

A Model of ChatGPT

Tom Cunningham

2025-06-14

Abstract

We present a simple Bayesian model of how large language models such as ChatGPT and humans answer factual questions. Each question is represented as a high-dimensional binary vector and its true answer is a linear form in an unknown weight vector. Agents interpolate from previously-answered questions drawn from private (human) or public (AI) data. We derive closed-form expressions for posterior beliefs, prediction error, the value of additional examples, and the incremental value of consulting ChatGPT before making a decision. The framework yields testable implications about which users will use ChatGPT for which tasks, specifically ChatGPT will be used when encountering a question logically distinct from the questions in a user's own prior experience, and that novel component is represented in the space of questions in the public dataset (e.g on the internet).

I first describe a simple model of an agent answering a question, so we can formalize the sense in which the answer will be better when the question is more similar to questions that the agent already knows the answer to. We can then also talk about what types of questions it will be worthwhile asking for advice from another agent (ChatGPT). This gives us a simple model of which types of questions will be worthwhile asking ChatGPT demand, based on the nature of the question, the prior experience of the agent, and the prior experience of ChatGPT (AKA its training data).

Formally, each \mathbf{q} is a vector of binary attributes, with a true scalar answer a , that is linear in the attributes. An agent's estimate of the answer will be an interpolation based on previously-seen questions and answers $(\mathbf{q}^i, a^i)_{i=1, \dots, n}$. This framework extends an [earlier model](#) I developed for a different purpose.

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 ChatGPT. The key statistics of the setting will determine the incremental value of consulting ChatGPT:

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 and by task:

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

Table 1: Conjectures about adoption by occupation

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.

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:

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

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

Table 2: Conjectures about adoption by task

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.

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

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 $\mathbf{P}_{\mathbf{Q}} = \mathbf{Q}'(\mathbf{Q}\mathbf{Q}')^{-1}\mathbf{Q}$, which is the projection matrix onto the row space of \mathbf{Q} . Then $\Sigma_{|\mathbf{a}} = \sigma^2(\mathbf{I} - \mathbf{P}_{\mathbf{Q}})$. The expected squared error is:

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

Since $\mathbf{I} - \mathbf{P}_{\mathbf{Q}}$ is an idempotent projection matrix, $\mathbf{q}'(\mathbf{I} - \mathbf{P}_{\mathbf{Q}})\mathbf{q} = \mathbf{q}'(\mathbf{I} - \mathbf{P}_{\mathbf{Q}})'(\mathbf{I} - \mathbf{P}_{\mathbf{Q}})\mathbf{q} = \|(\mathbf{I} - \mathbf{P}_{\mathbf{Q}})\mathbf{q}\|^2$. Thus,

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

□

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

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

where the expectation is taken over new questions \mathbf{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 \mathbf{q} . Second, we average this error over the distribution of all possible questions.

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

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

where $\mathbf{P}_Q = \mathbf{Q}'(\mathbf{Q}\mathbf{Q}')^{-1}\mathbf{Q}$ is the projection matrix onto the row-span of \mathbf{Q} .

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

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

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

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

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

□

Proposition 4 (Two-stage updating with agents 1 and 2). *Consider two agents who share an isotropic prior $\mathbf{w} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_p)$.*

- *Agent 1 observes data $(\mathbf{Q}_1, \mathbf{a}_1)$ and forms the posterior mean*

$$\hat{\mathbf{w}}_1 = \mathbf{Q}_1^\top (\mathbf{Q}_1 \mathbf{Q}_1^\top)^{-1} \mathbf{a}_1, \quad \mathbf{P}_1 := \mathbf{Q}_1^\top (\mathbf{Q}_1 \mathbf{Q}_1^\top)^{-1} \mathbf{Q}_1.$$

- *Agent 2 observes data $(\mathbf{Q}_2, \mathbf{a}_2)$ and forms the posterior mean*

$$\hat{\mathbf{w}}_2 = \mathbf{Q}_2^\top (\mathbf{Q}_2 \mathbf{Q}_2^\top)^{-1} \mathbf{a}_2, \quad \mathbf{P}_2 := \mathbf{Q}_2^\top (\mathbf{Q}_2 \mathbf{Q}_2^\top)^{-1} \mathbf{Q}_2.$$

- *A new question $\mathbf{q} \in \{-1, 1\}^p$ arrives. Agent 1 announces the estimate $\hat{a}_1 = \mathbf{q}'\hat{\mathbf{w}}_1$. Let*

$$\mu_2 = \mathbf{q}'\hat{\mathbf{w}}_2, \quad \sigma_2^2 = \sigma^2 \mathbf{q}'(\mathbf{I} - \mathbf{P}_2)\mathbf{q}, \quad (1)$$

$$\mu_{2,1} = \mathbf{q}'\mathbf{P}_1\hat{\mathbf{w}}_2, \quad \sigma_{21} = \sigma^2 \mathbf{q}'(\mathbf{I} - \mathbf{P}_2)\mathbf{P}_1^\top \mathbf{q}, \quad (2)$$

$$\sigma_{1|2}^2 = \sigma^2 \mathbf{q}'\mathbf{P}_1(\mathbf{I} - \mathbf{P}_2)\mathbf{P}_1^\top \mathbf{q}, \quad \kappa = \frac{\sigma_{21}}{\sigma_{1|2}^2}. \quad (3)$$

Then, the posterior distribution of the true answer $a = \mathbf{q}'\mathbf{w}$ for agent 2 before seeing \hat{a}_1 is $N(\mu_2, \sigma_2^2)$, and after observing \hat{a}_1 it is

$$a \mid \hat{a}_1, \mathbf{a}_2 \sim N(\mu_{2|1}, \sigma_{2|1}^2), \quad \mu_{2|1} = \mu_2 + \kappa(\hat{a}_1 - \mu_{2,1}), \quad \sigma_{2|1}^2 = \sigma_2^2 - \kappa \sigma_{21}.$$

Proof. The estimate of agent 1 is a linear function of the true weights \mathbf{w} , since $\mathbf{a}_1 = \mathbf{Q}_1\mathbf{w}$, so $\hat{a}_1 = \mathbf{q}'\hat{\mathbf{w}}_1 = \mathbf{q}'\mathbf{Q}_1^\top(\mathbf{Q}_1\mathbf{Q}_1^\top)^{-1}\mathbf{a}_1 = \mathbf{q}'\mathbf{P}_1\mathbf{w}$. Conditioning on agent 2's data, the posterior for \mathbf{w} is $N(\hat{\mathbf{w}}_2, \Sigma_2)$ with $\Sigma_2 = \sigma^2(\mathbf{I} - \mathbf{P}_2)$. The pair (a, \hat{a}_1) is therefore jointly Gaussian, since both are linear functions of \mathbf{w} . Their joint distribution conditional on agent 2's data has:

- $E[a|\mathbf{a}_2] = \mathbf{q}'\hat{\mathbf{w}}_2 = \mu_2$
- $E[\hat{a}_1|\mathbf{a}_2] = \mathbf{q}'\mathbf{P}_1\hat{\mathbf{w}}_2 = \mu_{2,1}$
- $\text{Var}(a|\mathbf{a}_2) = \mathbf{q}'\sigma^2(\mathbf{I} - \mathbf{P}_2)\mathbf{q} = \sigma_2^2$
- $\text{Var}(\hat{a}_1|\mathbf{a}_2) = \mathbf{q}'\mathbf{P}_1\sigma^2(\mathbf{I} - \mathbf{P}_2)\mathbf{P}_1^\top\mathbf{q} = \sigma_{1|2}^2$
- $\text{Cov}(a, \hat{a}_1|\mathbf{a}_2) = \mathbf{q}'\sigma^2(\mathbf{I} - \mathbf{P}_2)\mathbf{P}_1^\top\mathbf{q} = \sigma_{21}$

So the covariance matrix of (a, \hat{a}_1) conditional on agent 2's data is:

$$\begin{pmatrix} \sigma_2^2 & \sigma_{21} \\ \sigma_{21} & \sigma_{1|2}^2 \end{pmatrix}$$

For any joint Gaussian vector, the conditional distribution of the first component given the second is again Gaussian with

$$\mu_{2|1} = \mu_2 + \frac{\sigma_{21}}{\sigma_{1|2}^2}(\hat{a}_1 - \mu_{2,1}), \quad \sigma_{2|1}^2 = \sigma_2^2 - \frac{\sigma_{21}^2}{\sigma_{1|2}^2}.$$

Identifying $\kappa = \sigma_{21}/\sigma_{1|2}^2$ yields the stated result. □

Proposition 5 (Conditions for valuable two-stage updating). *In the setting of Proposition 4, consulting agent 1 provides value to agent 2 if and only if:*

$$\mathbf{q}'(\mathbf{I} - \mathbf{P}_2)\mathbf{P}_1^\top\mathbf{q} \neq 0$$

When this condition holds:

- The posterior mean changes: $\mu_{2|1} \neq \mu_2$
- The posterior variance decreases: $\sigma_{2|1}^2 < \sigma_2^2$

When this condition fails, consulting agent 1 provides no additional information: $\mu_{2|1} = \mu_2$ and $\sigma_{2|1}^2 = \sigma_2^2$.

Proof. From Proposition 4, the change in the posterior mean is $\mu_{2|1} - \mu_2 = \kappa(\hat{a}_1 - \mu_{2,1})$, and the change in posterior variance is $\sigma_2^2 - \sigma_{2|1}^2 = \kappa\sigma_{21}$, where $\kappa = \frac{\sigma_{21}}{\sigma_{1|2}^2}$ and $\sigma_{21} = \sigma^2 \mathbf{q}'(\mathbf{I} - \mathbf{P}_2)\mathbf{P}_1^\top \mathbf{q}$.
If $\sigma_{21} = 0$, then $\kappa = 0$, so $\mu_{2|1} = \mu_2$ and $\sigma_{2|1}^2 = \sigma_2^2$.
If $\sigma_{21} \neq 0$, then $\kappa \neq 0$ (since $\sigma_{1|2}^2 \geq 0$ with equality only when $\mathbf{P}_1(\mathbf{I} - \mathbf{P}_2) = \mathbf{0}$, which implies $\sigma_{21} = 0$). In this case, both the mean and variance will generally change unless $\hat{a}_1 = \mu_{2,1}$, which occurs with probability zero.
Therefore, two-stage updating provides value if and only if $\sigma_{21} = \sigma^2 \mathbf{q}'(\mathbf{I} - \mathbf{P}_2)\mathbf{P}_1^\top \mathbf{q} \neq 0$. \square

The intuition behind Proposition 5 is straightforward: **it is worthwhile to consult another agent if and only if the component of the question that you don't understand overlaps with the other agent's area of expertise.**

More precisely:

- $(\mathbf{I} - \mathbf{P}_2)\mathbf{q}$ represents the *residual* of the question after projecting it onto agent 2's own experience. This is the part of the question that agent 2 finds novel or unfamiliar.
- $\mathbf{P}_1^\top \mathbf{q}$ represents the component of the question that lies within agent 1's area of expertise (the row space of their experience matrix \mathbf{Q}_1).
- The condition $\mathbf{q}'(\mathbf{I} - \mathbf{P}_2)\mathbf{P}_1^\top \mathbf{q} \neq 0$ requires that these two components are not orthogonal—there must be some overlap between what agent 2 doesn't know and what agent 1 does know.

This formalizes the intuitive notion that collaboration is valuable when agents have *complementary* rather than identical or completely unrelated knowledge. If agent 1's expertise is orthogonal to the unfamiliar aspects of the question for agent 2, then agent 1's opinion provides no useful information. Conversely, if there is overlap between agent 2's knowledge gaps and agent 1's strengths, then consultation becomes valuable.

In the context of the ChatGPT model, this suggests that an AI assistant is most valuable for questions where:

1. The question contains elements that are novel to the human user (large $\|(\mathbf{I} - \mathbf{P}_2)\mathbf{q}\|$)
2. These novel elements fall within the AI's training domain (non-zero projection onto the AI's knowledge space)
3. The user's existing knowledge and the AI's training data are complementary rather than redundant

3 References
