each agent i possesses an information set \mathcal{D}_i , which consists of n_i questions they have previously encountered, along with their true answers. We can represent this information as a pair $(\mathbf{Q}_i, \mathbf{a}_i)$:

- \mathbf{Q}_i is an $n_i \times p$ matrix where each row is a question vector. Let the j -th question for agent i be $\mathbf{q}'_{i,j}$, so that:

$$\mathbf{Q}_i = \begin{bmatrix} \mathbf{q}'_{i,1} \\ \vdots \\ \mathbf{q}'_{i,n_i} \end{bmatrix} = \begin{bmatrix} q_{i,1,1} & \cdots & q_{i,1,p} \\ \vdots & \ddots & \vdots \\ q_{i,n_i,1} & \cdots & q_{i,n_i,p} \end{bmatrix}$$

- \mathbf{a}_i is an $n_i \times 1$ vector of the corresponding answers. The answers are generated according to the true model:

$$\mathbf{a}_i = \mathbf{Q}_i \mathbf{w}$$

Beliefs. All agents share a common prior belief about the state of the world, assuming the weights \mathbf{w} are drawn from a multivariate Gaussian distribution:

$$\mathbf{w} \sim N(\mathbf{0}, \Sigma)$$

where Σ is a $p \times p$ positive-semidefinite covariance matrix. A common assumption we will use is an isotropic prior, where $\Sigma = \sigma^2 \mathbf{I}_p$ for some scalar $\sigma^2 > 0$. This implies that, a priori, the weights are uncorrelated and have equal variance.

Given their information set \mathcal{D}_i , agent i forms a posterior belief about \mathbf{w} . When a new question \mathbf{q}_{new} arises, the agent uses their posterior distribution to form an estimate of the answer, $\hat{a}_{\text{new}} = \mathbf{q}'_{\text{new}} \mathbb{E}[\mathbf{w} \mid \mathcal{D}_i]$.

2 Propositions

Proposition 1 (Posterior over \mathbf{w} given \mathbf{Q} and \mathbf{a}). *The agent's posterior mean and variance will be:*

$$\begin{aligned}\hat{\mathbf{w}} &= \Sigma \mathbf{Q}^\top (\mathbf{Q} \Sigma \mathbf{Q}^\top)^{-1} \mathbf{a} \\ \Sigma_{|\mathbf{a}} &= \Sigma - \Sigma \mathbf{Q}^\top (\mathbf{Q} \Sigma \mathbf{Q}^\top)^{-1} \mathbf{Q} \Sigma.\end{aligned}$$

Proof. The derivation follows from the standard formula for conditional Gaussian distributions. We begin by defining the joint distribution of the weights \mathbf{w} and the answers \mathbf{a} . The weights and answers are jointly Gaussian:

$$\begin{pmatrix} \mathbf{w} \\ \mathbf{a} \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \Sigma & \Sigma \mathbf{Q}' \\ \mathbf{Q} \Sigma & \mathbf{Q} \Sigma \mathbf{Q}' \end{pmatrix} \right)$$

where the covariance terms are derived as follows:

- $Cov(\mathbf{w}, \mathbf{w}) = \Sigma$ (prior covariance)
- $Cov(\mathbf{a}, \mathbf{a}) = Cov(\mathbf{Q}\mathbf{w}, \mathbf{Q}\mathbf{w}) = \mathbf{Q} Cov(\mathbf{w}, \mathbf{w}) \mathbf{Q}' = \mathbf{Q} \Sigma \mathbf{Q}'$
- $Cov(\mathbf{w}, \mathbf{a}) = Cov(\mathbf{w}, \mathbf{Q}\mathbf{w}) = Cov(\mathbf{w}, \mathbf{w}) \mathbf{Q}' = \Sigma \mathbf{Q}'$

The conditional mean $E[\mathbf{w}|\mathbf{a}]$ is given by the formula:

$$E[\mathbf{w}|\mathbf{a}] = E[\mathbf{w}] + Cov(\mathbf{w}, \mathbf{a}) Var(\mathbf{a})^{-1} (\mathbf{a} - E[\mathbf{a}])$$

Substituting the values from our model ($E[\mathbf{w}] = \mathbf{0}$, $E[\mathbf{a}] = \mathbf{0}$):

$$\hat{\mathbf{w}} = \mathbf{0} + (\Sigma \mathbf{Q}') (\mathbf{Q} \Sigma \mathbf{Q}')^{-1} (\mathbf{a} - \mathbf{0}) = \Sigma \mathbf{Q}' (\mathbf{Q} \Sigma \mathbf{Q}')^{-1} \mathbf{a}$$

This gives us the posterior mean of the weights. The posterior covariance is given by:

$$Var(\mathbf{w}|\mathbf{a}) = Var(\mathbf{w}) - Cov(\mathbf{w}, \mathbf{a}) Var(\mathbf{a})^{-1} Cov(\mathbf{a}, \mathbf{w}) = \Sigma - \Sigma \mathbf{Q}' (\mathbf{Q} \Sigma \mathbf{Q}')^{-1} \mathbf{Q} \Sigma.$$

□

Proposition 2 (Expected error for a given question). *The expected squared error for a new question \mathbf{q} is:*

$$\mathbb{E}[(\mathbf{q}'(\mathbf{w} - \hat{\mathbf{w}}))^2] = \mathbf{q}' \Sigma_{|\mathbf{a}} \mathbf{q}$$

For an isotropic prior where $\Sigma = \sigma^2 \mathbf{I}$, the error is proportional to the squared distance of \mathbf{q} from the subspace spanned by the previously seen questions in \mathbf{Q} :

$$\mathbb{E}[(\mathbf{q}'(\mathbf{w} - \hat{\mathbf{w}}))^2] = \sigma^2 \|\mathbf{I} - \mathbf{P}_{\mathbf{Q}}\mathbf{q}\|^2$$

where $\mathbf{P}_{\mathbf{Q}}$ is the projection matrix onto the row-span of \mathbf{Q} .

Proof. The prediction error is $\mathbf{q}'\mathbf{w} - \mathbf{q}'\hat{\mathbf{w}} = \mathbf{q}'(\mathbf{w} - \hat{\mathbf{w}})$. The expected squared error is the variance of this prediction error.

$$\begin{aligned}\mathbb{E}[(\mathbf{q}'(\mathbf{w} - \hat{\mathbf{w}}))^2] &= \mathbb{E}[\mathbf{q}'(\mathbf{w} - \hat{\mathbf{w}})(\mathbf{w} - \hat{\mathbf{w}})' \mathbf{q}] \\ &= \mathbf{q}' \mathbb{E}[(\mathbf{w} - \hat{\mathbf{w}})(\mathbf{w} - \hat{\mathbf{w}})'] \mathbf{q} \\ &= \mathbf{q}' Var(\mathbf{w} | \mathbf{a}) \mathbf{q} = \mathbf{q}' \Sigma_{|\mathbf{a}} \mathbf{q}\end{aligned}$$

This proves the first part of the proposition. For the second part, we assume an isotropic prior $\Sigma = \sigma^2 \mathbf{I}$. Substituting this into the expression for $\Sigma_{|\mathbf{a}}$ from Proposition 1:

$$\begin{aligned}\Sigma_{|\mathbf{a}} &= \sigma^2 \mathbf{I} - (\sigma^2 \mathbf{I}) \mathbf{Q}' (\mathbf{Q} (\sigma^2 \mathbf{I}) \mathbf{Q}')^{-1} \mathbf{Q} (\sigma^2 \mathbf{I}) \\ &= \sigma^2 \mathbf{I} - \sigma^4 \mathbf{Q}' (\sigma^2 \mathbf{Q} \mathbf{Q}')^{-1} \mathbf{Q} \\ &= \sigma^2 \mathbf{I} - \sigma^4 (\sigma^2)^{-1} \mathbf{Q}' (\mathbf{Q} \mathbf{Q}')^{-1} \mathbf{Q} \\ &= \sigma^2 (\mathbf{I} - \mathbf{Q}' (\mathbf{Q} \mathbf{Q}')^{-1} \mathbf{Q})\end{aligned}$$

Let $P_Q = Q'(QQ')^{-1}Q$, which is the projection matrix onto the row space of Q . Then $\Sigma|_a = \sigma^2(I - P_Q)$. The expected squared error is:

$$\mathbb{E}[(q'(\mathbf{w} - \hat{\mathbf{w}}))^2] = q'\sigma^2(I - P_Q)q = \sigma^2 q'(I - P_Q)q$$

Since $I - P_Q$ is an idempotent projection matrix, $q'(I - P_Q)q = q'(I - P_Q)'(I - P_Q)q = \|(I - P_Q)q\|^2$. Thus,

$$\mathbb{E}[(q'(\mathbf{w} - \hat{\mathbf{w}}))^2] = \sigma^2 \|(I - P_Q)q\|^2$$

□

Proposition 3 (Error decreases with more independent questions). *The average expected squared error over all possible new questions q decreases linearly with the number of linearly independent questions in the training set Q . Specifically, with an isotropic prior $\Sigma = \sigma^2 I$, the average error is:*

$$\mathbb{E}_q[\text{error}(q)] = \sigma^2(p - \text{rank}(Q))$$

where the expectation is taken over new questions q with i.i.d. components drawn uniformly from $\{-1, 1\}$.

Proof. The proof proceeds in two steps. First, we write the expression for the error for a given new question q . Second, we average this error over the distribution of all possible questions.

1. **Predictive error for a fixed q .** From Proposition 2, the expected squared error for a specific new question q , given an isotropic prior $\Sigma = \sigma^2 I$, is:

$$\text{error}(q) = \mathbb{E}[(q'(\mathbf{w} - \hat{\mathbf{w}}))^2] = \sigma^2 q'(I - P_Q)q$$

where $P_Q = Q'(QQ')^{-1}Q$ is the projection matrix onto the row-span of Q .

2. **Average over random new questions.** We now take the expectation of this error over the distribution of new questions q . The components of q are i.i.d. uniform on $\{-1, 1\}$, which implies that $\mathbb{E}[q] = \mathbf{0}$ and $\mathbb{E}[qq'] = I_p$. The average error is:

$$\begin{aligned} \mathbb{E}_q[\text{error}(q)] &= \mathbb{E}_q[\sigma^2 q'(I - P_Q)q] \\ &= \sigma^2 \mathbb{E}_q[\text{tr}(q'(I - P_Q)q)] \\ &= \sigma^2 \mathbb{E}_q[\text{tr}((I - P_Q)qq')] \\ &= \sigma^2 \text{tr}((I - P_Q)\mathbb{E}_q[qq']) \\ &= \sigma^2 \text{tr}(I - P_Q) \\ &= \sigma^2(\text{tr}(I) - \text{tr}(P_Q)) \end{aligned}$$

The trace of the identity matrix is p . The trace of a projection matrix is the dimension of the subspace it projects onto, so $\text{tr}(P_Q) = \text{rank}(Q)$. Thus, the average error is:

$$\mathbb{E}_q[\text{error}(q)] = \sigma^2(p - \text{rank}(Q))$$

Since the rank of Q increases with each linearly independent question added, the average error decreases linearly until $\text{rank}(Q) = p$, at which point it becomes zero.

□

Proposition 4 (Posterior in two-stage estimation). *We consider a two-stage process. First, an agent (the "computer," C) with training data (Q_C, \mathbf{a}_C) forms an estimate for the answer to a new question q . Second, another agent (the "human," H) with their own training data (Q_H, \mathbf{a}_H) observes the computer's estimate and updates their own belief.*

The human has a prior over the weights $\mathbf{w} \sim N(\mathbf{0}, \Sigma)$. After observing their own data, the human's posterior for \mathbf{w} is $N(\hat{\mathbf{w}}_H, \Sigma_H)$, where from Proposition 1:

$$\begin{aligned} \hat{\mathbf{w}}_H &= \Sigma Q_H^\top (Q_H \Sigma Q_H^\top)^{-1} \mathbf{a}_H \\ \Sigma_H &= \Sigma - \Sigma Q_H^\top (Q_H \Sigma Q_H^\top)^{-1} Q_H \Sigma \end{aligned}$$

The human's initial estimate for the answer to a new question \mathbf{q} is $\mu_H = \mathbf{q}'\hat{\mathbf{w}}_H$ with variance $\sigma_H^2 = \mathbf{q}'\Sigma_H\mathbf{q}$.

The computer has its own training data $(\mathbf{Q}_C, \mathbf{a}_C)$. It provides an estimate $\hat{a}_C = \mathbf{q}'\hat{\mathbf{w}}_C$ for the true answer $a = \mathbf{q}'\mathbf{w}$. The human observes \hat{a}_C and updates their posterior for a . We assume the computer's observations may be noisy, such that $\mathbf{a}_C = \mathbf{Q}_C\mathbf{w} + \boldsymbol{\epsilon}_C$ with $\boldsymbol{\epsilon}_C \sim N(0, s_C^2\mathbf{I})$.

We analyze the human's final posterior for a under different assumptions about what the human knows about the computer's process.

Proposition 4.1 (Updating with minimal information). Assume the human has no knowledge of the computer's training set \mathbf{Q}_C but believes the computer's estimate is unbiased with a known mean squared error τ^2 . That is, $\hat{a}_C = a + \eta$, where $\eta \sim N(0, \tau^2)$ and is independent of \mathbf{w} .

Upon observing \hat{a}_C , the human's posterior for a is:

$$a \mid \hat{a}_C \sim N(\mu_H + \alpha(\hat{a}_C - \mu_H), (1 - \alpha)\sigma_H^2)$$

where $\alpha = \frac{\sigma_H^2}{\sigma_H^2 + \tau^2} \in [0, 1]$. The human's new estimate is a weighted average of their own initial estimate and the computer's estimate. The weight α placed on the computer's estimate is higher when the computer is believed to be more accurate (smaller τ^2) or when the human's own estimate is more uncertain (larger σ_H^2).

Proposition 4.2 (Knowledge of computer's questions). Assume the human knows the computer's training questions \mathbf{Q}_C and its noise level s_C^2 , but not the observed answers \mathbf{a}_C .

The human can model the computer's estimate as $\hat{a}_C = \mathbf{q}'\mathbf{P}\mathbf{w} + \mathbf{q}'\boldsymbol{\zeta}$, where $\mathbf{P} = \Sigma\mathbf{Q}'_C(\mathbf{Q}_C\Sigma\mathbf{Q}'_C + s_C^2\mathbf{I})^{-1}\mathbf{Q}_C$ and $\boldsymbol{\zeta} = \Sigma\mathbf{Q}'_C(\mathbf{Q}_C\Sigma\mathbf{Q}'_C + s_C^2\mathbf{I})^{-1}\boldsymbol{\epsilon}_C$.

The pair (a, \hat{a}_C) is jointly Gaussian, conditional on the human's data. The posterior for a is:

$$a \mid \hat{a}_C \sim N(\mu_H + \kappa(\hat{a}_C - \mu_C), \sigma_H^2 - \kappa\sigma_{HC})$$

where:

- $\mu_C = \mathbf{q}'\mathbf{P}\hat{\mathbf{w}}_H$ (human's expectation of computer's estimate)
- $\sigma_{HC} = \mathbf{q}'\Sigma_H\mathbf{P}'\mathbf{q}$ (covariance)
- $\sigma_C^2 = \mathbf{q}'\mathbf{P}\Sigma_H\mathbf{P}'\mathbf{q} + \mathbf{q}'\Sigma_\zeta\mathbf{q}$ (variance of computer's estimate)
- $\Sigma_\zeta = s_C^2\mathbf{P}\Sigma^{-1}\mathbf{P}'$
- $\kappa = \frac{\sigma_{HC}}{\sigma_C^2}$ (the weight on the computer's prediction error)

The weight κ depends on the covariance structure, which is influenced by the overlap between the subspaces spanned by \mathbf{Q}_H and \mathbf{Q}_C .

Proposition 4.3 (Limiting cases). The framework of Proposition 4.2 nests two extreme cases:

1. **Oracle Trust:** If the human believes the computer's estimate is perfect (e.g., $s_C^2 \rightarrow 0$ and \mathbf{Q}_C spans the relevant subspace), then $\kappa \rightarrow \sigma_H^2/(\mathbf{q}'\mathbf{P}\Sigma_H\mathbf{P}'\mathbf{q})$, and the posterior variance collapses towards zero. In the simplified Kalman model, if $\tau^2 \rightarrow 0$, then $\alpha \rightarrow 1$, and the human adopts the computer's answer, $a \mid \hat{a}_C \rightarrow N(\hat{a}_C, 0)$.
2. **Total Skepticism:** If the human believes the computer provides no information (e.g., $\sigma_{HC} \rightarrow 0$ because \mathbf{Q}_C is irrelevant to \mathbf{q}), then $\kappa \rightarrow 0$. In the Kalman model, if $\tau^2 \rightarrow \infty$, then $\alpha \rightarrow 0$. In both cases, the human ignores the computer's estimate and reverts to their original posterior, $a \mid \hat{a}_C \sim N(\mu_H, \sigma_H^2)$.

3 Related Literature

3.1 Agrawal et al. (2018) "Exploring the Impact of Artificial Intelligence: Prediction versus Judgment"

https://www.nber.org/system/files/working_papers/w24626/w24626.pdf

3.2 Kleinberg et al. (2017) "Human Decisions and Machine Predictions"

4 References
