

Opening the Black Box^{*}

A Statistical Theory of the Value of Data

Giovanni Colla Rizzi[†]

Latest Draft

Abstract

This paper develops a theory of the value of data for prediction. An agent chooses a sample of individuals and a subset of their observable characteristics (covariates) to estimate the parameters of a data-generating process and predict an outcome for a target individual based on her characteristics. I distinguish between covariates collected on the sample (training data) and covariates collected on the target individual (prediction data).

The main findings are: (i) training covariates exhibit economies of scope, as the value of one covariate is higher when others are also observed; (ii) the value of an additional training covariate is inverted-U-shaped in the sample size, so training covariates and observations are complements when data are scarce but become substitutes when data are abundant; and (iii) the value of a prediction covariate for the target individual is strictly increasing in the sample size.

These patterns have three policy implications. Mergers between firms holding different covariates can be privately profitable yet reduce welfare, especially when data are scarce (e.g., under strict privacy rules). Allowing firms to sell covariate bundles is always procompetitive because it removes double marginalization, whereas bundling observations can be anticompetitive when data are abundant. Finally, a data seller may profitably exclude one of several competing prediction providers even when this lowers total welfare.

JEL CLASSIFICATION: C11, D83, L12, O33.

KEYWORDS: Value of Data, Prediction, Economies of Scope, Data Markets.

^{*}I thank Patrick Rey for his guidance, and Jean-Pierre Florens and Doh-Shin Jeon for their suggestions. I also thank Jad Beyhum, Michele Biscaglia, Zhijun Chen, Krishna Dasaratha, Alexandre de Cornière, Eric Gautier, Johannes Gessner, Andrei Hagiu, Johannes Horner, Marco Iansiti, Bruno Jullien, Hiroaki Kaido, Simon Lortscher, Friedrich Lucke, Leonardo Madio, Giovanni Morzenti, Juan Ortner, Alessio Ozanne, Christoph Reidl, Andrew Rhodes, David Salant, Enrico Mattia Salonia, Maximilian Schaefer, Sara Shahanaghi, Tim Simcoe, Alex Smolin, Emanuele Tarantino, Ehsan Valavi, Marshall Van Alstyne, Davide Viviano, and Julian Wright, as well as participants at the European Association for Research in Industrial Economics conference (Valencia, 2025), the Questrom Digital Platforms Seminar, and seminars at TSE and Boston University.

[†]Toulouse School of Economics, University of Toulouse Capitole, France. E-mail: giovanni.rizzi@tse-fr.eu

1 Introduction

In digital markets, data drives competitive advantage, notably because it allows firms to make better predictions: Google and Meta use data to predict ad clicks, Amazon and Uber to predict the demand of goods and rides, and Spotify and Netflix to predict choices of songs and movies. Indeed, over the last decade, many digital markets have become dominated by a few firms that control vast amounts of data, attracting the attention of competition authorities. According to the 2021 U.S. House Report, for instance, “data advantages [...] can reinforce dominance and serve as a barrier to entry,”¹ a warning echoed in the EU Data Act proposal² and in major policy reports.³

These concerns are partly fueled by earlier statements from industry leaders: as Google’s CEO declared in 2009, “Scale is the key. We just have so much scale in terms of the data we can bring to bear.”⁴ More recently, however, technology firms have downplayed policymakers’ claims, arguing that data has diminishing returns because of the Law of Large Numbers, and that data concentration reflects the technological gap between firms rather than entry barriers.⁵

To assess these arguments, we must understand whether data exhibit economies of scale and whether there are economies of scope between different varieties of data. To do so, I study a model in which firm may collect: (i) data on many individuals (observations), and (ii) many attributes per individual (covariates). The firm must choose both how many people to sample and which features of them to measure. These choices determine how precisely it can estimate the unknown relationship between attributes and outcomes.

Therefore, I develop a theory of the value of data for an agent predicting a target variable generated by a linear process of infinitely many covariates. I distinguish between the value of additional *observations* (e.g., individuals) and additional *covariates* (e.g., their attributes). Furthermore, I distinguish between the *training covariates* the agent uses to learn the parameters of the process (e.g., tall people tend to like basketball), and the *prediction covariates* she uses to apply the learned parameters to the individual of interest (e.g., Martha is tall, so she probably likes basketball).

Specifically, the agent chooses how many observations and which training and prediction covariates to collect to minimize the out-of-sample prediction error, subject to data collection costs. I first characterize the optimal predictor for a given dataset. Building on this, I derive closed-form expressions for the value of data, showing that returns depend on how the variance is distributed across covariates. Three main insights on the economies of scope and scale of data emerge.

¹U.S. House of Representatives (2020), Investigation of Competition in Digital Markets, 117-40, pp. 36–38.

²European Commission (2023), Data Act proposal, COM(2023) 193 final.

³Digital Platforms (2019) and UK Competition and Markets Authority (2019).

⁴Schmidt (2009), “How Google Plans to Stay Ahead in Search,” Bloomberg, 2 October.

⁵Varian (2018) and Bajari et al. (2019).

First, there are *economies of scope in training*. Reducing regression noise by observing training covariates has a positive and accelerating effect on the precision of the parameter estimates: when the estimation is already relatively precise (low regression noise), further reductions in noise yield disproportionately large gains in precision. I name this convexity the *House Party Effect*, by analogy to a conversation in a noisy room.⁶

Second, *training covariates and observations are complements when they are scarce, but become substitutes when they are abundant*. Observations reduce the estimator variance arising from dimensionality and regression noise. Additional covariates increase dimensionality, which raises the value of observations, but also reduce regression noise, which lowers it. Adding a covariate increases the marginal value of an observation for an agent with few covariates, as the dimensionality effect dominates; the opposite is true for an agent with many covariates, as the noise-reduction effect dominates.

Third, *training data and prediction covariates are complements*. When the agent collects a new training covariate or observations, it increases the precision of the estimates of the parameters of the prediction covariates: there are positive spillovers of parameter estimate precision across training and prediction covariates.

I then allow the agent to select the most valuable covariates, namely those with the highest variance, and show that selection induces diminishing returns.⁷ When covariate variances are similar, the House Party Effect dominates, and returns to covariates are increasing. Instead, when variance is concentrated in a few covariates, the declining marginal variance from selection prevails, and returns may be diminishing.

Finally, I show that these findings generate sharp and policy-relevant implications in three applications, beginning with data-driven acquisitions. Even when firms operate in completely independent prediction markets, an incumbent with a richer set of covariates (e.g., Google) optimally acquires an entrant with distinct covariates (e.g., Fitbit). This incentive arises purely from economies of scope in covariates: merging datasets raises the marginal value of every covariate, producing a statistical analogue of a “natural monopsony” in data. As a result, incumbents acquire entrants even when there are no overlaps in products, users, or advertising markets, and even when entry by the smaller firm would be socially valuable due to knowledge spillovers, experimentation, or technological diversity. The application therefore illustrates a central implication of the theory: data consolidation is privately optimal well beyond the boundaries of traditional IO theories of killer acquisitions, and the tension between private and social incentives is exacerbated when the number of collected covariates is limited—for example due to privacy regulation or compliance costs.

⁶At a loud party with 100 guests (i.e., when there are 100 unobserved covariates), the departure of 5 guests (i.e., observing 5 covariates) hardly makes a difference in hearing one’s interlocutor (i.e., improves prediction accuracy), but when there are only 5 other guests (i.e., 5 unobserved covariates), the same reduction makes the conversation perfectly clear (i.e., prediction perfectly accurate), as the room falls silent except for the speakers themselves.

⁷The mechanism is analogous to Ricardo (1817) Law of Diminishing Returns for heterogeneous land quality: as the agent collects the most variable covariates first, the variance of the marginal covariate is diminishing.

Second, the model delivers sharp predictions for data pools, a policy tool explicitly encouraged in the EU Data Act and in several European common data-space initiatives. In my application, datasets used by prediction firms can be fragmented either in observations (i.e., firms pool same covariates on different users) or on covariates (i.e., firms pool different covariates on the same users). Pooling along these dimensions has fundamentally different welfare implications. Pooling covariates on the same individuals combines complementary features and always eliminates Cournot double marginalization between brokers, generating procompetitive efficiencies. By contrast, pooling samples of different individuals combines substitute datasets: when data is abundant competition between brokers is fierce and pooling is collusive.

Third, I apply the framework to exclusive licensing of training datasets in prediction markets. Motivated by recent agreements such as the 2024 Reddit–OpenAI deal, I study a setting in which a data seller holds a non-rival training dataset and licenses it to two competing prediction firms. Each firm can combine the licensed training data with its own proprietary covariates on target users, so that licensed and proprietary data are strict complements in generating prediction accuracy. An exclusive license allows the seller to commit not to supply the rival, raising the incumbent buyer’s willingness to pay, but it also changes the rival’s incentives to invest in proprietary data. The model shows that exclusivity is privately profitable whenever competition for prediction buyers is sufficiently intense, yet it can be socially harmful in an intermediate region where datasets are rich enough that both firms would invest under non-exclusive access. In that region, exclusivity amplifies incumbency advantages and forecloses a rival that would otherwise combine its own proprietary data with the shared training dataset, providing a microfoundation for regulatory scrutiny of data-licensing agreements such as Reddit–OpenAI.

Beyond these economic insights, the model also clarifies a puzzling pattern that has recently attracted attention in machine learning: in some predictive tasks, making a model more complex can first worsen and then improve its statistical performance. Using the equivalence between my Bayesian predictor and an optimally tuned ridge estimator, I show that the optimal degree of regularization—the force that shrinks coefficients to control overfitting—need not increase steadily with the number of regressors. When additional covariates are genuinely informative relative to the residual noise, expanding the model can actually reduce the optimal penalty, lowering rather than raising the degree of shrinkage. This provides a simple, closed-form mechanism for why prediction error can fall again in very rich models, a phenomenon sometimes called “double descent” in the machine-learning literature.

Roadmap. Section 2 introduces the data-generating process and the agent’s dataset-design problem. Section 3 derives the optimal predictor and characterizes the posterior distribution, highlighting the role of training information and the misspecification penalty. Section 4 analyzes the value of a dataset design, develops the core comparative statics on economies of

scale and scope, and studies optimal covariate selection under data-collection costs. Section 5 provides a frequentist interpretation of the Bayesian predictor and shows how the model yields a simple mechanism for non-monotonic regularization. Section 6 applies the framework to three policy-relevant environments—data-driven acquisitions, data pools, and exclusive licensing. Section 7 concludes. Technical proofs and extensions are collected in Section A and Section B.

Related Literature. There is a rich information design literature on the value of data, starting with Bergemann, Bonatti, and Smolin (2018) and continuing with Jones and Tonetti (2020), Bergemann, Bonatti, and Gan (2022), Bergemann and Bonatti (2024), and Acemoğlu et al. (2022). While this literature elegantly connects data value to the choice of a probability distribution by an agent, I link it directly to the realization of a dataset, allowing me to disentangle the statistical forces at play. Methodologically, my work is related to Montiel Olea et al. (2022), Iyer and Ke (2024), and Dasaratha, Ortner, and Zhu (2025), who analyze competition between different models with different covariates. In contrast, I jointly model covariates and observations, and distinguish training from prediction data, which allows me to derive structural non-convexities generating increasing returns and complementarities. Another related paper is Strzalecki (2024) interprets deviations from Bayesian updating through a misspecification parameter that is taken as primitive. In contrast, my model delivers an analogous distortion as "classic" Bayesian behavior under incomplete data collection, providing a microfoundation for this feature and reconnecting it to Bayesian updating.

I also contribute to the broad literature on economies of scope and scale to data. Whereas most models fix covariates and study returns to observations as Bajari et al. (2019) and Goldfarb and Tucker (2011), my framework endogenizes covariate collection. This extension rationalizes empirical findings on complementarities in Schaefer and Sapi (2023) and economies of scope in Carballa-Smichowski et al. (2025b). Schaefer (2025) develops a complementary frequentist approach and shows that the distribution of the covariate eigenspectrum shapes returns to scale, coherently with our results. Allcott et al. (2025) run a structural model to estimate returns to scale of additional observations in search. They finding diminishing returns and evidence of limited complementarities across different queries. Radner and Stiglitz (1984) attributes increasing returns to information costs, while I show they may emerge independently of costs.

My work provides microfoundations for two strands of literature that take increasing returns to data as assumptions: the IO literature on platforms and the macroeconomics literature on data as a production input. Prior work explains increasing returns through feedback between data and demand (Hagiu and Wright (2023), Prüfer and Schottmüller (2021), Farboodi and Veldkamp (2025), Aral, Brynjolfsson, and Wu (2008), and Cong, He, and Yu (2021)) or by assuming complementarities across datasets Carballa-Smichowski et al. (2025a), De Corniere and Taylor (2025), and Calzolari, Cheysson, and Rovatti (2025). I show instead that prediction

accuracy alone generates increasing returns due to the statistical structure of data, independent of demand feedback.

The applications of the model contribute more broadly to the IO literature on digital markets. De Corniere and Taylor (2025) shows that the pro- or anti-competitive effect of collecting more data only depends on the supply of data-driven services rather than its demand. Furthermore, Cornière and Taylor (2024) studies data-driven mergers developing a theory of harm of mergers which rely on cross-market effects. Both works complement the applications of my paper as they treat data as an undifferentiated good to which my results readily apply. Several papers on the economics of patents may be applied to datasets. Lerner and Tirole (2004) deals with complementarity/substitutability of patents and the private and social value of commercializing them jointly in pools. Gu, Madio, and Reggiani (2021) also study broker pools in a context without double marginalization and conclude they are anticompetitive when datasets are substitutes. I connect these two papers and microfound the welfare considerations on the statistical properties of datasets, complementing the findings in Nocke, Peitz, and Stahl (2007) by showing another source of potential social benefits of vertical integration and the avoidance of platform fragmentation. Katz and Shapiro (1986) show that exclusive deals are optimal when an inventor sells patents to competing firms, and this results in a suboptimal dissemination of patents; Aghion and Bolton (1987) shows that exclusive deals may serve as a strategic commitment device to deter entry by raising rivals' costs, allowing incumbents to extract rents from buyers at the expense of social welfare. My model develops insights from both papers to nonrival data markets with complementarities and endogenous investment in proprietary data.

Finally, I contribute to the literature on theoretical machine learning by developing a simple framework to study data scaling laws, providing a simple explanation of the phenomenon of double descent, i.e., that maximum likelihood-based algorithms generalize well even when overparametrized as explored in Hastie et al. (2020), Nakkiran et al. (2021), and Belkin et al. (2019).

2 Model Setup

An agent (she) seeks to predict a target variable of a target individual (he). To do so, she can collect (i) training data, consisting of covariates and past target variables for a sample of individuals; and (ii) prediction covariates for the target individual. The agent chooses the size of the training sample and which covariates to use in training and prediction, trading off prediction accuracy and data collection costs.

2.1 Data-Generating Process

The agent must predict a *target variable* $y \in \mathbb{R}$ for a *target individual*. For each individual i , the target variable is determined by

$$y^i = \sum_{k \in \mathbb{N}} \beta_k x_k^i,$$

where $x_k^i \in \mathbb{R}$ is the realization of covariate k for individual i , and β_k is its (unknown) impact on y . I normalize $\text{Var}[y] = 1$.

Parameters Parameters are mutually independent and independent of covariates, with prior distribution⁸

$$\beta_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1).$$

The prior variance reflects the agent’s uncertainty on the impact of covariate k on the target variable. Because the parameters are shared across individuals, they are learned through across-user learning in the sense of Hagiu and Wright (2023).

Covariates Covariates are mutually independent and i.i.d. across individuals, with prior distribution⁹

$$x_k^i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_k^2), \quad \sigma_k^2 \in (0, 1),$$

where I index covariates in decreasing order of variance,

$$\sigma_1^2 \geq \sigma_2^2 \geq \sigma_3^2 \geq \dots$$

Because the prior parameter variance is identical across covariates, σ_k^2 captures the amount of information contained in each covariate, i.e., the fraction of the variance of y due to covariate k . Moreover, let

$$S(\mathcal{K}) \equiv \sum_{k \in \mathcal{K}} \sigma_k^2,$$

denote the **cumulative variance of a covariate set** $\mathcal{K} \subseteq \mathbb{N}$, which is the fraction of the variance of y due to covariates in \mathcal{K} . Note that $S(\emptyset) = 0$ and $S(\mathbb{N}) = 1$.¹⁰

Notation For future reference, let $\mathcal{K}^c \equiv \mathbb{N} \setminus \mathcal{K}$ denote the complement of \mathcal{K} , $\mathbf{x}^i \equiv (x_k^i)_{k \in \mathbb{N}}$ denote the covariate vector, and $\boldsymbol{\beta} \equiv (\beta_k)_{k \in \mathbb{N}}$ denote the parameter vector. Furthermore, for any vector \mathbf{v} and set \mathcal{K} , let $\mathbf{v}_{\mathcal{K}} \equiv (v_k)_{k \in \mathcal{K}}$ denote the subvector of \mathbf{v} with indices in \mathcal{K} .

⁸Normalizing $\text{Var}[\beta_k] = 1$ is without loss of generality: we may recover any variance τ^2 by $\tilde{\sigma}_k^2 \equiv \tau^2 \sigma_k^2$.

⁹Appendix B relaxes the independence assumption by allowing arbitrary covariance structures. After using Principal Component Analysis to orthogonalize the covariates, all results carry through with the same notation and proofs.

¹⁰The upper bound follows from the independence of $\{\beta_k\}$ and $\{x_k^i\}$ and from $\text{Var}[y] = 1$.

2.2 Dataset Design and Predictor

I will distinguish between the way the agent collects the data (dataset design choice) and the way she uses the data (predictor choice).

Dataset Design To determine a dataset design, the agent chooses a set of *training covariates*, $\mathcal{T} \subset \mathbb{N}$, and the size of a sample of n individuals. Their realization is the *training data*,

$$(\mathbf{y}, \mathbf{X}),$$

where $\mathbf{y} \in \mathbb{R}^n$ and $\mathbf{X} \equiv \{\mathbf{x}_{\mathcal{T}}^i\}_{i=1}^n \in \mathbb{R}^{n \times |\mathcal{T}|}$. For any matrix of covariates \mathbf{X} , let $\mathcal{P}_{\mathbf{X}}(n, \mathcal{T})$, denote its prior distribution defined under the data-generating process above.

Furthermore, for the target individual, the agent collects a set of *prediction covariates*, $\mathcal{P} \subset \mathbb{N}$. Their realization is the vector

$$\mathbf{x}_{\mathcal{P}} \in \mathbb{R}^{|\mathcal{P}|}.$$

To summarize, the agent chooses a dataset design $(n, \mathcal{T}, \mathcal{P})$ defined as follows:

Definition 1. A *dataset design* is a triple $(n, \mathcal{T}, \mathcal{P})$, where $n \in \mathbb{N}$ is the sample size and $\mathcal{T}, \mathcal{P} \subset \mathbb{N}$ are the training and prediction covariate sets. Its *size* is $(n, |\mathcal{T}|, |\mathcal{P}|)$, where $|\mathcal{T}|$ is the *training (covariate) size* and $|\mathcal{P}|$ is the *prediction (covariate) size*.

The dataset design is the data collection strategy of the agent and reflects the choice of the experiments, sampling techniques and sensors which generate the data. As the model is linear, the dataset design also determines the number of parameters estimated in learning and, hence, the effective complexity of the training algorithm.

A given design $(n, \mathcal{T}, \mathcal{P})$, will induce an array of random variables, whose realization is a *dataset* comprising the training data and prediction covariates:

$$\mathbf{D} \equiv ((\mathbf{y}, \mathbf{X}_{\mathcal{T}}), \mathbf{x}_{\mathcal{P}}),$$

where we let $\mathcal{P}_{\mathbf{D}}(n, \mathcal{T}, \mathcal{P})$, denote the prior distribution of a dataset of design design $(n, \mathcal{T}, \mathcal{P})$, defined under the data-generating process in Section 2.1.

Predictor The agent's predictions will depend on the realization of the dataset. A *predictor* is a measurable map $\hat{y} : \mathbb{R}^{n \times (1+|\mathcal{T}|)} \times \mathbb{R}^{|\mathcal{P}|} \rightarrow \mathbb{R}$ that maps a dataset \mathbf{D} to a prediction

$$\mathbf{D} \mapsto \hat{y}(\mathbf{D}).$$

Equivalently, the prediction is the realization of a random variable which is a deterministic function of the dataset, i.e. the predictor. The choice of a predictor reflects how the agent plans to use the data it will observe to make a prediction.

2.3 Agent's Problem

Agent Utility Denoting the agent's prediction by $\hat{y} \in \mathbb{R}$, the realized prediction error is $y - \hat{y}$, and the agent suffers a quadratic loss so that

$$u(y, \hat{y}) = -(y - \hat{y})^2.$$

Therefore, agent is penalized proportionally more for larger prediction errors. Because of its tractability, quadratic loss is a standard assumption in econometrics and machine learning. Furthermore it has several theoretical properties discussed in Brier (1950) and axiomatized in Selten (1998).¹¹The agent will attempt to minimize her expected loss based on her knowledge on the data-generating process at the time of prediction. To increase her knowledge, she will choose a design which will produce in a dataset.

Design Cost I assume the cost of a design $(n, \mathcal{T}, \mathcal{P})$ depends only on the size of the sample and the covariate sets, and is additively separable:

$$C(n, \mathcal{T}, \mathcal{P}) = C_n(n) + C_t(|\mathcal{T}|) + C_p(|\mathcal{P}|).$$

This reduced form cost imposes a minimal structure on results, making the qualitative predictions of the model depend exclusively on the improvements in prediction accuracy stemming from data collection rather than from the function form of the cost of design. We assume that the choice of a predictor is costless.

Timing The agent's objective is to maximize expected utility minus dataset design cost. The model unfolds in three steps:

1. **Dataset Design Problem:** the agent chooses a dataset design $(n, \mathcal{T}, \mathcal{P})$ and pays the associated cost;
2. **Inference Problem:** given the dataset design, the agent chooses a predictor $\hat{y}(D)$; and
3. **Dataset Realization:** nature draws a dataset D from the data-generating process \mathcal{P} ; the prediction and the target variable are realized and losses occur.

As customary we will solve the model backwards. I will give a broad overview of the problem the agent faces, and will defer the solution to Sections 3 and 4.

¹¹Namely, among all incentive-compatible scoring rules—that is, rules for which truthful probabilistic forecasts maximize expected scores—the quadratic rule is uniquely characterized (up to a positive linear transformation) by three axioms: symmetry, elongation invariance, and neutrality.

2.3.1 Inference Problem

The agent will choose the predictor so that conditional on the realization of a dataset D , maximizes her expected utility, conditional on D , i.e.

$$U_D(\hat{y}) \equiv \mathbb{E}_{y|D} [u(y, \hat{y}(D)) | D].$$

Therefore, agent chooses the **optimal predictor**

$$\hat{y}^*(D) = \arg \max_{\hat{y}} U_D(\hat{y}).$$

For quadratic loss, $\hat{y}^*(D)$ exists and is unique.

2.3.2 Dataset Design Problem

The agent must choose which dataset design allows her to maximize her utility, having as outside option the choice not to collect any data and to make the ex-ante optimal predictor $\hat{y}^*(\emptyset)$. We will define the net ex-ante expected utility she derives from a design $(n, \mathcal{T}, \mathcal{P})$ as the value thereof.

Definition 2. The **value of a dataset design** $(n, \mathcal{T}, \mathcal{P})$ is

$$V(n, \mathcal{T}, \mathcal{P}) \equiv \underbrace{\mathbb{E}_{D \sim \mathcal{P}_D(n, \mathcal{T}, \mathcal{P})} [U_D(\hat{y})]}_{\text{Ex-ante Expected Utility of } (n, \mathcal{T}, \mathcal{P})} - \underbrace{\mathbb{E}_y [U_\emptyset(\hat{y})]}_{\text{Ex-ante Outside Option}}.$$

The value of a design $(n, \mathcal{T}, \mathcal{P})$ is the maximal increase in ex-ante expected utility that that agent can achieve when using the optimal predictor $\hat{y}^*(D)$.

Maximizing the utility minus data collection costs amounts to maximizing the value of a design minus its costs. Therefore, the agent will solve

$$\max_{n, \mathcal{T}, \mathcal{P}} \left\{ V(n, \mathcal{T}, \mathcal{P}) - (C_n(n) + C_t(|\mathcal{T}|) + C_p(|\mathcal{P}|)) \right\},$$

which, I will show, always has a solution. We can distinguish two subproblems.

Covariate Selection Problem Because costs depend only on the dimensions $|\mathcal{T}|$ and $|\mathcal{P}|$, for given a design size (n, t, p) , the choice of covariates dose not affect their cost. Therefore, I fix (n, t, p) and let the agent choose

$$\left(\tilde{\mathcal{T}}(n, t, p), \tilde{\mathcal{P}}(n, t, p) \right) \equiv \arg \max_{|\mathcal{T}| \leq t, |\mathcal{P}| \leq p} V(n, \mathcal{T}, \mathcal{P}),$$

which, we will show, always exist. The **optimal constrained covariate sets** $\tilde{\mathcal{T}}(n, t, p)$, $\tilde{\mathcal{P}}(n, t, p)$ are respectively the best training covariate set of size t and prediction covariate set of size p .

Accuracy Problem Having found the optimal constrained covariate sets, I obtain a micro-founded measure of the increase in prediction accuracy deriving from the best design of a given size (n, t, p) under the optimal predictor.

Definition 3. The *(surplus prediction) accuracy* of a design of size (n, t, p) is

$$A(n, t, p) \equiv \max_{\mathcal{T} \leq t, \mathcal{P} \leq p} V(n, \mathcal{T}, \mathcal{P}) = V\left(n, \tilde{\mathcal{T}}(n, t, p), \tilde{\mathcal{P}}(n, t, p)\right).$$

The agent now solves a three-dimensional optimization problem to select the design size which maximizes the difference between the surplus accuracy and the data collection cost:

$$(n^*, t^*, p^*) = \arg \max_{n, t, p} \left\{ A(n, t, p) - (C_n(n) + C_t(t) + C_p(p)) \right\},$$

which, we will show, always exist. Given the *optimal design size* (n^*, t^*, p^*) , the *optimal design* $(n^*, \mathcal{T}^*, \mathcal{P}^*)$ follows from

$$(\mathcal{T}^*, \mathcal{P}^*) \equiv \left(\tilde{\mathcal{T}}(n^*, t^*, p^*), \tilde{\mathcal{P}}(n^*, t^*, p^*) \right).$$

3 Inference

In this section we characterize the agent's choice of the optimal predictor for a given dataset design $(n, \mathcal{T}, \mathcal{P})$.

3.1 The Optimal Predictor

Before observing the dataset generating by the design, the agent can choose a predictor to specify how to use it to make predictions. Under quadratic loss, the optimal predictor is the posterior mean of the target variable, a well known Bayesian result that motivates the following lemma.

Lemma 1 (Optimal Predictor). *The optimal predictor is*

$$\hat{y}^*(D) = \mathbb{E}[y \mid D] = \mathbf{x}'_p \mathbb{E}[\boldsymbol{\beta}_p \mid (y, \mathbf{X}_{\mathcal{T}})].$$

Because the target variables are independent across individuals conditional on $\boldsymbol{\beta}$, the training data affect the optimal predictor only through the posterior distribution of $\boldsymbol{\beta}$. Furthermore, since covariates are mutually independent and independent from $\boldsymbol{\beta}$, the agent does not condition her predictions on the unobserved covariates of the target individual (those in $\mathbb{N} \setminus \mathcal{P}$), because their prior mean is zero.

Lemma 1 has a straightforward consequence for the value of a dataset design, which will depend on the training information on each parameter, which is thus defined:

Definition 4. The *(expected) training information on parameter k* is

$$\tau_k^2(n, \mathcal{T}) \equiv \underbrace{\text{Var}[\beta_k]}_{=1} - \mathbb{E}_{\mathcal{P}_{(y, X_{\mathcal{T}})(n, \mathcal{T})}} [\text{Var}[\beta_k \mid (y, X_{\mathcal{T}})]] .$$

The training information is the expected reduction of uncertainty on parameter k , which is fixed for the design $(n, \mathcal{T}, \mathcal{P})$. Armed with Definition 4, we can use Lemma 1 to characterize the value of the design $(n, \mathcal{T}, \mathcal{P})$.

Lemma 2. *The value of a dataset design $(n, \mathcal{T}, \mathcal{P})$ is*

$$V(n, \mathcal{T}, \mathcal{P}) = \text{Var}_{\mathcal{P}_D(n, \mathcal{T}, \mathcal{P})} [\hat{y}^*(D)] = \sum_{k \in \mathcal{P}} \sigma_k^2 \tau_k^2(n, \mathcal{T}).$$

The value of a dataset design is the ex-ante variance of the optimal predictor $\hat{y}^*(D)$, taken over the prior distribution of the dataset $\mathcal{P}_D(n, \mathcal{T}, \mathcal{P})$. Intuitively, the variance of the optimal predictor reflects its sensitivity to the realizations of the dataset. Therefore, data designs which yield optimal predictors that “pick up” a large fraction of the variance in the dataset have a higher value. Furthermore, the variance of the optimal predictor is the linear combination of the prediction covariate variance σ_k^2 and the training information on its parameter τ_k^2 . As the training information depends on the expected posterior variance, we will characterize the posterior distribution in the following section.

3.2 Posterior Distribution

To characterize the posterior distribution, observe that because individuals are independent conditionally on β and covariates are independent, knowing the prediction covariates but not the realization of the target variable gives no information on β . Hence, the posterior will only depend on the training data $(y, X_{\mathcal{T}})$.

Recall that, for each individual i ,

$$y^i = \sum_{k \in \mathbb{N}} \beta_k x_k^i = \sum_{k \in \mathcal{T}} \beta_k x_k^i + \sum_{k \in \mathbb{N} \setminus \mathcal{T}} \beta_k x_k^i,$$

so that, in matrix notation,

$$y = X_{\mathcal{T}} \beta_{\mathcal{T}} + \varepsilon_{\mathcal{T}}, \quad \varepsilon_{\mathcal{T}} \equiv X_{\mathcal{T}^c} \beta_{\mathcal{T}^c}. \quad (1)$$

Because covariates are mutually independent across k and i.i.d. across individuals, and independent of β , the noise term $\varepsilon_{\mathcal{T}}$ is independent of $X_{\mathcal{T}}$ conditional on β . Moreover, for each fixed parameter vector β ,

$$\varepsilon_{\mathcal{T}}^i = \sum_{k \in \mathbb{N} \setminus \mathcal{T}} \beta_k x_k^i \mid \beta \sim \mathcal{N}\left(0, v(\beta)\right), \quad v(\beta) \equiv \sum_{k \in \mathbb{N} \setminus \mathcal{T}} \sigma_k^2 \beta_k^2,$$

and the $\varepsilon_{\mathcal{T}}^i$ are independent across i . Hence, the exact likelihood under the data-generating process in Section 2.1 is

$$y \mid \boldsymbol{\beta}, X_{\mathcal{T}} \sim \mathcal{N}\left(X_{\mathcal{T}}\boldsymbol{\beta}_{\mathcal{T}}, v(\boldsymbol{\beta}_{\mathcal{T}^c})\mathbf{I}_n\right), \quad v(\boldsymbol{\beta}_{\mathcal{T}^c}) = \sum_{k \in \mathbb{N} \setminus \mathcal{T}} \sigma_k^2 \beta_k^2. \quad (2)$$

This differs from the standard conjugate Bayesian linear regression (e.g., see DeGroot (2005) and Berger (1990)), in which the error variance is a fixed constant: here, the variance $v(\boldsymbol{\beta})$ is a nuisance parameter that depends on $\boldsymbol{\beta}_{\mathcal{T}^c}$, so the same parameter vector $\boldsymbol{\beta}$ appears in both the mean and the variance of the likelihood. As a result, the Gaussian prior in Section 2.1 is no longer conjugate to the exact likelihood (2).

Tail variance and moment matching. The following result shows that the residual variance is almost surely deterministic for most covariate sets.

Lemma 3. *For any nested sequence $(\mathcal{T}_m)_{m \geq 1}$ with $\mathcal{T}_m \subset \mathcal{T}_{m+1}$ and $\bigcup_m \mathcal{T}_m = \mathcal{T}$,*

$$\sum_{k \in \mathbb{N} \setminus \mathcal{T}_m} \sigma_k^2 \beta_k^2 \xrightarrow[m \rightarrow \infty]{a.s.} \text{Var}(\varepsilon_{\mathcal{T}}) = 1 - S(\mathcal{T}_m). \quad (3)$$

For almost every draw of $\boldsymbol{\beta}$ from the prior, the random tail variance $v(\boldsymbol{\beta})$ concentrates around the deterministic tail sum $1 - S$ along nested truncations of the training covariate set, i.e. the impact of the residual omitted covariates on y is a white noise with variance $1 - S$. This shows that the random tail variance typically lies very close to its deterministic counterpart $1 - S$ and provides a natural justification for treating the omitted covariates as homoskedastic noise. As the cumulative variance of omitted covariates is fixed by the design, we will denote it by at

$$S \equiv S(\mathcal{T}).$$

Working Gaussian likelihood. Motivated by Lemma 3, I approximate the exact likelihood (2) by a homoskedastic Gaussian likelihood that matches the conditional mean and the unconditional variance of the data-generating process,

Assumption 1 (Working Gaussian likelihood). *For the purpose of inference, the agent uses the working likelihood*

$$y \mid \boldsymbol{\beta}_{\mathcal{T}}, X_{\mathcal{T}} \sim \mathcal{N}\left(X_{\mathcal{T}}\boldsymbol{\beta}_{\mathcal{T}}, (1 - S)\mathbf{I}_n\right).$$

This working likelihood coincides with the exact likelihood (2) when $v(\boldsymbol{\beta}) = 1 - S$, and Lemma 3 shows that the discrepancy between $v(\boldsymbol{\beta})$ and $1 - S$ vanishes almost surely along nested refinements of the training covariate set.

Under Assumption 1, the posterior distribution of $\boldsymbol{\beta}_{\mathcal{T}}$ is Gaussian, see DeGroot (2005) and Berger (1990). We will state the classic result in terms of a key variable:

Definition 5. The *misspecification penalty* is

$$\lambda \equiv \frac{1 - S}{n}.$$

Intuitively, this captures the fraction of variability in the target variable which is due to the regression residual rather than $\beta_{\mathcal{T}}$ rescaled by the sample size n , which the agent uses to average out that residual.

Proposition 1 (Posterior Distribution). *The posterior distribution of the parameter vector is given by:*

- For untrained parameters,

$$\beta_{\mathcal{T}^c} \mid (y, X_{\mathcal{T}}) \sim \mathcal{N}(\mathbf{0}_{|\mathcal{T}^c|}, I_{|\mathcal{T}^c|});$$

- For trained parameters ,

$$\beta_{\mathcal{T}} \mid (y, X_{\mathcal{T}}) \sim \mathcal{N}\left(\underbrace{(X'_{\mathcal{T}}X_{\mathcal{T}} + n\lambda \cdot I_{|\mathcal{T}|})^{-1} X'_{\mathcal{T}}y}_{\equiv \mathbb{E}[\beta_{\mathcal{T}} \mid (y, X_{\mathcal{T}})]}, \underbrace{\left(\frac{1}{n\lambda} X'_{\mathcal{T}}X_{\mathcal{T}} + I_t\right)^{-1}}_{\equiv \text{Var}[\beta_{\mathcal{T}} \mid (y, X_{\mathcal{T}})]}\right),$$

with $\beta_{\mathcal{T}}$ independent from $\beta_{\mathcal{T}^c}$.

Because parameters are independent, training $\beta_{\mathcal{T}}$ provides no information on the untrained parameters $\beta_{\mathcal{T}^c}$, whose prior mean is $\mathbf{0}$. The posterior mean of the trained parameters lies between the OLS estimator $(X'_{\mathcal{T}}X_{\mathcal{T}})^{-1} X'_{\mathcal{T}}y$ and the prior mean $\mathbf{0}_{|\mathcal{T}|}$, and λ reflects the relative weight given to the prior mean. Intuitively, if $\lambda = 0$, then all the variation in the training data is due to the regression parameters $\beta_{\mathcal{T}}$ (i.e., the regression is perfectly specified). Therefore, the posterior mean coincides with the OLS estimator, which is most sensitive to the randomness in the training data, because it maximizes the likelihood. Instead, if $\lambda \rightarrow \infty$, then all the variation in the training data is due to the residual (i.e., it is totally uninformative concerning $\beta_{\mathcal{T}}$) so the posterior mean is equal to the prior.

The rescaled penalty $n\lambda = 1 - S$ in the posterior mean is equivalent to the misspecification parameter of Strzalecki (2024). In contrast to his model, where λ reflects an exogenous concern for misspecification, my framework endogenizes λ as the weight a Bayesian agent assigns to the prior because of the omitted parameters. The agent's implicit discounting of the likelihood thus emerges endogenously from the dataset design, rather than exogenously because of a preference for robustness.

Proposition 1 allows us to characterize the training information for any parameter k .

Corollary 1. *The k -training information is*

$$\tau_k^2(\lambda; n, \mathcal{T}) = \left(1 - \mathbb{E}_{X_{\mathcal{T}} \sim \mathcal{P}_{X_{\mathcal{T}}}(n, \mathcal{T})} \left[\left(\frac{1}{n\lambda} X_{\mathcal{T}}' X_{\mathcal{T}} + I_{|\mathcal{T}|} \right)^{-1} \right]_{kk} \right) \mathbf{1}(k \in \mathcal{T}).$$

The k -training information is decreasing in the misspecification penalty λ . Furthermore, it is increasing in σ_k^2 because the latter increases the expected variation of the training data in the space spanned by covariate k . However, it is also increasing in σ_j^2 with $j \in \mathcal{T} \setminus \{k\}$ because of higher-order interactions between covariates. Note that if there is no misspecification penalty, i.e., $\lambda = 0$, the agent learns parameter k perfectly and $\tau_k^2(0; n, \mathcal{T}) = 1$; conversely, if the training data is completely uninformative, i.e., $\lambda \rightarrow \infty$, the agent learns nothing on the parameters and $\lim_{\lambda \rightarrow \infty} \tau_k^2(\lambda; n, \mathcal{T}) = 0$.

4 Dataset Design

4.1 The Value of a Dataset Design

We characterize the value of a generic dataset design $(n, \mathcal{T}, \mathcal{P})$. To do so we use Lemma 2 and Corollary 1, endogenizing the misspecification penalty, because it now depends on the design though

$$\lambda(n, \mathcal{T}) \equiv \frac{1 - S(\mathcal{T})}{n}.$$

Proposition 2 (Value of Dataset Design). *The value of a dataset design $(n, \mathcal{T}, \mathcal{P})$ is*

$$V(n, \mathcal{T}, \mathcal{P}) = \sum_{k \in \mathcal{P} \cap \mathcal{T}} \sigma_k^2 \tau_k^2(\lambda(n, \mathcal{T}); n, \mathcal{T}).$$

Intuitively, the value of a design is the linear combination of covariate variances σ_k^2 and the information on its parameter, $\tau_k^2(\lambda(n, \mathcal{T}); n, \mathcal{T})$. Note that adding a covariate j to $\mathcal{T} \setminus \mathcal{P}$ has two positive spillover effects on the value of a covariate $k \in \mathcal{P}$:

1. **First-order spillovers**, mediated by the misspecification penalty and deriving from the fact that $\tau_k^2(\lambda; n, \mathcal{T})$ is increasing in λ and that $\lambda(n, \mathcal{T}) < \lambda(n, \mathcal{T} \cup \{j\})$: the reduction in the regression noise from j ; and
2. **Higher-order spillovers**, which derive from $\tau_k^2(\lambda; n, \mathcal{T} \cup \{j\}) > \tau_k^2(\lambda; n, \mathcal{T})$, meaning that the mixed moments of the empirical variance/covariance matrix help estimation.

Because of these higher-order spillovers, the values of covariates are connected in a complicated web which is analytically intractable. The following result shows that if the training size (i.e., the trained parameter count) grows slower than the sample size, the direct spillovers are negligible, yielding a tractable shrinkage factor approximation $\bar{\tau}_k^2(\lambda(n, \mathcal{T}))$.

Proposition 3 (Training Information Approximation). *The following asymptotic approximation holds*

$$V(n, \mathcal{T}, \mathcal{P}) = \sum_{k \in \mathcal{P} \cap \mathcal{T}} \sigma_k^2 \bar{\tau}_k^2(\lambda(n, \mathcal{T})) + O\left(\sqrt{\frac{|\mathcal{T}|}{n}} + \frac{|\mathcal{T}|}{n}\right),$$

where

$$\bar{\tau}_k^2(\lambda) \equiv \frac{\sigma_k^2}{\sigma_k^2 + \lambda}.$$

When the training size is small relative to the sample size, all cross-covariate spillovers and dimensionality effects are fully captured by the misspecification penalty. This approximation is typically violated in modern machine-learning architectures such as neural networks or large language models, which feature far more parameters than training tokens and therefore exhibit substantial higher-order spillovers. Therefore, the following assumption amounts to taking a conservative stance on the magnitude of positive spillovers.

Assumption 2 (Low-dimensional Regime). $\frac{|\mathcal{T}|}{n} \rightarrow 0$.

Under Assumption 2, the value of a design $(n, \mathcal{P}, \mathcal{T})$, becomes separable in the training covariates, conditional on λ , which absorbs all cross-covariate spillovers.¹² We are now ready to highlight four key properties which drive the returns to data:

- **Training Covariate Spillovers:** the misspecification penalty $\lambda(n, \mathcal{T})$ is decreasing in $k \in \mathcal{T}$, so adding a training covariate improves estimates of all trained parameters:

$$-\Delta_{\mathcal{T}}^k \lambda(n, \mathcal{T}) \equiv -[\lambda(n, \mathcal{T}) - \lambda(n, \mathcal{T} \setminus \{k\})] = \frac{\sigma_k^2}{n} > 0;$$

- **Sampling Spillovers:** the misspecification penalty $\lambda(n, \mathcal{T})$ is decreasing in n , so adding an observation improves estimates of all trained parameters:

$$-\Delta_n \lambda(n, \mathcal{T}) \equiv -[\lambda(n, \mathcal{T}) - \lambda(n-1, \mathcal{T})] = \frac{1 - S(\mathcal{T})}{n(n-1)} > 0;$$

- **Spillover Substitution:** the sampling spillovers are decreasing in $k \in \mathcal{T}$ and covariate spillovers are decreasing in n :

$$|\Delta_n \lambda(n, \mathcal{T})| - |\Delta_n \lambda(n, \mathcal{T} \setminus \{k\})| = -\frac{\sigma_k^2}{n(n-1)} < 0.$$

- **House Party Effect:** the shrinkage factor $\bar{\tau}_k^2(\lambda)$ is convex in the misspecification penalty λ . This means that reductions in λ (which correspond to more precise estimation,

¹²Relaxing this assumption would require random-matrix methods to analyze the regime $t/n \rightarrow \gamma > 0$. However, closed-form solutions can then be obtained only under homoskedasticity, i.e. when all covariates have identical variance $\sigma_k^2 = \sigma^2 > 0$ for some fixed σ^2 . Because heterogeneity in covariate variances is central to the problem of covariate selection, we maintain the small- t/n approximation.

or equivalently to larger datasets) generate *increasing gains* the value of data as each marginal reduction in noise improves predictive accuracy by more than the previous one:

$$\frac{\partial^2}{\partial \lambda^2} \bar{\tau}_k^2(\lambda) = \frac{2\sigma_k^2}{(\sigma_k^2 + \lambda)^3} > 0.$$

The name “House Party Effect” refers to a useful analogy: two people are having a conversation at a crowded house party. When there are one hundred other people in the room (analogous to one hundred residual covariates contributing to noise), removing five people (observing five covariates/observations) has virtually no effect on how easily the two people can talk (on the precision of estimates). As the crowd thins (as we observe more and more covariates/observations), however, each person who leaves has a stronger impact on the overall noise level, e.g., once only five other people remain, their departure dramatically improves audibility as communication becomes perfectly clear. The same logic applies to the residual variance in the regression model: due to the convexity of $\bar{\tau}_k^2(\lambda)$, when the residual noise λ is large, small reductions in λ have negligible impact; but when residual noise is small, additional reductions sharply improve accuracy. Note that since $\frac{\partial^3}{\partial \lambda^3} \bar{\tau}_k^2(\lambda) < 0$ the House Party Effect strengthens as λ decreases, i.e. $\frac{\partial^2}{\partial \lambda^2} \bar{\tau}_k^2(\lambda)$ increases as n becomes large and \mathcal{T} expands.

4.2 Properties of the Value of a Dataset Design

We will now perform some comparative statics on the value of a dataset design to highlight the contribution of covariates and observations to training and prediction.

Properties of Prediction Covariates. For any \mathcal{P} , define the *contribution of prediction covariate* $k \in \mathcal{P}$ as

$$\Delta_k^{\mathcal{P}}(n, \mathcal{P}, \mathcal{T}) \equiv V(n, \mathcal{P}, \mathcal{T}) - V(n, \mathcal{P} \setminus \{k\}, \mathcal{T}).$$

Proposition 4 (Contribution of Prediction Covariate). *Fix (n, \mathcal{T}) . For any \mathcal{P} and any $k \in \mathcal{P} \cap \mathcal{T}$,*

$$\Delta_k^{\mathcal{P}}(n, \mathcal{P}, \mathcal{T}) = \sigma_k^2 \bar{\tau}_k^2(\lambda(n, \mathcal{T})) > 0,$$

which is independent of $\mathcal{P} \setminus \{k\}$. Furthermore, it is strictly increasing in \mathcal{T} and n .

This result has two consequences for the economics of prediction.

Corollary 2. *There are no spillover effects among training covariates.*

As prediction covariates enter in an additively separable way, the value each of them generates is independent from the collection other prediction covariate

Corollary 3. *Prediction covariates are complements to the training covariates and observations.*

Because of the training covariate spillovers and the sampling spillovers, collecting additional training data reduces the misspecification penalty $\lambda(n, \mathcal{T})$ for all the prediction covariates, meaning that there are complementarities between training and prediction covariates. The complementarity between prediction covariates and observations matches the empirical results of Schaefer and Sapi (2023) and Lee and Wright (2023), which however do not make the distinction between training and prediction covariates. The findings are also coherent with Wilson (1975): better information can be leveraged across the entire scale of production, so the non-rival nature of information generates complementarities.

Properties of Training Covariates. We next turn to the training covariates, define the *contribution of training covariate* $k \in \mathcal{T}$ as

$$\Delta_k^\mathcal{T}(n, \mathcal{P}, \mathcal{T}) \equiv V(n, \mathcal{P}, \mathcal{T}) - V(n, \mathcal{P}, \mathcal{T} \setminus \{k\}).$$

The following result characterizes the contribution of a training covariate.

Proposition 5 (Contribution of Training Covariate). *Fix (n, \mathcal{P}) . For any \mathcal{T} and any $k \in \mathcal{T} \setminus \mathcal{P}$,*

$$\Delta_k^\mathcal{T}(n, \mathcal{P}, \mathcal{T}) = \sum_{m \in \mathcal{P} \cap \mathcal{T}} \sigma_m^2 \int_{-\frac{\sigma_k^2}{n}}^0 -\frac{\partial}{\partial \lambda} \bar{\tau}_m^2(\lambda(n, \mathcal{T} \setminus \{k\}) + u) du > 0,$$

which is increasing in \mathcal{P} .

Proposition 11 shows that a training covariate $k \in \mathcal{T} \setminus \mathcal{P}$ contributes value only through its effect on the misspecification penalty. Adding k reduces $\lambda(n, \mathcal{T})$ by exactly σ_k^2/n , and because $-\partial_\lambda \bar{\tau}_m^2(\lambda) > 0$ for all $m \in \mathcal{P} \cap \mathcal{T}$, this reduction increases the predictive value of every prediction covariate m . The expression in Proposition 11 therefore aggregates these spillovers across prediction covariates. The contribution is strictly increasing in \mathcal{P} because more prediction covariates benefit from the same improvement in λ . The next result studies how the contribution of k changes when an additional training covariate j is already present. The double difference

$$\Delta_{jk}^{\mathcal{T}\mathcal{T}}(n, \mathcal{P}, \mathcal{T}) \equiv \Delta_k^\mathcal{T}(n, \mathcal{P}, \mathcal{T}) - \Delta_k^\mathcal{T}(n, \mathcal{P}, \mathcal{T} \setminus \{j\}),$$

captures exactly this interaction. A positive double difference means that training covariates are *complements*: the marginal value of k is higher when j is already included in the training design.

Using Proposition 11, we have the following result.

Proposition 1 (Economies of Scope in Training). *For $j, k \in \mathcal{T} \setminus \mathcal{P}$, so the double difference is*

$$\Delta_{jk}^{\mathcal{T}\mathcal{T}}(n, \mathcal{P}, \mathcal{T}) = \sum_{m \in \mathcal{P} \cap \mathcal{T} \setminus \{j\}} \sigma_m^2 \int_{-\sigma_k^2/n - \sigma_j^2/n}^0 \int_{-\sigma_k^2/n - \sigma_j^2/n}^0 \frac{\partial^2}{\partial \lambda^2} \bar{\tau}_m^2(\lambda(n, \mathcal{T} \setminus \{k\}) + u + s) ds du > 0.$$

$\Delta_{jk}^{\mathcal{T}\mathcal{T}}(n, \mathcal{P}, \mathcal{T})$ is strictly positive because $\bar{\tau}_m^2(\lambda)$ is convex in λ . Convexity (the House Party Effect) implies that a reduction in λ is more valuable when λ is already low. Therefore, when j is added first, the resulting decrease in λ increases the marginal gain from subsequently adding k . This is exactly the economies of scope captured by the double difference.

These results have important economic implications. They show that when the sample size is fixed, training covariates exhibit *economies of scope*: the marginal value of an additional training covariate is higher in richer training designs. This provides a theoretical foundation for the convex returns to covariates documented empirically by Carballea-Smichowski et al. (2025b). In environments where firms compete to acquire new covariates—for example through data auctions or licensing markets—this convexity creates a natural force towards concentration. If two firms bid for an additional covariate, the firm with the larger existing training set values the new covariate more and will win the auction. As a result, ownership of training covariates tends to become increasingly asymmetric over time, even in the absence of product-market overlap or traditional market power.

Proposition 6 (Complementarity/Substitutability in Training). *Fix (n, \mathcal{P}) . For any \mathcal{T} and any $k \in \mathcal{T} \setminus \mathcal{P}$,*

$$\frac{d}{dn} \Delta_k^{\mathcal{T}}(n, \mathcal{P}, \mathcal{T}) = \frac{\sigma_k^2}{n^2} \sum_{m \in \mathcal{P} \cap \mathcal{T}} \sigma_m^2 \left(\underbrace{\frac{1 - S(\mathcal{T} \setminus \{k\})}{\sigma_k^2} \int_{-\sigma_k^2/n}^0 \frac{\partial^2}{\partial \lambda^2} \bar{\tau}_k^2(\lambda(n, \mathcal{T} \setminus \{k\}) + u) du}_{(+)\text{ HPE} \times \text{Sampling spillovers}} + \underbrace{\frac{\partial}{\partial \lambda} \bar{\tau}_m^2(\lambda(n, \mathcal{T}))}_{(-)\text{ Spillover substitution}} \right) > 0$$

if and only if

$$n < \tilde{n}(\mathcal{T}, \mathcal{P}) \in (0, \infty),$$

which is decreasing in $\mathcal{T} \setminus \{k\}$ and \mathcal{P} .

Corollary 4. *Training covariates and observations are complements when they are scarce and substitutes when they are abundant.*

There are two opposite forces: on the one hand, the House Party Effect implies an additional observation raises the value of an additional training covariate, because it increases the marginal value of noise reduction; on the other hand, by spillover substitution the reduction in misspecification penalty brought forth by an additional training covariate is decreasing in the sample size. Both effects strengthen as n and \mathcal{T} expand; initially, the House Party Effect dominates, therefore, training data dimensions are complements, but eventually the spillover

substitution will prevail and they become substitutes. This implies that the complementarity between covariates and observations empirically documented in Schaefer and Sapi (2023) and Lee and Wright (2023) may not hold when datasets are large.

Properties of Sample Size. We now analyze the impact of a larger sample on the value of an additional observation. For any $n \geq 1$, the *marginal value of the n -th observation* is

$$\Delta^n(n, \mathcal{P}, \mathcal{T}) \equiv V(n, \mathcal{P}, \mathcal{T}) - V(n-1, \mathcal{P}, \mathcal{T}).$$

Proposition 7 (Marginal Value of Observation). *Fix $(\mathcal{T}, \mathcal{P})$. For any \mathcal{T} and any $k \in \mathcal{T} \setminus \mathcal{P}$,*

$$\Delta^n(n, \mathcal{P}, \mathcal{T}) = \sum_{m \in \mathcal{P} \cap \mathcal{T}} \sigma_m^2 \int_{-\frac{\lambda(n-1, \mathcal{T})}{n}}^0 -\frac{\partial}{\partial \lambda} \bar{\tau}_m^2(\lambda(n-1, \mathcal{T}) + u) du,$$

which is decreasing in n , increasing in \mathcal{P} and inverted U-shaped in \mathcal{T} .

Interestingly, there are two opposite forces driving the marginal value of an observation: on the one hand, the House Party Effect makes additional reductions in misspecification penalty λ more valuable by decreasing the argument of $\partial_\lambda \bar{\tau}_m^2$; on the other, the percentage reduction in λ brought forth by an additional observation is rapidly diminishing, which is reflected in the fact that the length of the integration interval is proportional to $1/n$, an instance of the Law of Large Numbers. The Law of Large Numbers always dominates the House Party Effect, so each additional observation improves predictions less and less.

Corollary 5. *There are diminishing returns to the sample size.*

This finding is consistent with Goldfarb and Tucker (2011), Bajari et al. (2019) and Schaefer and Sapi (2023). When covariate sets are fixed, data ownership has a natural tendency to become more distributed: if two agents were to participate in a second-price auction for an additional observation, the agent with a smaller sample would have a higher marginal value for that observation and would therefore acquire it, reducing the asymmetry in data ownership.

4.3 Constrained Covariate Selection

We will now characterize the optimal dataset design $(n, \mathcal{P}, \mathcal{T})$ under Assumption 2, with an additively separable data collection cost:

$$C_n(n) + C_t(t) + C_p(p) + F.$$

We will therefore solve the Covariate Selection Problem and the Accuracy Production Problem.

Proposition 8 (Optimal Covariate Selection). *The optimal constrained covariate sets are*

$$\tilde{\mathcal{T}}(t) = \bigcup_{k=1}^t \sigma_k^2, \quad \tilde{\mathcal{P}}(p) = \bigcup_{k=1}^p \sigma_k^2.$$

Recall that we ordered covariates in a sequence of decreasing variance $\{\sigma_k^2\}_{k \in \mathbb{N}}$, so the agent collects the most variable covariates first because $V(n, \mathcal{T}, \mathcal{P})$ is increasing in all σ_k^2 , directly for $k \in \mathcal{P}$ and through $\lambda(n, \mathcal{T})$ for $k \in \mathcal{T}$.

We will henceforth abuse notation and denote the misspecification penalty and the cumulative variance of the constrained covariate sets as

$$\lambda(n, t) \equiv \frac{1 - S(t)}{n}, \quad \text{and} \quad S(t) \equiv \sum_{k=1}^t \sigma_k^2.$$

Using these definitions, we can characterize the accuracy of a dataset design of a given size.

Proposition 9 (Accuracy). *Under Assumption 2, the accuracy of a data design of size (n, t, p) is*

$$A(n, p, t) = \sum_{k=1}^p \sigma_k^2 \bar{\tau}_k^2(\lambda(n, t)).$$

The accuracy reflects the interplay between the economies of scope and the diminishing returns deriving from the optimal covariate selection. The latter is the statistical equivalent of the Law of Diminishing Returns in Ricardo (1817): the agent treats covariates as production factors of heterogeneous quality and each additional covariate marginal productivity (variance) falls as one moves down the ordered list. Therefore there is a **covariate selection effect** which implies that $S(t)$ is concave in t . Consequently,

$$\frac{\partial^2}{\partial t^2} \lambda(n, t) < 0$$

meaning that collecting training covariates of higher index (lower variance) brings forth smaller reductions of the misspecification penalty.

4.4 Economies of Scale in Accuracy

Define the **marginal value of the p -th prediction covariate** as

$$\Delta^p(n, p, t) \equiv A(n, p, t) - A(n, p-1, t).$$

Proposition 10 (Marginal Value of Prediction Size). *Fix (n, t) . For any $p \in [1, t]$,*

$$\Delta^p(n, p, t) = \sigma_p^2 \bar{\tau}_p^2(\lambda^*(n, t)) > 0,$$

which is decreasing in p .

Given that the contribution of each covariate is independent from that of any other, only the selection effect is present and therefore the marginal value of the p -th prediction covariate is decreasing.

Corollary 6. *There are diminishing returns to prediction covariate size.*

Similarly, define the **marginal value of the t -th training covariate** as

$$\Delta^t(n, p, t) \equiv A(n, p, t) - A(n, p, t - 1).$$

Proposition 11 (Marginal Value of Training Size). *Fix (n, p) . For any $t > p$,*

$$\Delta^t(n, p, t) = \sum_{m \in \mathcal{P} \cap \mathcal{T}} \sigma_m^2 \int_{-\frac{\sigma_t^2}{n}}^0 -\partial_\lambda \bar{\tau}_m^2(\lambda(n, t - 1) + u) du.$$

Now increasing t decreases the interval of integration by $\frac{\sigma_{t+1}^2 - \sigma_t^2}{n}$. Furthermore, λ is convex decreasing in t , so higher t s reduce the misspecification penalty less. These two effects counter the House Party Effect and may induce diminishing returns to covariates if the selection effects are large enough, i.e. if the variance is concentrated in few covariates.

In general, the expression of the accuracy for an arbitrary distribution of variance across covariates σ_k^2 will be complicated. We will focus on the specific case which yields closed form accuracies.

4.4.1 Homoskedastic Covariates

Assume there is a set $\mathcal{I} \subset \mathbb{N}$ such that covariates are informative if and only if they are in \mathcal{I} ,

$$\sigma_k^2 = \sigma^2 \mathbf{1}(k \in \mathcal{I}),$$

so under the normalization $\text{Var}[y] = 1$, $|\mathcal{I}| = 1/\sigma^2$. Under optimal covariates selection, the agent will observe the informative ones first, so the cumulative variance is

$$S(k) = \min\{k\sigma^2, 1\}.$$

Proposition 12 (Closed Form Accuracy). *If $\sigma_k^2 = \sigma^2 \mathbf{1}(k \in \mathcal{I})$, the accuracy is*

$$A(n, p, t) = \frac{S(\min\{t, p\})}{\frac{1-S(t)}{n\sigma^2} + 1}.$$

We can perform some simple comparative statics.

Corollary 7. $A_{tt}(n, p, t) > 0$.

This implies that there are economies of scale to the training covariate size since there are no selection effects.

Proposition 13. *Provided $t \geq p$,*

$$A_{tn}(n, p, t) > 0 \iff t \leq \hat{t}(n) \equiv 1 - n\sigma^2.$$

Since $\hat{n}(k)$ is decreasing in k , covariates and observations are complements in accuracy when they are scarce and substitutes when they are abundant, as we had uncovered in Proposition 6.

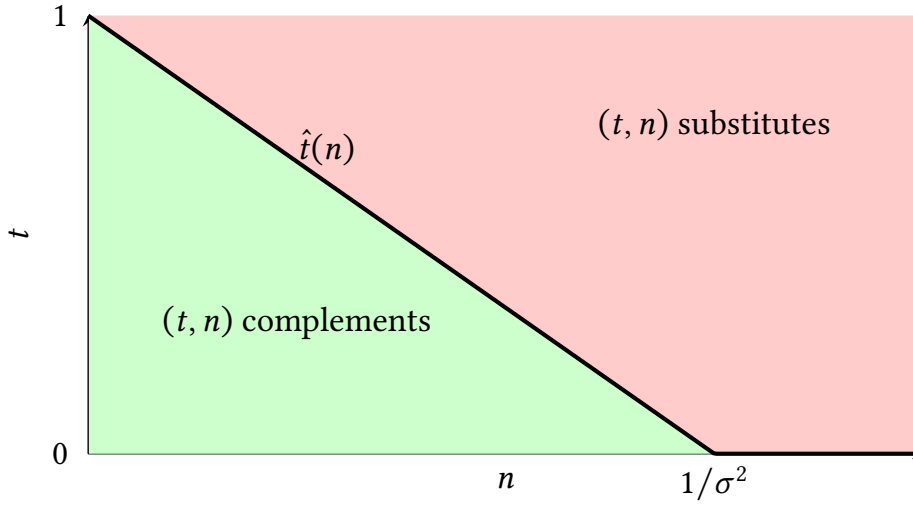


Figure 1: Covariates and observations are complements for all (n, t) below the decreasing threshold $t = \hat{t}(n)$, where additional observations raise the marginal value of training covariates. Above this curve, the two inputs become substitutes, as further observations reduce the incremental accuracy gain from enlarging the covariate set.

4.5 Comparative Statics

Assuming that an interior solution exists the comparative statics will be given by the substitution/complementarities I have highlighted in this section.

An increase in the cost of prediction covariates $C_p(\cdot)$ will reduce the amount of training covariates and observations as prediction and training are complements.

An increase in the cost of training covariates $C_t(\cdot)$ reduces the optimal number of training covariates and, via the optimality condition linking training and prediction scope, lowers the number of prediction covariates, thereby raising the marginal variance of the last prediction covariate purchased. Moreover, firms with few training covariates and observations —operating in the region where t and n are complementary— optimally cut back on both margins, while firms with abundant t and n —operating in the region where they are sub-

stitutes — will reduce training covariates but compensate cost shock on t by expanding their collection of observations.

An increase in the cost of observations $C_n(\cdot)$ has analogous effects: it lowers the optimal number of observations and prediction covariates and reshapes training choices. Firms with scarce t and n , where observations and training covariates are complementary, cut both margins; firms with abundant t and n , where they are substitutes, shift towards relatively more training covariates.

Therefore increases in training costs in the form of $C_n(\cdot)$ or $C_t(\cdot)$, amplify disparities in data collection between data-rich and data-poor agents.

5 Frequentist Interpretation

The reader who is not interested in the statistical properties of the model can safely skip to Section 6.

Ridge Regression Interpretation To give an intuitive interpretation of the results in Section 3, we will draw a connection to a common estimator in frequentist statistics, the ridge estimator. It is well known in the Bayesian statistics literature that optimal estimators of linear models, like that in Proposition 1, have a frequentist counterpart in the ridge estimator defined as:

$$\hat{\beta}_{\mathcal{T}}(\xi; (\mathbf{y}, \mathbf{X}_{\mathcal{T}})) \equiv \arg \min_{\mathbf{b} \in \mathbb{R}^t} \left\{ \frac{1}{n} \|\mathbf{y} - \mathbf{X}_{\mathcal{T}} \mathbf{b}\|_2^2 + \xi \|\mathbf{b}\|_2^2 \right\},$$

where $\xi \geq 0$ is the regularization parameter which penalizes the squared Euclidean distance of \mathbf{b} from the origin.¹³ Intuitively, the ridge estimator is a form of shrinkage estimator which is used to stabilize parameter estimates when maximum likelihood methods have too great variance. The ex-ante optimal ξ is the minimizer of the *expected mean squared error (MSE)* as

$$\xi^* \equiv \arg \min_{\xi} \mathbb{E}_{(\mathbf{y}, \mathbf{X}_{\mathcal{T}})} \left[\|\hat{\beta}_{\mathcal{T}}(\xi; (\mathbf{y}, \mathbf{X}_{\mathcal{T}})) - \beta_{\mathcal{T}}\|_2^2 \right],$$

where the expectation is taken with respect to the prior distribution of the training data $(\mathbf{y}, \mathbf{X}_{\mathcal{T}})$. We can define the *optimal ridge estimator* as

$$\hat{\beta}_{\mathcal{T}}^*(\mathbf{y}, \mathbf{X}_{\mathcal{T}}) \equiv \hat{\beta}_{\mathcal{T}}(\xi^*; (\mathbf{y}, \mathbf{X}_{\mathcal{T}})).$$

The following result characterizes the optimal regularization parameter ξ^* .

Proposition 14 (Ridge Estimator Interpretation). *The ex-ante optimal regularization param-*

¹³See DeGroot (2005) and Berger (1990).

eter is equivalent to the misspecification penalty

$$\xi^*(n, \mathcal{T}) = \lambda(n, \mathcal{T}) \equiv \frac{1 - S(\mathcal{T})}{n}.$$

This implies that there is an equivalence between the posterior mean parameters and the a priori optimal ridge estimator. The posterior mean parameters and the optimal ridge estimator are equivalent

$$\mathbb{E}[\boldsymbol{\beta}_{\mathcal{T}} \mid (\mathbf{y}, \mathbf{X}_{\mathcal{T}})] = \hat{\boldsymbol{\beta}}_{\mathcal{T}}^*(\mathbf{y}, \mathbf{X}_{\mathcal{T}}).$$

This result is coherent with the implied regularization in a Bayesian linear model characterized in Lindley and Smith (1972). Corollary 5 implies that a Bayesian agent acts like frequentist agent would if she treated the Bayesian prior as the true DGP and selected the ridge penalty that minimizes expected MSE across all possible samples generated from the prior distribution. Therefore, a Bayesian agent endogenously solves a bias–variance trade-off ex ante.

Optimal Regularization Since regularization applies uniformly to all parameters, the aggregate strength of regularization is defined as

$$\Lambda \equiv t\lambda = \frac{t(1 - S(t))}{n},$$

where t is the number of covariates and $S(t) \in [0, 1]$ is increasing and concave by Proposition 8. Conventional statistical intuition suggests that Λ should increase in t : richer models require stronger regularization to curb estimator variance.

However, Λ can decrease if we add a covariate to the model when t is large, specifically if

$$\Lambda(n, t + 1) \leq \Lambda(n, t) \iff t \times \sigma_t^2 \geq 1 - S(t),$$

Adding a covariate reduces the overall rate of regularization Λ if and only if dimensionality t and covariate variance σ_k^2 are large relative to the population regression noise $1 - S(t)$. This is because it affects the optimal rate of regularization through two opposing forces:

- an *intensive margin* (positive effect): the inclusion of one additional parameter mechanically increases the total regularization by $\lambda(n, t - 1) = [1 - S(t)]/n$;
- an *extensive margin* (negative effect): expanding the set of covariates reduces the optimal regularization applied to each parameter, captured by

$$-t \times [\lambda(n, t - 1) - \lambda(n, t)] = -\frac{t}{n} \times \sigma_k^2 < 0.$$

When the latter prevails, adding a covariate to a regression reduces the strength of the optimal regularization.

Double Descent Recent work by Belkin et al. (2019) and Nakkiran et al. (2021) have shown that when t becomes large, the optimal level of overall regularization Λ may actually decline—a phenomenon known as double descent. This finding has drawn significant attention as one of the central puzzles in modern machine learning theory. This decomposition offers a particularly transparent interpretation of double descent—arguably simpler than existing explanations such as Hastie et al. (2020)—and, to the best of my knowledge, the only one that does not rely on high-dimensional asymptotics.

Illustration of Regularization Dynamics To build intuition for how the optimal mis-specification penalty, I provide assume that the informativeness function follows the smooth concave form

$$S(t) = 1 - e^{-\rho t},$$

where $\rho > 0$ reflects the dispersion of variance across covariates, i.e., the strength of the selection effect, or equivalently measures the rate at which additional covariates reduce uncertainty. The implied penalty is therefore

$$\lambda(n, t) = \frac{e^{-\rho t}}{n},$$

and the aggregate strength of regularization is

$$\Lambda(t) = t\lambda(n, t) = \frac{te^{-\rho t}}{n}.$$

The figure below plots $\Lambda(t)$ for three values of ρ . This makes transparent the non-monotonicity at the heart of double descent.

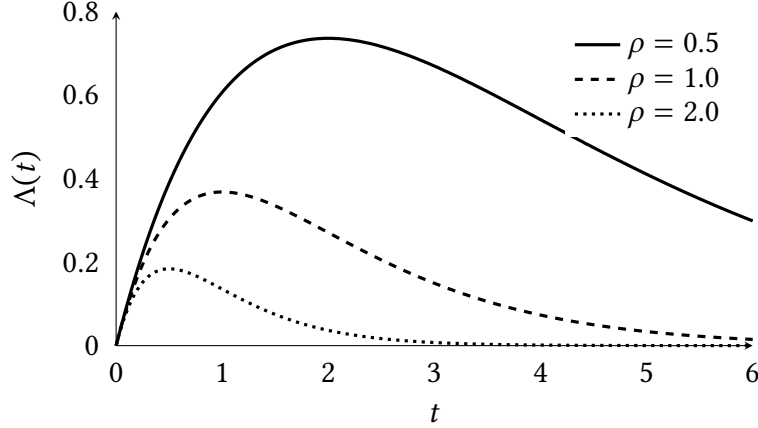


Figure 2: **Regularization as a function of model size.** The curves plot $\Lambda(t) = te^{-\rho t}/n$ for three values of ρ . For small t , the intensive margin dominates: adding a parameter increases total regularization. When t becomes large, the extensive margin dominates: informativeness grows slowly, $e^{-\rho t}$ declines rapidly, and the optimal level of regularization decreases. The turning point occurs at $t = 1/\rho$, where information accumulation and variance control exactly balance.

Intuition The function $te^{-\rho t}$ is maximized at $t = 1/\rho$. For $t < 1/\rho$, the intensive margin dominates, so $\Lambda(t)$ rises. For $t > 1/\rho$, the extensive margin dominates: informativeness grows too slowly to justify heavy regularization, and $\Lambda(t)$ declines.

This mirrors the empirical “double descent” pattern: once the model becomes sufficiently over-parameterized, the optimal amount of regularization falls even though the number of parameters continues to increase.

6 Applications

6.1 Natural Monopsony and Data-driven Acquisitions

In this application we show how a “natural monopsony” in data can emerge: a firm with more covariates has a strictly increasing and convex advantage in acquiring additional covariates. This leads to inevitable acquisitions even when products, users, and markets do not overlap. Mergers may be anti-competitive even when the merging firms do not compete for users, attention, advertising, or any traditional IO market.

6.1.1 Setup

Demand. For each prediction seller there is a unit mass of captive prediction buyers.¹⁴ Each buyer must choose an action $\hat{y} \in \mathbb{R}$, while the payoff-relevant state is $y \in \mathbb{R}$. Payoffs depend

¹⁴I assume there is no downstream competition to abstract from the usual “killer acquisition” motive: acquisitions are solely motivated by the acquisition of data.

on the squared error $(\hat{y} - y)^2$. Buyers share a common prior on y as specified in Section 2.1. If they do not buy any prediction, they choose the prior mean and obtain expected utility

$$\bar{u} = -1,$$

because of the normalization $\text{Var}(y) = 1$.

If a buyer purchases from seller i at price p_i , her *net* surplus relative to the outside option is

$$u(p_i, A_i) = A_i - p_i,$$

where $A_i \in [0, 1]$ is the value of the dataset design of seller i , which is equivalent to the buyers' willingness-to-pay for seller i 's prediction: each buyer strictly prefers i whenever $p_i < A_i$. We summarize individual and competitive demand by a reduced-form function $D(p_i, A_i)$, which is increasing in A_i and decreasing in p_i .

Supply. There are two potential sellers: an incumbent I ("Big Tech") and an entrant E ("Fit-bit"). Both observe the same sample size n . Seller $i \in \{I, E\}$ has access to a set \mathcal{K}_i of training and prediction covariates, with disjoint covariate sets $\mathcal{K}_I \cap \mathcal{K}_E = \emptyset$. To obtain closed-form expressions, assume all covariates have identical informativeness and impose symmetry by letting

$$|\mathcal{K}_I| = |\mathcal{K}_E| = k/2,$$

so that firms are symmetric.

The standalone value of seller i 's design is $A_i = V(n, \mathcal{K}_i)$.

The entrant E can either: (i) remain active and sell its own predictions, or (ii) accept a take-it-or-leave-it acquisition offer P from I . If I acquires E , it gains access to E 's covariates. The merged firm operates with sample size n and combined covariate sets, so post-merger design value is

$$A_2 = V(n, \mathcal{K}_I \cup \mathcal{K}_E).$$

Profits. Let $D(p, V)$ denote the demand faced by a monopolist offering a prediction with design value V at price p . There are two possible cases

- *If I acquires E :* Only the merged seller operates. Setting price p_2 , the incumbent's profit is

$$\Pi_I^{\text{acq}}(P, p_2) = p_2 D(p_2, A_2) - P,$$

while the entrant receives the acquisition payment

$$\Pi_E^{\text{acq}}(P) = P.$$

- *If there is no acquisition:* Both firms are active. As demands are not related because each

firm has a unit mass of captive customers, let $D_i(p_i, A_i)$ be seller i 's demands. Profits are

$$\Pi_i^{\text{no}}(p_i) = p_i D_i(p_i, A_i).$$

Planner welfare. The planner's objective is

$$W = \Pi_I + \Pi_E + CS - \xi \times \mathbf{1}(I \text{ acquires } E),$$

where CS is customer surplus and $\xi \geq 0$ is the social cost of reduced entry. Intuitively, ξ reflects positive spillovers, e.g., new data, experimentation on alternative prediction tasks, technological diversity, knowledge spillovers.

Timing. The game unfolds in two stages:

1. **Acquisition stage.** The incumbent I offers P to acquire E . The entrant accepts or rejects.
2. **Pricing stage.** If the acquisition occurs, the merged firm chooses p_2 . Otherwise, I and E simultaneously set (p_I, p_E) .

Payoffs are realized at the end. The solution concept is subgame perfect Nash equilibrium, obtained by backward induction.

6.1.2 Buyer Purchase

Prediction Sale Buyers purchase the prediction if and only if $V \geq p$. The resulting demand function is

$$D(p, V) = \mathbf{1}(p \leq V).$$

The seller therefore sets

$$p^* = V,$$

sells to all buyers, and obtains profit equal to the value of its dataset design:

$$\Pi_i^{\text{no}} = A_i.$$

Intuitively, with homogeneous buyers and no frictions, the platform can fully extract the expected gains from improved prediction.

Acquisition For the entrant E , the standalone outside option is the expected gains from improved prediction when it operates separately from I , earning

$$\Pi_E^{\text{no}} = A_1.$$

To acquire E , it is sufficient that the incumbent I offer it this amount. When I acquires E , it combines covariates and gains an increase of value

$$A_2 - A_1$$

compared to its outside option A_1 .

Proposition 15 (Harmful Acquisition). *The incumbent I always acquires the entrant E .*

Intuitively, the economies of scope stemming from Propositions 10 implies that

$$A_2 - \sum_{i \in \{I, E\}} A_i \geq 0.$$

Combining disjoint datasets makes each covariate more informative. This is due to two things: first, as far as prediction covariates \mathcal{K}_i can benefit from the covariate spillovers of the training covariates \mathcal{K}_{-i} (and vice versa); second, due to the House Party Effect, pushing training covariates from \mathcal{K}_i to $\mathcal{K}_I \cup \mathcal{K}_E$ yields supermodular gains. As the incumbent can fully exploit these complementarities, it will prevent entry by E even if demands are independent. This is the natural-monopsony effect: data-driven economies of scope push data into the hands of a single firm.

6.1.3 Planner's Problem

With homogeneous buyers and full surplus extraction, customer surplus is zero in equilibrium. Social welfare therefore equals industry profit net of any cost of foregone entry and spillovers:

$$W = \begin{cases} \Pi_2 - \xi, & \text{if } I \text{ acquires } E, \\ \Pi_I + \Pi_E, & \text{otherwise,} \end{cases}$$

where $\xi \geq 0$ denotes the social cost of eliminating the entrant as an independent firm.

Proposition 16. *The acquisition is harmful whenever*

$$A_2 - 2A_1 \geq \xi.$$

This condition is equivalent to

$$k \leq \tilde{k}(n, \xi) \equiv \frac{n+1}{2} \frac{\sqrt{\xi(8n+\xi)} - 3\xi}{n-\xi},$$

which is U-shaped in n and increasing in ξ .

The social planner compares the economies of covariate integration reflected in $A_2 - 2A_1$ and the spillovers ξ to evaluate whether the loss of spillovers is worth the gains in predic-

tion accuracy it generates. However, the incumbent does not internalize the loss of positive spillovers, generating a tension between private and social incentives for acquisition. Although scope economies make data combinations privately valuable, they also cause the incumbent to systematically over-acquire entrants. The incumbent captures all private complementarities but does not internalize the entrant's contribution to future innovation, data generation, or technological diversity. The resulting acquisitions are therefore privately efficient but potentially socially harmful—too many mergers occur from society's perspective.

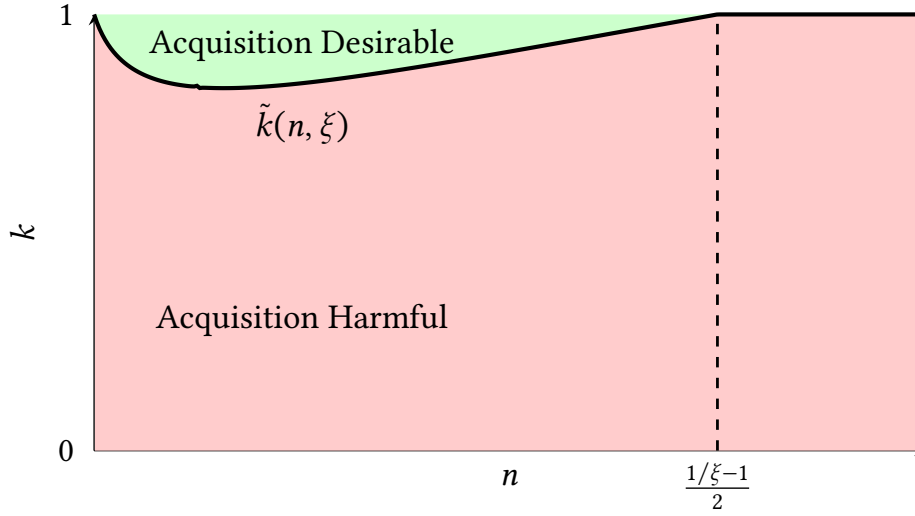


Figure 3: An acquisition is harmful for all (n, k) below the U-shaped threshold $k = \tilde{k}(n, \xi)$, where the loss of innovation spillovers exceeds the scope economies from combining datasets.

Regulation and Remedies Privacy regulation can worsen this tension. By raising compliance costs, the EU's General Data Protection Regulation (GDPR) lowers k , thereby dampening economies of scope and reducing the interval $[0, \tilde{k}(n, \xi)]$ in which the acquisition is socially desirable. Conversely, open-data initiatives raise k and broaden the interval in which economies of scope make acquisition desirable. These results suggest that privacy law, data governance, and merger policy must be coordinated: measures that restrict data access may unintentionally increase the likelihood of harmful data-driven acquisitions.

A potential remedy to potentially harmful acquisitions are FRAND APIs (Fair, Reasonable, and Non-Discriminatory Application Programming Interfaces). FRAND APIs are a regulated mechanism that guarantees standardized, non-exclusive access to a firm's data on fair technical and pricing terms. It preserves interoperability by allowing rivals to continue accessing the entrant's data after a merger. In this model, a FRAND API mitigates the dynamic loss ξ by keeping the entrant's data and experimentation available to the ecosystem, while still allowing the incumbent to capture the scope economies $\Delta(n, k)$. Therefore, the model provides

support to the FRAND access obligations adopted in EU digital regulation as merger remedies, most notably in the Digital Markets Act and the Data Act.

6.2 Data Pools

Data owners often form partnerships to pool their datasets and sell access jointly. For example, BMW, Mercedes-Benz, and Audi co-founded the platform *Here Mobility Data Marketplace*, which aggregates GPS, speed, and road-condition data from connected cars. These agreements can benefit society as they address an instance of the fundamental *complements problem* originally noted by Cournot (1838): if independent monopolists sell datasets that are complements, the resulting double marginalization leads to inefficiently high prices. Consequently, policymakers have encouraged pooling or sharing of datasets to overcome these frictions in the EU Data Act and in the European strategy for common data spaces. However, several papers have warned that brokers might use data sharing agreements to collude.

To assess the merits of these arguments, we develop a model in the spirit of Lerner and Tirole (2004): it turns out that the economics of combining datasets closely matches the economics of combining patents.

6.2.1 Setup

Data Owners Consider two data owners, each holding a dataset of identical informativeness. The complete dataset consists of n observations and k covariates, with the training and prediction covariates coinciding. However, the data may be split between the two owners either along the n (sample size) dimension or along the k (covariate) dimension. All parties are symmetrically informed about the informativeness of each dataset.

If the data are split along the sample dimension, each owner holds n observations and all k covariates. Pooling thus doubles the sample size. Conversely, if the data are split along the covariate dimension, each owner holds all n observations but only half the covariates, so pooling expands the covariate space. These two cases capture two distinct sources of complementarity: statistical precision (more n) and informational richness (more k).

Data Buyers There is a continuum of potential buyers (e.g., prediction firms) who can purchase access to one or both datasets and combine them without cost. Buyers are heterogeneous and indexed by $\theta \in [\underline{\theta}, \bar{\theta}]$, representing their adoption cost or opportunity cost of using the prediction technology.

A buyer of type θ who purchases access to $q \in \{1, 2\}$ datasets obtains a gross surplus

$$U_q = A_q - \theta,$$

where A_q denotes the predictive value of having access to q datasets. Specifically,

$$A_q \equiv \begin{cases} A\left(\frac{2n}{q}, k\right), & \text{if the data are split along the } n \text{ dimension,} \\ A\left(n, \frac{2k}{q}\right), & \text{if the data are split along the } k \text{ dimension.} \end{cases}$$

Since $A(n, k)$ is increasing in both arguments, combining datasets strictly improves predictive accuracy. Intuitively, pooling along the n dimension increases the number of observations available for training, which reduces estimation error, whereas pooling along the k dimension increases the number of predictive features, which broadens the scope of prediction.

The heterogeneity parameter θ is distributed according to

$$G(\theta) = \theta^\alpha, \quad \alpha \in [0, 1],$$

implying a strictly increasing hazard rate $\frac{g(\theta)}{1-G(\theta)}$ with $g = G'$. The parameter α thus governs the curvature of demand: when α is low, heterogeneity is large and demand is relatively inelastic; when α is high, buyers are more homogeneous and demand becomes more elastic.

Given a price P for access to a bundle of $q \in \{1, 2\}$ datasets, only buyers with $A_q - \theta \geq P$ make the purchase. Hence, the corresponding demand function is

$$D(P; A_q) = \Pr(A_q - \theta \geq P) = (A_q - P)^\alpha.$$

This formulation implies that the elasticity of demand increases in α : higher α corresponds to a market in which adoption falls faster as price rises. Equivalently, α can be interpreted as the *semi-elasticity of demand*, describing how responsive adoption is to changes in price.

6.2.2 Pooling Price

When data owners form a pool, they coordinate pricing and behave as a single monopolist offering a bundled dataset. This situation mirrors the *pool pricing benchmark* in Section I.b of Lerner and Tirole (2004), where a patent pool chooses the package price P to maximize joint revenue given demand $D(P - V_2)$. The logic is analogous here: the pooled data bundle yields predictive value V_2 , and all buyers face a single posted price P .

Lemma 4 (Optimal Pool Price). *The optimal pool price is*

$$P^* = \frac{A_2}{\alpha + 1}.$$

The pool acts as a monopolist setting the joint-access price that equates marginal revenue to zero, just as in Equation (1) in Lerner and Tirole (2004). The resulting price is decreasing in the demand semi-elasticity α : when α is high, buyers are more sensitive to price, leading the

pool to charge less. Conversely, when α is low, demand is inelastic, so the pool can extract a larger fraction of the total value V_2 .

Economically, pooling internalizes the complementarity between the dataset, analogous to Cournot's double marginalization problem in complementary goods. When data are sold separately, each owner fails to account for the positive externality that lower prices have on the other's sales. By setting a single joint price, the pool eliminates this inefficiency and behaves as an integrated monopolist.

6.2.3 Covariate Silos

If the datasets are split in two sets of covariates, the datasets are complements by Proposition ?? . We can use the results in Section II in Lerner and Tirole (2004) to characterize the unique symmetric equilibrium in the case in which brokers selling complementary data do not form a pool.

Lemma 5 (Covariate Fragmentation Price). *If the brokers have distinct covariates each B_i prices at*

$$p_i = \frac{A_2}{2 + \alpha},$$

and the buyers will buy from both brokers.

Applying Proposition 1 in Lerner and Tirole (2004) directly yields the following result.

Proposition 17. *A pool of buyers with distinct covariates is always procompetitive.*

Covariate silos are an instance of the classic complements problem of Cournot (1838). Each broker prices without internalizing that lowering its price increases the value of the other dataset. Pooling eliminates this double marginalization and behaves exactly like integration in complementary patents in Lerner and Tirole (2004).

6.2.4 Observation Silos

If the datasets are split in two sets of covariates, the datasets are substitutes by Proposition 7. We can use the results in Section II in Lerner and Tirole (2004) characterizes the unique symmetric equilibrium in the case in which the brokers selling substitute data do not form a pool.

Lemma 6 (Sample Fragmentation Price). *If the brokers have different observations on the same covariates,*

$$p_i = \min \left\{ A_2 - A_1, \frac{A_2}{2 + \alpha} \right\},$$

and the buyers will buy from both brokers.

Applying Proposition 1 in Lerner and Tirole (2004) directly yields the following characterization of when fragmented sample pooling is procompetitive.

Proposition 18 (Welfare of Observation Pooling). *A pool of observations is procompetitive if and only if*

$$k < \hat{k}(n, \alpha) \equiv 1 - \frac{n}{2\alpha}.$$

The condition $k \leq \hat{k}(n, \alpha)$ mirrors the distinction between the *competition margin* and the *demand margin* in Lerner and Tirole (2004). When k is low, the incremental predictive value of the second block of observations, $A_2 - A_1$, remains substantial. Buyers therefore strictly prefer to acquire both datasets, so each broker must set a price low enough not to be excluded. Fragmented pricing is thus constrained by the competition margin, and the two datasets act as effective complements at the equilibrium prices. Separate pricing then reproduces the classic Cournot double-marginalization problem, and pooling internalizes this externality, lowering the total price.

When $k > \hat{k}(n, \alpha)$, the incremental value $A_2 - A_1$ becomes small. Buyers are nearly indifferent between acquiring one or two datasets, so each broker can raise its price without being displaced. Pricing is now governed by the demand margin, and a pool relaxes this competitive pressure. In this region, pooling raises the joint-access price and is therefore anti-competitive.

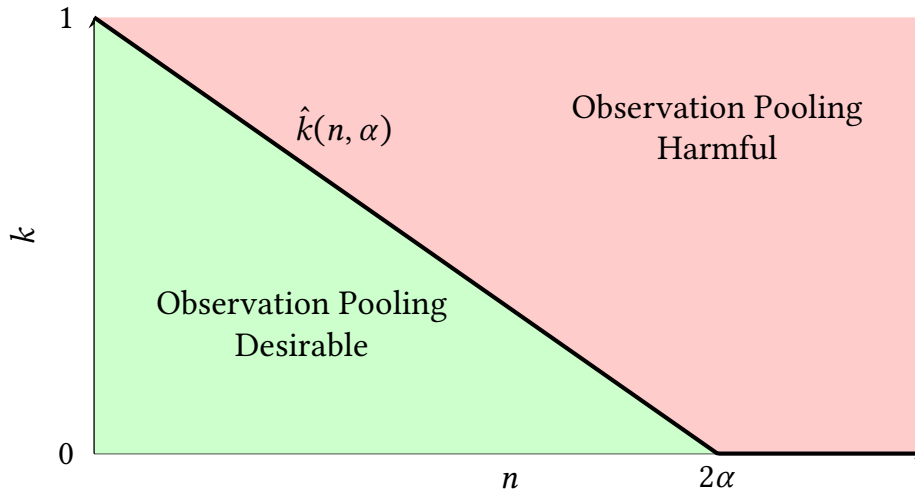


Figure 4: The threshold $k = \hat{k}(n, \alpha)$ separates the region in which fragmented pricing is constrained by the competition margin (pooling lowers prices) from the region in which it is constrained by the demand margin (pooling raises prices).

Because $\hat{k}(n, \alpha)$ is decreasing in n , the model offers a sharp statistical insight: observation pooling becomes anti-competitive when datasets are large. As n grows, posterior variance is already small, so doubling the sample yields only a negligible reduction in prediction error. The marginal value $A_2 - A_1$ of the second dataset therefore collapses. In this case, fragmentation disciplines prices through the competition margin—each broker must keep its price close to $A_2 - A_1$ to avoid being dropped—while pooling removes this discipline. As a result, pooling

observations raises the price of access and is anti-competitive in the upper-right corner of the (n, k) plane.

Regulation of Data Pools The model qualifies current policy discussions on data sharing. European legislation such as the EU Data Act and the European strategy for common data spaces broadly promotes data pooling and data spaces, which can overcome the Cournot complements problem in fragmented data markets. However, these policy initiatives do not distinguish between sources of dataset heterogeneity (covariates vs. observations) as a condition for pooling to be procompetitive. The model shows that this distinction is essential. Pooling heterogeneous user attributes (different covariates) always resolves a pure complements problem and is therefore procompetitive, fully consistent with the intuition in the Data Act Recitals that data portability and interoperability reduce market power. In contrast, pooling additional observations on the same attributes may be anticompetitive when $k > \hat{k}(n, \alpha)$: in this region the second dataset adds little predictive value, fragmented pricing would be disciplined by the competition margin, and a pool merely relaxes competitive pressure.

These results imply that policies encouraging sharing of data on similar attributes may unintentionally raise prices in markets where data is abundant. Conversely, when different attributes are shared, the model confirms the procompetitive rationale underlying EU data-space initiatives. More generally, the model highlights the need for analytical tools within merger control and data-access regulation to understand the level of heterogeneity of the shared data, a distinction absent from existing regulatory guidance.

6.3 Data Exclusivity

Recent exclusive data-licensing agreements—such as Reddit’s 2024 deals with OpenAI and Google—highlight broader concerns that proprietary access to datasets may distort competition in AI and prediction markets. Because data are non-rival, a data seller faces a dynamic commitment problem akin to that of a durable-good monopolist: once it has licensed the dataset to one firm, it is tempted to also license it to the rival, eroding the first buyer’s advantage. As in Katz and Shapiro (1986) and Aghion and Bolton (1987), exclusivity can serve as a contractual commitment device that mitigates opportunism by the seller but introduces a welfare trade-off: it softens business-stealing while depressing investment by excluded firms when data and proprietary inputs are complements. The model predicts that profitable exclusivity may become socially undesirable when datasets are abundant and product-market rivalry is intense. In such settings, exclusivity amplifies incumbency advantages and can deter entry—offering a micro-foundation for regulatory scrutiny of data-sharing agreements such as the Reddit–OpenAI deal.

6.3.1 Setup

There are three players: a data seller S (e.g., a platform holding user-generated data) and two prediction firms F_1 and F_2 . The firms compete to sell predictions to customers whose utility depends on the accuracy of the prediction.

Prediction Buyers. There is a unit mass of customers, divided into:

- a mass $s \in [0, 1]$ of *shoppers*, who can buy from either firm;
- a mass $(1 - s)/2$ of *captive customers* for each firm, who can only buy from that firm.

The parameter s captures the intensity of competition: when $s = 0$, all customers are captive and firms behave as local monopolists; when $s = 1$, all customers are shoppers and the market is fully competitive. This simple structure captures the idea that data quality matters only in relative terms, since shoppers migrate toward the firm offering the more accurate prediction.

Customer utility from firm $i \in \{1, 2\}$ is

$$u_i = A_i - p_i,$$

where A_i denotes the accuracy of the prediction and p_i is the price charged. Shoppers buy from the firm with the highest net utility $A_i - p_i$, while captive customers always purchase from their incumbent firm.

Data Buyers. We assume covariates are ex-ante identical and that training and prediction covariates coincide. Each firm $i \in \{1, 2\}$ starts with no data and can improve its prediction quality through two channels:

1. *Training Data Licensing.* Firm i can buy a license for the seller's training dataset of k covariates and the realization of the target variable for n observations, paying a fee f . This choice is represented by a binary variable $\ell_i \in \{0, 1\}$.
2. *Proprietary Data Collection.* Firm i can collect the k covariates for the target user, paying a fixed cost $c > 0$, represented by a binary choice $r_i \in \{0, 1\}$.

A firm's data strategy is $(\ell_i, r_i) \in \{0, 1\}^2$, the prediction quality of firm i is

$$A_i \equiv A(\ell_i, r_i) = \bar{A} \times \mathbf{1}(\ell_i r_i = 1),$$

where

$$\bar{A} \equiv \frac{k}{\frac{1-k}{n} + 1}.$$

Intuitively, in order to make a prediction, the firm must license the training data and invest to collect proprietary data on the target individual.

Pricing and Profits. Firms can price discriminate between captive customers and shoppers. Denote prices by p_i^c (captives) and p_i^s (shoppers). Let D_i^s denote firm i 's share of shoppers, with $D_1^s + D_2^s = 1$. Firm i 's revenue from prediction is then

$$\pi_i = p_i^c \frac{1-s}{2} + p_i^s s D_i^s,$$

and total profit is

$$\Pi_i = \pi_i - c r_i - f \ell_i.$$

The cost terms capture the trade-off between acquiring data and improving prediction quality: collecting proprietary data entails the fixed cost c , while licensing requires the payment f to the data seller.

Data Seller. The data seller S has a monopoly over the training dataset and sets a take-it-or-leave-it license fee f .¹⁵ The seller's profit is

$$\Pi_S = f \sum_{i \in \{1,2\}} \ell_i.$$

Because data are non-rival, S could in principle sell to both firms without loss of quality, but doing so may erode the exclusivity premium paid by the first buyer. This creates the central commitment problem: after selling to one firm, the seller is tempted to license to the rival as well, thereby reducing the first buyer's willingness to pay ex ante.

Social Planner. A benevolent social planner evaluates total welfare as

$$W = \Pi_S + \Pi_1 + \Pi_2 + CS,$$

where customer surplus is given by

$$CS = \sum_{i \in \{1,2\}} \left[(A_i - p_i^c) \frac{1-s}{2} + (A_i - p_i^s) s D_i^s \right].$$

This allows us to assess how exclusivity affects welfare through both prices and investment incentives.

Timing. Information is complete. The game unfolds in four stages:

1. **Data Pricing.** The data seller S sets a publicly observed license fee f ;
2. **Data Collection.** Each firm F_i chooses $(\ell_i, r_i) \in \{0, 1\}^2$, which are publicly observed;

¹⁵Allowing S to set different fees for F_1 and F_2 would not affect the analysis as far as exclusivity is concerned.

3. **Prediction Pricing.** Firms simultaneously set (p_i^c, p_i^s) . Customers choose the firm offering higher utility,
4. **Realization.** Nature draws the data and the target variable, prediction losses and profits are realized for S , F_1 , and F_2 .

We solve the game by backward induction and characterize the Subgame Perfect Equilibrium (SPE).

6.3.2 Prediction Pricing

We analyze the pricing subgame, taking prediction accuracies (A_1, A_2) as given from the data collection. Firms can price-discriminate between captive customers and shoppers.

Captive pricing. Let the outside option yield utility 0. Assuming they are expected utility maximizers, captive customers of F_i buy if and only if $A_i - p_i \geq 0$. Therefore, the unique optimal price is

$$p_i^c = A_i \quad \text{for each } i \in \{1, 2\}. \quad (4)$$

The firm extracts the full surplus of captives.

Shopper pricing. Shoppers buy from the firm that offers the higher net utility; if $A_i - p_i > A_j - p_j^s$ they all buy from F_i , if $A_i - p_i < A_j - p_j^s$ they all buy from F_j , and if equal any split yields the same payoff for firms and customers.¹⁶ Write $(x)^+ \equiv \max\{x, 0\}$. The next lemma characterizes the unique trembling-hand perfect equilibrium for shoppers.

Lemma 7 (Shopper-price equilibrium). *Fix (A_1, A_2) . The pricing subgame in shopp[Shopper-price equilibrium]er prices admits a unique equilibrium under trembling-hand perfection:*

$$p_i^s = (A_i - A_j)^+, \quad i \neq j, \quad (5)$$

so that (i) if $A_i > A_j$, firm i sets $p_i^s = A_i - A_j$, firm j sets $p_j^s = 0$, and all shoppers buy from i ; (ii) if $A_i = A_j$, both set $p_i^s = p_j^s = 0$ and shoppers can be split arbitrarily.

Combining Equation (4) and Lemma 7, the equilibrium profit of firm i given data strategies (ℓ_i, r_i) and (ℓ_j, r_j) (which pin down (A_i, A_j)) is

$$\pi(\ell_i, r_i; \ell_j, r_j) = \frac{1-s}{2} A_i + s(A_i - A_j)^+, \quad j \neq i. \quad (6)$$

The first term is captive revenue at the value price; the second term is shopper revenue, which equals the quality advantage when the firm is better and zero otherwise.

¹⁶This tie-breaking rule is without loss for the equilibrium characterization; a trembling-hand refinement will select the limit outcome stated below.

6.3.3 Cost Regions

The trade-off between exclusivity and non-exclusivity depends on the equilibrium of the proprietary data investment subgame, which in turn depends on the cost of collecting proprietary data c relative to the accuracy gain A . Two regimes arise:

- **Low-cost region:** $c \in [0, \bar{c})$, where

$$\bar{c} \equiv \frac{1-s}{2}A.$$

Here proprietary data are cheap enough that any licensed firm will always invest.

- **High-cost region:** $c \in [\bar{c}, \bar{\bar{c}})$, where

$$\bar{\bar{c}} \equiv \frac{1+s}{2}A.$$

Here investment occurs only when a firm obtains a monopoly over shoppers, leading to mixed strategies when both license.

- **No investment region:** $c \geq \bar{\bar{c}}$. Here investment never occurs.

These regions generate two distinct exclusivity incentives for the seller and two distinct welfare benchmarks for the planner, which we will analyze separately.

6.3.4 Low-cost Region

Proprietary Data Investment. If $c < \bar{c}$, the cost of investing is less than the increase in revenue on a firm's captive customers. Therefore investment will occur if and only if the license is purchased.

Proposition 19. *If $c < \bar{c}$, for any licensing choice (ℓ_1, ℓ_2) there is a unique Nash Equilibrium*

$$(r_1^*, r_2^*) = (\ell_1, \ell_2).$$

The intuition is straightforward: the rival's competition in the shopper segment is irrelevant because the incentives to invest in the captive segment are sufficiently large relative to the cost of investing. Furthermore, because investing in proprietary data is costly but gives no benefit if the license is acquired, in case a firm does not license it will also avoid investing.

Data Pricing. The seller chooses these fees to maximize its profit, taking into account that the fees determine how many firms choose to license the data. The following result characterizes the maximal fee f_d compatible with $d \in \{1, 2\}$ firms purchasing the data.

Lemma 8. *If $c < \bar{c}$,*

$$f_2 = \frac{1-s}{2}\bar{A} - c, \quad \text{and} \quad f_1 = \frac{1+s}{2}\bar{A} - c.$$

Intuitively, if both firm license the training data they will both invest and will only be able to extract the surplus of the captive customers $\frac{1-s}{2}\bar{A}$. Instead, if only one licenses, it will extract the full surplus of the shoppers as well, giving it an additional demand of sA , because by Proposition 19 it will invest, and since $f_1 > f_2$ the rival will not license as it is only profitable for a single firm to license.

If the profit from selling to one firm is higher than selling to both, i.e., $f_1 > 2f_2$, the data seller will pursue *de facto exclusivity*, setting a price so high hat only one firm will choose to purchase the data and the other will avoid entry.

Proposition 20. *The data seller will choose $f^* = f_1$, i.e., it will pursue de facto exclusivity, if $c \in [\underline{c}, \bar{c})$, where*

$$\underline{c} \equiv \frac{1-3s}{2}\bar{A} \in [0, \bar{c}].$$

Therefore, whenever competition is fierce, i.e. $s > 1/3$, the seller will always prefer de facto exclusivity. If $s \leq 1/3$, the seller will prefer de facto exclusivity if the cost of collecting prediction covariates is large compared to \bar{A} . Since $\bar{A} = A(n, k)$, which is increasing in both arguments, exclusivity is profitable if datasets are sufficiently small.

Social Welfare. I assume the planner can choose between de facto exclusivity and real non exclusivity but cannot influence the proprietary data investment equilibrium.

Proposition 21. *If $c < \bar{c}$, the planner prefers real non-exclusivity to de facto exclusivity.*

The cost is too low for the motive of avoiding excessive entry and overinvestment in data collection ot make it socially desirable to pursue de facto exclusivity. The following result characterizes an interval in which de facto excusivity is profitable for the data seller even if it harms social welfare.

Corollary 8. *The seller chooses de facto exclusivity even though it is harmful if and only if $c \in (\underline{c}, \bar{c})$.*

The length of the interval is $\bar{c} - \underline{c} = s\bar{A}$, so harmful exclusivity is more likely if competition is fierce and, given $\bar{A} \equiv A(n, k)$, which is increasing in both arguments, if data is abundant.

6.3.5 High-cost Region

I now consider what happens when the cost of proprietary data collection is $c \in (\bar{c}, \bar{\bar{c}})$. This means the cost of investing is more than the surplus of a firm's captive customers but possibly less than the cumulative surplus of a firm's captive customers and the shoppers.

Proprietary Data Investment. If $c \in (\bar{c}, \bar{\bar{c}})$, a firm will collect prediction covariates with probability one only if it can sell as a monopolist to the shoppers, meaning that it benefits from de facto exclusivity. If both firms purchase the license, there is no pure-strategy equilibrium in the proprietary data-collection subgame.

Proposition 22. *We distinguish two subcases:*

- **De facto exclusivity:** *Without loss of generality, let $\ell_1 = 1$, $\ell_2 = 0$, the unique NE is*

$$r_1 = 1, \quad r_2 = 0.$$

- **Both firms license:** $\ell_1 = \ell_2 = 1$. *The unique symmetric mixed strategy equilibrium has each firm investing with probability*

$$\xi^* = \frac{A(1+s) - 2c}{2As}, \quad \xi^* \in [0, 1].$$

In equilibrium, profits are zero: $\Pi_i^{\text{mix}} = 0$.

The mixed strategy equilibrium investment probability is increasing in s if and only if $c > \bar{A}/2$ (i.e., when proprietary data are relatively costly) and decreasing in s when $c < \bar{A}/2$.

Data Pricing. As the profit of the mixed strategy equilibrium is zero, firms will have no willingness to pay for the license unless there is de facto exclusivity. Under de facto exclusivity, only one firm licenses and invests, while the rival neither licenses nor invests.

Proposition 2. *If $c \in (\bar{c}, \bar{\bar{c}})$, the data seller will always choose $f^* = f_1$, i.e., it will always pursue de facto exclusivity, where*

$$f_1 = \frac{1+s}{2}A_2 - c.$$

Social welfare. The following result characterizes the social planner optimum, by comparing the pure strategy equilibrium of the investment subgame and the mixed strategy non-exclusive one.

Proposition 23. *If $c \in (\bar{c}, \bar{\bar{c}})$, the planner prefers de facto exclusivity to real non-exclusivity.*

The overall investment in the real non exclusive case is so low that the planner prefers to grant exclusivity in order to ensure that investment occurs.

6.3.6 Summary.

6.3.7 Summary

The equilibrium interaction between exclusivity, investment, and welfare is governed by the fixed cost of proprietary data collection c , relative to the accuracy premium

$$A \equiv A(k, n) = \frac{k}{\frac{1-k}{n} + 1}.$$

Three cost thresholds determine firms' willingness to collect proprietary covariates:

$$\underline{c} \equiv \frac{1-3s}{2}A, \quad \bar{c} \equiv \frac{1-s}{2}A, \quad \bar{\bar{c}} \equiv \frac{1+s}{2}A.$$

For fixed (c, s) we can invert these equations and obtain three endogenous thresholds in covariate richness,

$$\underline{k}(n, c, s), \quad \bar{k}(n, c, s), \quad \bar{\bar{k}}(n, c, s),$$

defined implicitly by

$$c = \frac{1-3s}{2}A(n, \underline{k}), \quad c = \frac{1-s}{2}A(n, \bar{k}), \quad c = \frac{1+s}{2}A(n, \bar{\bar{k}}).$$

For $s < 1/3$ these satisfy

$$\underline{k}(n, c, s) < \bar{k}(n, c, s) < \bar{\bar{k}}(n, c, s) \quad \text{for all } (n, c).$$

These thresholds partition the (n, k) -space into four regions:

- **No investment:** If $k < \underline{k}(n, c, s)$, proprietary data collection is never profitable, even under exclusivity. No firm invests, and licensing decisions have no effect.
- **Beneficial exclusivity:** If $k \in [\underline{k}(n, c, s), \bar{k}(n, c, s))$, investment occurs only under exclusivity, because a firm collects proprietary data if and only if it can monopolize shoppers. Both the seller and the planner prefer exclusivity, as it ensures that at least one firm invests.
- **Harmful exclusivity:** If $k \in [\bar{k}(n, c, s), \bar{\bar{k}}(n, c, s))$, both firms would invest even without exclusivity, but the seller prefers to restrict access to extract a higher fee. The planner prefers non-exclusivity to keep both firms active. Exclusivity is privately profitable but socially harmful.
- **No exclusivity:** If $k \geq \bar{\bar{k}}(n, c, s)$, proprietary data are sufficiently cheap relative to dataset quality that both firms always invest. Both the seller and the planner prefer non-exclusivity, so broad data access is jointly optimal.

The figure below plots these three thresholds and the resulting four regions for a fixed pair (c, s) with $s < 1/3$.

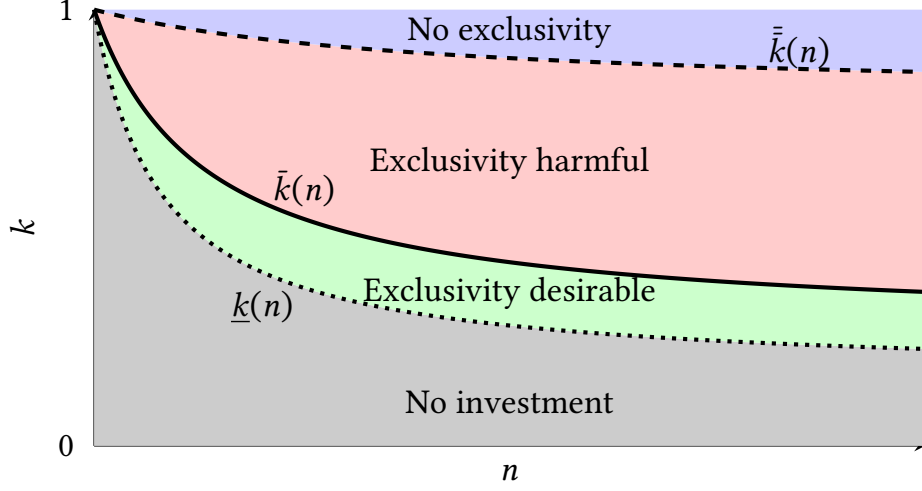


Figure 5: The three thresholds $k = \underline{k}(n)$, $k = \bar{k}(n)$, and $k = \bar{\bar{k}}(n)$ partition the (n, k) -space into four regions. Below $\underline{k}(n)$, investment never occurs. Between $\underline{k}(n)$ and $\bar{k}(n)$, exclusivity is jointly optimal because it induces investment. Between $\bar{k}(n)$ and $\bar{\bar{k}}(n)$, exclusivity is privately profitable but socially harmful. Above $\bar{\bar{k}}(n)$, both the seller and the planner prefer non-exclusivity.

7 Conclusion

This paper develops a general framework for understanding the value of data in prediction by explicitly modeling covariates. The analysis shows how economies of scope across covariates, interactions between covariates and observations, and complementarities between training and prediction can generate increasing returns, offering a microfoundation for the rich-get-richer effects often observed in data-driven markets.

These forces have direct implications for policy and strategy. Prediction technologies may display natural monopsony characteristics, as concentrating covariates within one firm can raise efficiency. Privacy regulation that fragments data supply may inadvertently reinforce monopsony power, creating a trilemma between privacy, competition, and efficiency. The framework also highlights that not all data pooling agreements are alike: pooling lists of users with the same covariates can be anticompetitive, whereas pooling different covariates on similar users raise welfare by eliminating double marginalization. Exclusivity deals, such as those signed between AI labs and data providers, may profitably foreclose entry by depriving rivals of essential complements. For firms, the results imply that prediction entails substantial sunk costs: early on, investment should balance user acquisition and attribute enrichment, while specialization and integration become optimal at a larger scale.

More broadly, the analysis cautions against treating data as homogeneous. Policies promoting open data without regard to dataset composition may miss crucial efficiency margins, whereas access remedies such as FRAND-priced APIs or federated learning preserve economies of scope.

My work opens two natural avenues for future research. The first is empirical. I aim to develop a methodology to test my results on real datasets. While the existing empirical literature¹⁷ provides partial support to my findings, it suffers from two limitations: (i) most studies focus on a single dataset, whereas uncovering general properties requires comparing multiple datasets along common dimensions; and (ii) no existing work systematically tests all the properties identified in my model. Once these empirical properties are validated, my framework could serve as the foundation for a practical formula for data valuation, in the spirit of the Black–Scholes–Merton formula for derivatives.¹⁸ The second avenue is theoretical. Embedding my static model into a dynamic Wald sampling framework would allow me to microfound data-enabled learning and analyze when feedback loops generate convergent data-collection strategies versus when they diverge.

Finally, the framework invites a broader research agenda: in his seminal critique of central planning, Hayek (1945) emphasized that “knowledge... never exists in concentrated form but solely as the dispersed bits... which all the separate individuals possess”. Today, users’ online activity transforms such dispersed knowledge into datasets that can be centralized, recombined, and monetized. My analysis shows that statistical properties of prediction create intrinsic incentives for such concentration. The concentration of data in servers controlled by a few large firms raises a broader question: do prediction algorithms substitute for, or complement, the market mechanism? Is the rise of data the panacea to market failures deriving from asymmetric information and search frictions, or is it the first step to the fall of the market? I leave this foundational question open to future research.

¹⁷See Bajari et al. (2019), Schaefer and Sapi (2023), Lee and Wright (2023), Yoganarasimhan (2020), and Carballa-Smichowski et al. (2025b)

¹⁸See Black and Scholes (1973), Merton (1973)

References

- Acemoglu, Daron et al. (2022). “Too much data: Prices and inefficiencies in data markets.” In: *American Economic Journal: Microeconomics* 14.4, pp. 218–256.
- Aghion, Philippe and Patrick Bolton (1987). “Contracts as a Barrier to Entry.” In: *The American Economic Review* 77.3, pp. 388–401.
- Allcott, Hunt et al. (2025). *Sources of market power in web search: Evidence from a field experiment*. Tech. rep. National Bureau of Economic Research.
- Aral, Sinan, Erik Brynjolfsson, and DJ Wu (2008). “Which came first, IT or productivity? The virtuous cycle of investment and use in enterprise systems.” In:
- Bajari, Patrick et al. (2019). “The impact of big data on firm performance: An empirical investigation.” In: *AEA papers and proceedings* 109, pp. 33–37.
- Belkin, Mikhail et al. (2019). “Reconciling modern machine-learning practice and the classical bias–variance trade-off.” In: *Proceedings of the National Academy of Sciences* 116.32, pp. 15849–15854.
- Bergemann, Dirk and Alessandro Bonatti (2024). “Data, Competition, and Digital Platforms.” In: *American Economic Review* 114.8, pp. 2553–2595.
- Bergemann, Dirk, Alessandro Bonatti, and Tan Gan (2022). “The economics of social data.” In: *The RAND Journal of Economics* 53.2, pp. 263–296.
- Bergemann, Dirk, Alessandro Bonatti, and Alex Smolin (2018). “The Design and Price of Information.” In: *American Economic Review* 108.1, pp. 1–48.
- Berger, James O. (1990). *Statistical decision theory*. Springer, pp. 277–284.
- Black, Fischer and Myron Scholes (1973). “The Pricing of Options and Corporate Liabilities.” In: *Journal of Political Economy* 81.3, pp. 637–654.
- Brier, Glenn W. (1950). “Verification of Forecasts Expressed in Terms of Probability.” In: *Monthly Weather Review* 78.1, pp. 1–3.
- Calzolari, Giacomo, Anatole Cheysson, and Riccardo Rovatti (2025). “Machine data: market and analytics.” In: *Management Science*.
- Carballa-Smichowski, Bruno et al. (2025a). *Data Sharing or Analytics Sharing?* TSE Working Paper 25-1615. Toulouse School of Economics.
- Carballa-Smichowski, Bruno et al. (2025b). “Economies of scope in data aggregation: Evidence from health data.” In: *Information Economics and Policy* 71, p. 101146.
- Cong, Lin William, Zhiguo He, and Changhua Yu (2021). “Data as Capital.” In: *Review of Financial Studies* 34.6, pp. 2895–2936.
- Cornière, Alexandre de and Greg Taylor (2024). “Data-Driven Mergers.” In: *Management Science* 70.9, pp. 6473–6482.
- Cournot, Antoine Augustin (1838). *Recherches sur les principes mathématiques de la théorie des richesses*. English translation: *Researches into the Mathematical Principles of the Theory of Wealth*, Macmillan, 1897. Paris: Hachette.

- Dasaratha, Krishna, Juan Ortner, and Chengyang Zhu (2025). "Markets for Models." In: *arXiv preprint arXiv:2503.02946*.
- De Corniere, Alexandre and Greg Taylor (2025). "Data and Competition: A Simple Framework." In: *Forthcoming, RAND Journal of Economics*.
- DeGroot, Morris H. (2005). *Optimal statistical decisions*. John Wiley & Sons.
- Digital Platforms, Stigler Committee on (2019). *Final Report*. Stigler Center for the Study of the Economy and the State, University of Chicago Booth School of Business.
- Farboodi, Maryam and Laura Veldkamp (2025). *A model of the Data Economy*. Tech. rep. R&R, Review of Economic Studies.
- Goldfarb, Avi and Catherine Tucker (2011). "Privacy Regulation and Online Advertising." In: *Management Science* 57.1, pp. 57–71.
- Gu, Yiquan, Leonardo Madio, and Carlo Reggiani (Sept. 2021). "Data brokers co-opetition." In: *Oxford Economic Papers* 74.3, pp. 820–839.
- Hagiu, Andrei and Julian Wright (2023). "Data-enabled learning, network effects, and competitive advantage." In: *The RAND Journal of Economics* 54.4, pp. 638–667.
- Hastie, Trevor et al. (2020). "Surprises in High-Dimensional Ridgeless Least Squares Interpolation." In.
- Hayek, Friedrich A. (1945). "The Use of Knowledge in Society." In: *American Economic Review* 35. Reprinted in F.A. Hayek (ed.), *Individualism and Economic Order*. London: Routledge and Kegan Paul, pp. 519–530.
- Iyer, Ganesh and Tianshu Ke (2024). "Competition and Algorithmic Complexity in Predictive Analytics." In: *Marketing Science* 43.2, pp. 215–233.
- Jacot, Arthur, Franck Gabriel, and Clement Hongler (2018). "Neural Tangent Kernel: Convergence and Generalization in Neural Networks." In: *Advances in Neural Information Processing Systems*. Ed. by S. Bengio et al. Vol. 31. Curran Associates, Inc.
- Jolliffe, I. T. (2002). *Principal Component Analysis*. Springer Series in Statistics. New York: Springer-Verlag.
- Jones, Charles I. and Christopher Tonetti (2020). "Nonrivalry and the Economics of Data." In: *American Economic Review* 110.9, 2819–58.
- Kaplan, Jared et al. (2020). "Scaling laws for neural language models." In: *arXiv preprint arXiv:2001.08361*.
- Katz, Michael L. and Carl Shapiro (Aug. 1986). "How to License Intangible Property*." In: *The Quarterly Journal of Economics* 101.3, pp. 567–589.
- Lee, Gunhaeng and Julian Wright (2023). "Recommender systems and the Value of User Data." In: *National University of Singapore Working Paper*.
- Lerner, Josh and Jean Tirole (2004). "Efficient Patent Pools." In: *American Economic Review* 94.3, 691–711.
- Lewis, Patrick et al. (2020). "Retrieval-augmented generation for knowledge-intensive nlp tasks." In: *Advances in neural information processing systems* 33, pp. 9459–9474.

- Lindley, D. V. and A. F. M. Smith (1972). “Bayes Estimates for the Linear Model.” In: *Journal of the Royal Statistical Society: Series B (Methodological)* 34.1, pp. 1–18.
- Liu, Nelson F et al. (2023). “Lost in the middle: How language models use long contexts.” In: *arXiv preprint arXiv:2307.03172*.
- MacKay, David J. C. (May 1992). “Bayesian Interpolation.” In: *Neural Computation* 4.3, pp. 415–447.
- Merton, Robert C. (1973). “Theory of Rational Option Pricing.” In: *The Bell Journal of Economics and Management Science* 4.1, pp. 141–183.
- Montiel Olea, José Luis et al. (Apr. 2022). “Competing Models.” In: *The Quarterly Journal of Economics* 137.4, 2419–2457.
- Nakkiran, Preetum et al. (2021). “Deep double descent: Where bigger models and more data hurt.” In: *Journal of Statistical Mechanics: Theory and Experiment* 2021.12, p. 124003.
- Nocke, Volker, Martin Peitz, and Konrad Stahl (Dec. 2007). “Platform Ownership.” In: *Journal of the European Economic Association* 5.6, pp. 1130–1160.
- Prüfer, Jens and Christoph Schottmüller (2021). “Competing with big data.” In: *The Journal of Industrial Economics* 69.4, pp. 967–1008.
- Radner, Roy and Joseph Stiglitz (1984). “A Nonconcavity in the Value of Information.” In: *Bayesian models in economic theory* 5, pp. 33–52.
- Ricardo, D. (1817). *On the Principles of Political Economy and Taxation*. John Murray.
- Schaefer, Maximilian (2025). *When Should we Expect Non-Decreasing Returns from Data in Prediction Tasks?*
- Schaefer, Maximilian and Geza Sapi (2023). “Complementarities in learning from data: Insights from general search.” In: *Information Economics and Policy* 65, p. 101063.
- Schmidt, Eric (Oct. 2, 2009). *How Google Plans to Stay Ahead in Search*. Accessed: 2025-11-12. Bloomberg. URL: <https://www.bloomberg.com/news/articles/2009-10-02/how-google-plans-to-stay-ahead-in-search>.
- Selten, Reinhard (1998). “Axiomatic Characterization of the Quadratic Scoring Rule.” In: *Experimental Economics* 1.1, pp. 43–62.
- Strzalecki, Tomasz (2024). *Variational Bayes and non-Bayesian Updating*.
- UK Competition and Markets Authority (2019). *Unlocking Digital Competition: Report of the Digital Competition Expert Panel*. London: Competition and Markets Authority.
- Varian, Hal (2018). “Artificial Intelligence, Economics, and Industrial Organization.” In: *The Economics of Artificial Intelligence: An Agenda*. NBER Chapters. National Bureau of Economic Research, Inc, pp. 399–419.
- Vershynin, Roman (2018). *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge: Cambridge University Press.
- Wilson, Robert (1975). “Informational Economies of Scale.” In: *Bell Journal of Economics* 6.1, pp. 184–195.

Yoganarasimhan, Hema (2020). “Search personalization using machine learning.” In: *Management Science* 66.3, pp. 1045–1070.

A Proofs

Lemma 1 (Optimal Predictor). *The optimal predictor is*

$$\hat{y}^*(D) = \mathbb{E}[y \mid D] = \mathbf{x}'_P \mathbb{E}[\boldsymbol{\beta}_P \mid (\mathbf{y}, X_{\mathcal{T}})].$$

Proof. Under squared loss, the Bayes optimal predictor is the conditional mean:

$$\hat{y}^*(D_{\mathcal{T}, P}) = \mathbb{E}[y \mid (\mathbf{y}, X), \mathbf{x}_P].$$

Write $y = \sum_{k \in \mathcal{K}} \beta_k x_k$. By the law of iterated expectations and independence of \mathbf{x} and $\boldsymbol{\beta}$,

$$\mathbb{E}[y \mid (\mathbf{y}, X), \mathbf{x}_P] = \sum_{k \in \mathcal{K}} \mathbb{E}[\beta_k x_k \mid (\mathbf{y}, X), \mathbf{x}_P] = \sum_{k \in \mathcal{P}} x_k \mathbb{E}[\beta_k \mid (\mathbf{y}, X)] + \sum_{k \notin \mathcal{P}} \mathbb{E}[\beta_k x_k \mid (\mathbf{y}, X), \mathbf{x}_P].$$

For $k \notin \mathcal{P}$, x_k is mean zero and independent of $((\mathbf{y}, X), \mathbf{x}_P, \boldsymbol{\beta})$, so $\mathbb{E}[\beta_k x_k \mid (\mathbf{y}, X), \mathbf{x}_P] = 0$. Thus

$$\mathbb{E}[y \mid (\mathbf{y}, X), \mathbf{x}_P] = \sum_{k \in \mathcal{P}} x_k \mathbb{E}[\beta_k \mid (\mathbf{y}, X)].$$

If $k \notin \mathcal{T}$, then β_k is not updated by (\mathbf{y}, X) and $\mathbb{E}[\beta_k \mid (\mathbf{y}, X)] = \mathbb{E}[\beta_k] = 0$. Hence, only indices in $\mathcal{T} \cap \mathcal{P}$ contribute, giving

$$\mathbb{E}[y \mid (\mathbf{y}, X), \mathbf{x}_P] = \mathbf{x}'_P \mathbb{E}[\boldsymbol{\beta}_P \mid (\mathbf{y}, X)],$$

which proves the claim. □

Lemma 3. *For any nested sequence $(\mathcal{T}_m)_{m \geq 1}$ with $\mathcal{T}_m \subset \mathcal{T}_{m+1}$ and $\bigcup_m \mathcal{T}_m = \mathcal{T}$,*

$$\sum_{k \in \mathbb{N} \setminus \mathcal{T}_m} \sigma_k^2 \beta_k^2 \xrightarrow[m \rightarrow \infty]{a.s.} \text{Var}(\boldsymbol{\varepsilon}_{\mathcal{T}}) = 1 - S(\mathcal{T}_m). \quad (3)$$

?proofname? Define $Y_k \equiv \sigma_k^2(\beta_k^2 - 1)$. Then $\mathbb{E}[Y_k] = 0$, the Y_k are independent, and

$$\text{Var}(Y_k) = \sigma_k^4 \text{Var}(\beta_k^2) \leq C \sigma_k^4 \quad \text{for some finite } C > 0$$

(e.g., $C = 2$ when $\beta_k \sim \mathcal{N}(0, 1)$). Because $0 \leq \sigma_k^2 \leq 1$ and $\sum_k \sigma_k^2 = 1$, we have $\sum_k \sigma_k^4 \leq \sum_k \sigma_k^2 = 1 < \infty$, hence

$$\sum_{k=1}^{\infty} \text{Var}(Y_k) \leq C \sum_{k=1}^{\infty} \sigma_k^4 < \infty.$$

By Kolmogorov's convergence theorem (a special case of the three-series theorem), the series $\sum_{k=1}^{\infty} Y_k$ converges almost surely. Therefore its *tails* vanish almost surely along any nested

complements:

$$R_m \equiv \sum_{k \in \mathcal{T}_m^c} Y_k = \sum_{k \in \mathcal{T}_m^c} \sigma_k^2 (\beta_k^2 - 1) \xrightarrow[m \rightarrow \infty]{a.s.} 0.$$

Now decompose

$$\sum_{k \in \mathcal{T}_m^c} \sigma_k^2 \beta_k^2 = \sum_{k \in \mathcal{T}_m^c} \sigma_k^2 + \sum_{k \in \mathcal{T}_m^c} \sigma_k^2 (\beta_k^2 - 1) = (1 - S(\mathcal{T}_m)) + R_m,$$

and use $R_m \rightarrow 0$ a.s. to conclude. Since $S(\mathcal{T}_m) \uparrow 1$, we also have $\sum_{k \in \mathcal{T}_m^c} \sigma_k^2 \beta_k^2 \xrightarrow{a.s.} 0$ and, when $1 - S(\mathcal{T}_m) > 0$,

$$\frac{\sum_{k \in \mathcal{T}_m^c} \sigma_k^2 \beta_k^2}{1 - S(\mathcal{T}_m)} \xrightarrow{a.s.} 1.$$

□

Proposition 1 (Posterior Distribution). *The posterior distribution of the parameter vector is given by:*

- For untrained parameters,

$$\boldsymbol{\beta}_{\mathcal{T}^c} \mid (\mathbf{y}, \mathbf{X}_{\mathcal{T}}) \sim \mathcal{N}(\mathbf{0}_{|\mathcal{T}^c|}, \mathbf{I}_{|\mathcal{T}^c|});$$

- For trained parameters ,

$$\boldsymbol{\beta}_{\mathcal{T}} \mid (\mathbf{y}, \mathbf{X}_{\mathcal{T}}) \sim \mathcal{N}\left(\underbrace{(X'_{\mathcal{T}} X_{\mathcal{T}} + n\lambda \cdot \mathbf{I}_{|\mathcal{T}|})^{-1} X'_{\mathcal{T}} \mathbf{y}}_{\equiv \mathbb{E}[\boldsymbol{\beta}_{\mathcal{T}} \mid (\mathbf{y}, \mathbf{X}_{\mathcal{T}})]}, \underbrace{\left(\frac{1}{n\lambda} X'_{\mathcal{T}} X_{\mathcal{T}} + \mathbf{I}_t\right)^{-1}}_{\equiv \text{Var}[\boldsymbol{\beta}_{\mathcal{T}} \mid (\mathbf{y}, \mathbf{X}_{\mathcal{T}})]}\right),$$

with $\boldsymbol{\beta}_{\mathcal{T}}$ independent from $\boldsymbol{\beta}_{\mathcal{T}^c}$.

?proofname? Because the prior is $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{|\mathbb{N}|})$ and only $\boldsymbol{\beta}_{\mathcal{T}}$ enters the likelihood, the joint posterior factorizes as

$$p(\boldsymbol{\beta} \mid (\mathbf{y}, \mathbf{X})) = p(\boldsymbol{\beta}_{\mathcal{T}} \mid (\mathbf{y}, \mathbf{X})) p(\boldsymbol{\beta}_{\mathcal{T}^c} \mid (\mathbf{y}, \mathbf{X})),$$

with $p(\boldsymbol{\beta}_{\mathcal{T}^c} \mid (\mathbf{y}, \mathbf{X})) = \mathcal{N}(\mathbf{0}, \mathbf{I}_{|\mathcal{T}^c|})$ since $\boldsymbol{\beta}_{\mathcal{T}^c}$ does not appear in the likelihood and the prior mean is zero. This proves part the second result.

For the first result, the likelihood is

$$\mathbf{y} \mid \boldsymbol{\beta}_{\mathcal{T}} \sim \mathcal{N}(X_{\mathcal{T}} \boldsymbol{\beta}_{\mathcal{T}}, u_{\mathcal{T}}^2 \mathbf{I}_n), \quad u_{\mathcal{T}}^2 \equiv 1 - S(\mathcal{T}),$$

and the prior is $\boldsymbol{\beta}_{\mathcal{T}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_t)$. Standard derivations in DeGroot (2005) or Berger (1990) show that by conjugacy (or completing the square), the posterior is Gaussian

$$\boldsymbol{\beta}_{\mathcal{T}} \mid (\mathbf{y}, \mathbf{X}) \sim \mathcal{N}(\boldsymbol{\mu}_{\mathcal{T}}, \boldsymbol{\Sigma}_{\mathcal{T}}),$$

with precision

$$\Sigma_{\mathcal{T}}^{-1} = \frac{1}{u_{\mathcal{T}}^2} \mathbf{X}'_{\mathcal{T}} \mathbf{X}_{\mathcal{T}} + \mathbf{I}_t,$$

and mean

$$\boldsymbol{\mu}_{\mathcal{T}} = \Sigma_{\mathcal{T}} \left(\frac{1}{u_{\mathcal{T}}^2} \mathbf{X}'_{\mathcal{T}} \mathbf{y} \right) = (\mathbf{X}'_{\mathcal{T}} \mathbf{X}_{\mathcal{T}} + u_{\mathcal{T}}^2 \mathbf{I}_t)^{-1} \mathbf{X}'_{\mathcal{T}} \mathbf{y}.$$

Therefore, the Bayes estimator (posterior mean) equals the stated ridge form, which proves the first result. \square

Proposition 14 (Ridge Estimator Interpretation). *The ex-ante optimal regularization parameter is equivalent to the misspecification penalty*

$$\xi^*(n, \mathcal{T}) = \lambda(n, \mathcal{T}) \equiv \frac{1 - S(\mathcal{T})}{n}.$$

?proofname? The ridge objective is strictly convex; its unique minimizer solves the first-order condition:

$$\frac{2}{n} \mathbf{X}'_{\mathcal{T}} (\mathbf{X}_{\mathcal{T}} \hat{\mathbf{b}} - \mathbf{y}) + 2\lambda \hat{\mathbf{b}} = \mathbf{0}.$$

Hence

$$\left(\frac{1}{n} \mathbf{X}'_{\mathcal{T}} \mathbf{X}_{\mathcal{T}} + \lambda \mathbf{I}_t \right) \hat{\mathbf{b}} = \frac{1}{n} \mathbf{X}'_{\mathcal{T}} \mathbf{y}, \quad \text{so} \quad \hat{\mathbf{b}} = (\mathbf{X}'_{\mathcal{T}} \mathbf{X}_{\mathcal{T}} + n\lambda \mathbf{I}_t)^{-1} \mathbf{X}'_{\mathcal{T}} \mathbf{y}.$$

Comparing with the posterior mean from Proposition 1,

$$\mathbb{E}[\boldsymbol{\beta}_{\mathcal{T}} \mid (\mathbf{y}, \mathbf{X})] = (\mathbf{X}'_{\mathcal{T}} \mathbf{X}_{\mathcal{T}} + u_{\mathcal{T}}^2 \mathbf{I}_t)^{-1} \mathbf{X}'_{\mathcal{T}} \mathbf{y},$$

we obtain equality when $n\lambda = u_{\mathcal{T}}^2$, i.e. $\lambda = \frac{u_{\mathcal{T}}^2}{n} = \frac{1-S(\mathcal{T})}{n}$. This proves the claim. \square

Lemma 9 (Optimal regularization under a random effects prior). *Consider the linear model*

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, u^2 \mathbf{I}_n),$$

where $\mathbf{X} \in \mathbb{R}^{n \times t}$ is of full column rank and

$$\frac{1}{n} \mathbf{X}' \mathbf{X} = \mathbf{V} \text{diag}(g_1, \dots, g_t) \mathbf{V}'.$$

Suppose the coefficients follow a random-effects prior

$$\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I}_t).$$

Then, under squared-error loss, the optimal ridge regularization parameter that minimizes the expected mean-squared error

$$\mathbb{E} \|\hat{\boldsymbol{\beta}}_{\lambda} - \boldsymbol{\beta}\|^2$$

is given by

$$\lambda^* = \frac{u^2}{n\tau^2}.$$

?proofname? The ridge estimator based on the penalized least squares criterion

$$\hat{\boldsymbol{\beta}}_\lambda = \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|^2 \right\}$$

has closed-form solution

$$\hat{\boldsymbol{\beta}}_\lambda = (\mathbf{G} + \lambda \mathbf{I}_t)^{-1} \frac{1}{n} \mathbf{X}' \mathbf{y}, \quad \mathbf{G} \equiv \frac{1}{n} \mathbf{X}' \mathbf{X} = \mathbf{V} \text{diag}(g_1, \dots, g_t) \mathbf{V}'.$$

In the eigenbasis of \mathbf{G} , let $\tilde{\boldsymbol{\beta}} = \mathbf{V}' \boldsymbol{\beta}$. Using the properties of the Gaussian prior and noise, the posterior mean satisfies

$$\mathbb{E}[\hat{\boldsymbol{\beta}}_\lambda \mid \boldsymbol{\beta}] = (\mathbf{G} + \lambda \mathbf{I}_t)^{-1} \mathbf{G} \boldsymbol{\beta}.$$

The mean-squared error can be decomposed as

$$\mathbb{E} \|\hat{\boldsymbol{\beta}}_\lambda - \boldsymbol{\beta}\|^2 = \sum_{j \in \mathcal{T}} \left[\left(\frac{\lambda}{g_j + \lambda} \right)^2 \mathbb{E} [\tilde{\beta}_j^2] + \frac{u^2}{n} \frac{g_j}{(g_j + \lambda)^2} \right].$$

Under the random-effects prior $\tilde{\beta}_j \sim \mathcal{N}(0, \tau^2)$, $\mathbb{E} \tilde{\beta}_j^2 = \tau^2$ for all j , so

$$\mathbb{E} \|\hat{\boldsymbol{\beta}}_\lambda - \boldsymbol{\beta}\|^2 = \sum_{j \in \mathcal{T}} \frac{\tau^2 \lambda^2 + \frac{u^2}{n} g_j}{(g_j + \lambda)^2}.$$

Differentiating with respect to λ gives

$$\frac{\partial}{\partial \lambda} \mathbb{E} \|\hat{\boldsymbol{\beta}}_\lambda - \boldsymbol{\beta}\|^2 = 2 \sum_{j \in \mathcal{T}} \frac{g_j \left(\tau^2 \lambda - \frac{u^2}{n} \right)}{(g_j + \lambda)^3}.$$

Setting the derivative to zero yields the first-order condition

$$\sum_{j \in \mathcal{T}} \frac{g_j \left(\tau^2 \lambda^* - \frac{u^2}{n} \right)}{(g_j + \lambda^*)^3} = 0.$$

Since all weights $g_j/(g_j + \lambda^*)^3 > 0$, this condition is satisfied if and only if

$$\tau^2 \lambda^* - \frac{u^2}{n} = 0,$$

which gives

$$\lambda^* = \frac{u^2}{n\tau^2}.$$

□

Lemma 2. *The value of a dataset design $(n, \mathcal{T}, \mathcal{P})$ is*

$$V(n, \mathcal{T}, \mathcal{P}) = \text{Var}_{\mathcal{D}_D(n, \mathcal{T}, \mathcal{P})} [\hat{y}^*(D)] = \sum_{k \in \mathcal{P}} \sigma_k^2 \tau_k^2(n, \mathcal{T}).$$

?proofname? Under squared loss, the posterior mean minimizes posterior risk, so the ex-ante (expected) value of the dataset equals the prior variance minus the posterior variance of y . Because $y = \mathbf{x}'_P \boldsymbol{\beta}_P$ and $\mathcal{P} \subseteq \mathcal{T}$,

$$\hat{y}^*(D_{\mathcal{T}, \mathcal{P}}^n) = \mathbf{x}'_{P \cap \mathcal{T}} \mathbb{E}[\boldsymbol{\beta}_{P \cap \mathcal{T}} \mid (\mathbf{y}, \mathbf{X})].$$

Taking the variance over the joint distribution of \mathbf{x}_P and (\mathbf{y}, \mathbf{X}) gives

$$\text{Var}[\hat{y}^*(D_{\mathcal{T}, \mathcal{P}}^n)] = \text{Var}(\mathbf{x}'_{P \cap \mathcal{T}} \mathbb{E}[\boldsymbol{\beta}_{P \cap \mathcal{T}} \mid (\mathbf{y}, \mathbf{X})]) = \sum_{k \in P \cap \mathcal{T}} \text{Var}(x_k) \text{Var}(\mathbb{E}[\beta_k \mid (\mathbf{y}, \mathbf{X})]),$$

using independence of covariates. Since $\text{Var}(x_k) = \sigma_k^2$ and $\mathbb{E}[\beta_k \mid (\mathbf{y}, \mathbf{X})] = \hat{\beta}_k^{\text{ridge}}(\lambda(n, \mathcal{T}), (\mathbf{y}, \mathbf{X}))$, we obtain

$$V(n, \mathcal{T}, \mathcal{P}) = \sum_{k \in P \cap \mathcal{T}} \sigma_k^2 \text{Var}_{(\mathbf{y}, \mathbf{X})} [\hat{\beta}_k^{\text{ridge}}(\lambda(n, \mathcal{T}), (\mathbf{y}, \mathbf{X}))] = \sum_{k \in P \cap \mathcal{T}} \sigma_k^2 \tau_k^2(\lambda(n, \mathcal{T})),$$

which proves the result. \square

Proposition 3 (Training Information Approximation). *The following asymptotic approximation holds*

$$V(n, \mathcal{T}, \mathcal{P}) = \sum_{k \in P \cap \mathcal{T}} \sigma_k^2 \bar{\tau}_k^2(\lambda(n, \mathcal{T})) + O\left(\sqrt{\frac{|\mathcal{T}|}{n}} + \frac{|\mathcal{T}|}{n}\right),$$

where

$$\bar{\tau}_k^2(\lambda) \equiv \frac{\sigma_k^2}{\sigma_k^2 + \lambda}.$$

?proofname? Write $\Sigma = \Sigma_{\mathcal{T}}$ and recall $\phi_k(\lambda; n, t, \sigma_{\mathcal{T}}^2) = \lambda \mathbb{E}[(G_n + \lambda \mathbf{I}_t)_{kk}^{-1}]$, where $\mathbf{G} \equiv \frac{1}{n} \mathbf{X}' \mathbf{X}$. For symmetric $\mathbf{A}, \mathbf{B} \succeq 0$, the resolvent identity gives

$$(\mathbf{A} + \lambda \mathbf{I})^{-1} - (\mathbf{B} + \lambda \mathbf{I})^{-1} = (\mathbf{A} + \lambda \mathbf{I})^{-1} (\mathbf{B} - \mathbf{A}) (\mathbf{B} + \lambda \mathbf{I})^{-1}.$$

With $\mathbf{A} = \mathbf{G}_n$ and $\mathbf{B} = \Sigma$, taking (k, k) entries and bounding by operator norms,

$$|((\mathbf{G}_n + \lambda \mathbf{I})^{-1} - (\Sigma + \lambda \mathbf{I})^{-1})_{kk}| \leq \|(\mathbf{G}_n + \lambda \mathbf{I})^{-1}\|_{\text{op}} \|\mathbf{G}_n - \Sigma\|_{\text{op}} \|(\Sigma + \lambda \mathbf{I})^{-1}\|_{\text{op}}.$$

Because $\mathbf{G}_n, \Sigma \succeq 0$, we have $\|(\mathbf{G}_n + \lambda \mathbf{I})^{-1}\|_{\text{op}} \leq \lambda^{-1}$ and $\|(\Sigma + \lambda \mathbf{I})^{-1}\|_{\text{op}} \leq (\sigma_{\min}^2 + \lambda)^{-1}$. Hence

$$|((\mathbf{G}_n + \lambda \mathbf{I})^{-1} - (\Sigma + \lambda \mathbf{I})^{-1})_{kk}| \leq \frac{1}{\lambda(\sigma_{\min}^2 + \lambda)} \|\mathbf{G}_n - \Sigma\|_{\text{op}}.$$

Taking expectations and multiplying by λ ,

$$\left| \lambda \mathbb{E}[(G_n + \lambda I)_{kk}^{-1}] - \lambda [(\Sigma + \lambda I)^{-1}]_{kk} \right| \leq \frac{1}{\sigma_{\min}^2 + \lambda} \mathbb{E} \|G_n - \Sigma\|_{\text{op}}.$$

Since $[(\Sigma + \lambda I)^{-1}]_{kk} = (\sigma_k^2 + \lambda)^{-1}$, this becomes

$$\left| \phi_k(\lambda; n, t, \sigma_{\mathcal{T}}^2) - \frac{\lambda}{\sigma_k^2 + \lambda} \right| \leq \frac{1}{\sigma_{\min}^2 + \lambda} \mathbb{E} \|G_n - \Sigma\|_{\text{op}},$$

where since covariates are ranked in descending order $\sigma_{\min}^2 = \sigma_{\max\{\mathcal{T}\}}^2$.

For Gaussian rows, the sample-covariance bound in Theorem 4.7.1. of Vershynin (2018) shows that there exists an absolute constant $C > 0$ such that

$$\mathbb{E} \|G_n - \Sigma\|_{\text{op}} \leq C \|\Sigma\|_{\text{op}} \left(\sqrt{\frac{t}{n}} + \frac{t}{n} \right),$$

and since covariates are ranked in descending order

$$\|\Sigma\|_{\text{op}} \equiv \sup_{\|x\|_2=1} \|\Sigma x\|_2 = \sup_{\|x\|_2=1} x' \Sigma x = \max_{k \in \mathcal{T}} \{\sigma_k^2\} = \sigma_{\min\{\mathcal{T}\}}^2.$$

Therefore,

$$\mathbb{E} \|G_n - \Sigma\|_{\text{op}} \leq,$$

which gives the claimed inequality. Putting together,

$$\left| \phi_k(\lambda; n, t, \sigma_{\mathcal{T}}^2) - \frac{\lambda}{\sigma_k^2 + \lambda} \right| \leq \frac{C \sigma_{\min\{\mathcal{T}\}}^2}{\sigma_{\max\{\mathcal{T}\}}^2 + \lambda} \left(\sqrt{\frac{t}{n}} + \frac{t}{n} \right).$$

The $O(\cdot)$ statement follows because σ_{\max}^2 and $(\sigma_{\min}^2 + \lambda)^{-1}$ are uniformly bounded as $\sigma_j^2 \in (0, 1)$ for all j and $\lambda \in (0, \infty)$.

Proposition 6 (Complementarity/Substitutability in Training). *Fix (n, \mathcal{P}) . For any \mathcal{T} and any $k \in \mathcal{T} \setminus \mathcal{P}$,*

$$\frac{d}{dn} \Delta_k^{\mathcal{T}}(n, \mathcal{P}, \mathcal{T}) = \underbrace{\frac{\sigma_k^2}{n^2} \sum_{m \in \mathcal{P} \cap \mathcal{T}} \sigma_m^2 \left(\frac{1 - S(\mathcal{T} \setminus \{k\})}{\sigma_k^2} \int_{-\sigma_k^2/n}^0 \frac{\partial^2}{\partial \lambda^2} \bar{\tau}_k^2(\lambda(n, \mathcal{T} \setminus \{k\}) + u) du \right)}_{(+)\text{ HPE} \times \text{Sampling spillovers}} + \underbrace{\frac{\partial}{\partial \lambda} \bar{\tau}_m^2(\lambda(n, \mathcal{T}))}_{(-)\text{ Spillover substitution}} > 0$$

if and only if

$$n < \tilde{n}(\mathcal{T}, \mathcal{P}) \in (0, \infty),$$

which is decreasing in $\mathcal{T} \setminus \{k\}$ and \mathcal{P} .

Fix $(\mathcal{P}, \mathcal{T})$ and let $k \in \mathcal{T} \setminus \mathcal{P}$. Recall

$$V(n, \mathcal{P}, \mathcal{T}) = \sum_{m \in \mathcal{P} \cap \mathcal{T}} \sigma_m^2 [1 - \bar{\phi}_m(\lambda^*(n, \mathcal{T}))], \quad \bar{\phi}_m(\lambda) = \frac{\lambda}{\sigma_m^2 + \lambda}, \quad \lambda^*(n, \mathcal{T}) = \frac{1 - S(\mathcal{T})}{n}.$$

Since $k \notin \mathcal{P}$, the prediction set does not change when we remove k , so

$$\Delta_k^\mathcal{T}(n, \mathcal{P}, \mathcal{T}) \equiv V(n, \mathcal{P}, \mathcal{T}) - V(n, \mathcal{P}, \mathcal{T} \setminus \{k\}) = \sum_{m \in \mathcal{P} \cap \mathcal{T}} \sigma_m^2 [\bar{\phi}_m(\lambda^*(n, \mathcal{T} \setminus \{k\})) - \bar{\phi}_m(\lambda^*(n, \mathcal{T}))].$$

Using

$$\lambda^*(n, \mathcal{T} \setminus \{k\}) = \frac{1 - S(\mathcal{T}) + \sigma_k^2}{n} = \lambda^*(n, \mathcal{T}) + \frac{\sigma_k^2}{n},$$

define

$$\lambda(n) \equiv \lambda^*(n, \mathcal{T}) = \frac{1 - S(\mathcal{T})}{n}, \quad \delta(n) \equiv \frac{\sigma_k^2}{n},$$

so that $\lambda^*(n, \mathcal{T} \setminus \{k\}) = \lambda + \delta$. Then

$$\Delta_k^\mathcal{T}(n, \mathcal{P}, \mathcal{T}) = \sum_{m \in \mathcal{P} \cap \mathcal{T}} \sigma_m^2 [\bar{\phi}_m(\lambda + \delta) - \bar{\phi}_m(\lambda)]. \quad (7)$$

Step 1: derivative decomposition.

Let

$$F_m(\lambda, \delta) \equiv \sigma_m^2 [\bar{\phi}_m(\lambda + \delta) - \bar{\phi}_m(\lambda)].$$

By the chain rule,

$$\frac{d}{dn} \Delta_k^\mathcal{T}(n, \mathcal{P}, \mathcal{T}) = \sum_{m \in \mathcal{P} \cap \mathcal{T}} (\partial_\lambda F_m(\lambda, \delta) \lambda' + \partial_\delta F_m(\lambda, \delta) \delta'),$$

where

$$\lambda' = \frac{d\lambda}{dn} = -\frac{1 - S(\mathcal{T})}{n^2}, \quad \delta' = \frac{d\delta}{dn} = -\frac{\sigma_k^2}{n^2}.$$

Since

$$\partial_\lambda F_m = \sigma_m^2 [\bar{\phi}'_m(\lambda + \delta) - \bar{\phi}'_m(\lambda)], \quad \partial_\delta F_m = \sigma_m^2 \bar{\phi}'_m(\lambda + \delta),$$

writing $[\bar{\phi}'_m(\lambda + \delta) - \bar{\phi}'_m(\lambda)] = \int_\lambda^{\lambda+\delta} \bar{\phi}''_m(u) du$, we obtain using the chain rule

$$\begin{aligned} \frac{d}{dn} \Delta_k^\mathcal{T} &= \sum_{m \in \mathcal{P} \cap \mathcal{T}} \sigma_m^2 \left(\lambda' \int_\lambda^{\lambda+\delta} \bar{\phi}''_m(u) du + \bar{\phi}'_m(\lambda + \delta) \delta' \right) \\ &= \sum_{m \in \mathcal{P} \cap \mathcal{T}} \sigma_m^2 \left(\frac{[1 - S(\mathcal{T})]}{n^2} \int_\lambda^{\lambda+\delta} \bar{\phi}''_m(u) du - \frac{\sigma_k^2}{n^2} \bar{\phi}'_m(\lambda + \delta) \right). \end{aligned} \quad (8)$$

This is the claimed decomposition: a positive term (HPE \times observation spillovers) minus a negative term (spillover substitution).

Step 2: signs of the two forces.

We use only concavity of $\bar{\phi}_m$. First, $\bar{\phi}'_m(\lambda) > 0$ and

$$\bar{\phi}''_m(\lambda) = -\frac{2\sigma_m^2}{(\sigma_m^2 + \lambda)^3} < 0,$$

so $\bar{\phi}'_m$ is strictly decreasing. Since $\delta > 0$,

$$\bar{\phi}'_m(\lambda) - \bar{\phi}'_m(\lambda + \delta) > 0.$$

Hence the first term in (8) is strictly positive:

$$\frac{\sigma_m^2[1 - S(\mathcal{T})]}{n^2} [\bar{\phi}'_m(\lambda) - \bar{\phi}'_m(\lambda + \delta)] > 0.$$

This captures that as n increases (so λ decreases), concavity of $\bar{\phi}_m$ amplifies the gain from the extra training covariate k (*House Party Effect* interacting with *observation spillovers*).

The second term is strictly negative since $\bar{\phi}'_m(\lambda + \delta) > 0$ and $\sigma_m^2, \sigma_k^2 > 0$:

$$\frac{\sigma_m^2 \sigma_k^2}{n^2} \bar{\phi}'_m(\lambda + \delta) > 0 \quad \Rightarrow \quad -\frac{\sigma_m^2 \sigma_k^2}{n^2} \bar{\phi}'_m(\lambda + \delta) < 0.$$

This is *spillover substitution*: as n increases, the direct impact $\delta = \sigma_k^2/n$ of covariate k on regularization shrinks.

Therefore

$$\frac{d}{dn} \Delta_k^\mathcal{T} > 0 \iff [1 - S(\mathcal{T})] A(n, \mathcal{T}, \mathcal{P}) > \sigma_k^2 B(n, \mathcal{T}, \mathcal{P}),$$

where

$$A(n, \mathcal{T}, \mathcal{P}) \equiv \sum_{m \in \mathcal{P} \cap \mathcal{T}} \sigma_m^2 [\bar{\phi}'_m(\lambda) - \bar{\phi}'_m(\lambda + \delta)] > 0, \quad B(n, \mathcal{T}, \mathcal{P}) \equiv \sum_{m \in \mathcal{P} \cap \mathcal{T}} \sigma_m^2 \bar{\phi}'_m(\lambda + \delta) > 0.$$

Step 3: existence and dependence of the threshold $\tilde{n}(\mathcal{T}, \mathcal{P})$.

Define

$$\Psi(n; \mathcal{T}, \mathcal{P}) \equiv [1 - S(\mathcal{T})] A(n, \mathcal{T}, \mathcal{P}) - \sigma_k^2 B(n, \mathcal{T}, \mathcal{P}),$$

so that $\frac{d}{dn} \Delta_k^\mathcal{T} > 0 \iff \Psi(n; \mathcal{T}, \mathcal{P}) > 0$.

For large n , we have $\lambda = \mathcal{O}(1/n)$ and $\delta = \mathcal{O}(1/n)$. Using the Taylor expansion

$$\bar{\phi}'_m(\lambda) = \frac{\sigma_m^2}{(\sigma_m^2 + \lambda)^2} = \frac{1}{\sigma_m^2} - \frac{2\lambda}{\sigma_m^4} + \mathcal{O}(\lambda^2),$$

we obtain uniformly over m :

$$\begin{aligned}\bar{\phi}'_m(\lambda) - \bar{\phi}'_m(\lambda + \delta) &= \frac{2\delta}{\sigma_m^4} + \mathcal{O}\left(\frac{1}{n^2}\right), \\ \bar{\phi}'_m(\lambda + \delta) &= \frac{1}{\sigma_m^2} + \mathcal{O}\left(\frac{1}{n}\right).\end{aligned}$$

Hence, as $n \rightarrow \infty$,

$$A(n, \mathcal{T}, \mathcal{P}) = 2\delta \sum_{m \in \mathcal{P} \cap \mathcal{T}} \frac{1}{\sigma_m^2} + \mathcal{O}\left(\frac{1}{n^2}\right) = \frac{2\sigma_k^2}{n} \sum_{m \in \mathcal{P} \cap \mathcal{T}} \frac{1}{\sigma_m^2} + \mathcal{O}\left(\frac{1}{n^2}\right),$$

and

$$B(n, \mathcal{T}, \mathcal{P}) = \sum_{m \in \mathcal{P} \cap \mathcal{T}} 1 + \mathcal{O}\left(\frac{1}{n}\right) = |\mathcal{P} \cap \mathcal{T}| + \mathcal{O}\left(\frac{1}{n}\right).$$

Therefore

$$\Psi(n; \mathcal{T}, \mathcal{P}) = \frac{2[1 - S(\mathcal{T})]\sigma_k^2}{n} \sum_{m \in \mathcal{P} \cap \mathcal{T}} \frac{1}{\sigma_m^2} - \sigma_k^2 |\mathcal{P} \cap \mathcal{T}| + \mathcal{O}\left(\frac{1}{n}\right).$$

For large n , the negative constant term dominates, so $\Psi(n; \mathcal{T}, \mathcal{P}) < 0$ and thus $\frac{d}{dn}\Delta_k^\mathcal{T} < 0$. For n close to 1 (or small enough), the positive term dominates by continuity of all components, so $\Psi(n; \mathcal{T}, \mathcal{P}) > 0$. Since Ψ is continuous in n , there exists at least one threshold $\tilde{n}(\mathcal{T}, \mathcal{P})$ such that

$$\frac{d}{dn}\Delta_k^\mathcal{T} > 0 \iff n < \tilde{n}(\mathcal{T}, \mathcal{P}).$$

Using the leading-order approximation $\Psi(n; \mathcal{T}, \mathcal{P}) \approx 0$ yields

$$\tilde{n}(\mathcal{T}, \mathcal{P}) \approx \frac{2[1 - S(\mathcal{T})]}{|\mathcal{P} \cap \mathcal{T}|} \sum_{m \in \mathcal{P} \cap \mathcal{T}} \frac{1}{\sigma_m^2}.$$

This expression makes the dependence on \mathcal{T} transparent: $\tilde{n}(\mathcal{T}, \mathcal{P})$ is approximately proportional to the *residual variance* $1 - S(\mathcal{T})$ and thus is *decreasing in* $S(\mathcal{T})$: richer training sets (larger $S(\mathcal{T})$) shift the peak to the left.

In summary, the derivative admits the decomposition (8); the first term (HPE \times observation spillovers) is positive, the second (spillover substitution) is negative, and their balance induces a unique threshold $\tilde{n}(\mathcal{T}, \mathcal{P})$, decreasing in $S(\mathcal{T})$, such that $\frac{d}{dn}\Delta_k^\mathcal{T} > 0$ if and only if $n < \tilde{n}(\mathcal{T}, \mathcal{P})$. \square

Proposition 12 (Closed Form Accuracy). *If $\sigma_k^2 = \sigma^2 \mathbf{1}(k \in \mathcal{I})$, the accuracy is*

$$A(n, p, t) = \frac{S(\min\{t, p\})}{\frac{1-S(t)}{n\sigma^2} + 1}.$$

?proofname? Let $\mathcal{I} \subset \mathbb{N}$ be the (finite) set of informative covariates. Assume:

$$\sigma_k^2 = s^2 > 0 \text{ for all } k \in \mathcal{I} \text{ and } \sigma_k^2 = 0 \text{ for } k \notin \mathcal{I};$$

$\text{Var}(y) = 1$, so $|\mathcal{I}|s^2 = 1$;

the agent is constrained to observe covariates in a fixed order, so any feasible training and prediction sets are prefixes of \mathbb{N} .

For $m \in \mathbb{N}$, let $S(m) \equiv \sum_{k=1}^m \sigma_k^2$ denote the cumulative variance, and for $t, p \geq 0$ write

$$S(t) \equiv S(\lfloor t \rfloor), \quad S(p) \equiv S(\lfloor p \rfloor),$$

viewing t and p as (possibly real) indexes along the ordered list.

Step 1: Optimal sets are prefixes at full capacity.

Under the fixed-order constraint, any feasible training set has the form

$$\mathcal{T} = \{1, \dots, m_T\} \quad \text{with } m_T \leq t,$$

and any prediction set

$$\mathcal{P} = \{1, \dots, m_P\} \quad \text{with } m_P \leq p.$$

Thus $\mathcal{P} \cap \mathcal{T} = \{1, \dots, \min\{m_T, m_P\}\}$.

Each term in V is nonnegative:

$$\sigma_k^2 \left[1 - \bar{\phi}_k(\lambda^*) \right] = \frac{\sigma_k^4}{\sigma_k^2 + \lambda^*} \geq 0.$$

Hence enlarging m_T or m_P weakly increases V . Therefore the supremum is attained at

$$m_T = t, \quad m_P = p,$$

so we may write

$$A(n, t, p) = V(n, \mathcal{T}, \mathcal{P}_p), \quad \mathcal{T}_t = \{1, \dots, t\}, \quad \mathcal{P}_p = \{1, \dots, p\}.$$

Step 2: Identify the intersection mass.

Since only $k \in \mathcal{I}$ are informative (with variance s^2), the total variance of informative covariates in the intersection is

$$\sum_{k \in \mathcal{P}_p \cap \mathcal{T}_t} \sigma_k^2 = S(\min\{t, p\}).$$

Because all informative coordinates have the same variance s^2 , the number of informative coordinates in the intersection is

$$M = \frac{S(\min\{t, p\})}{s^2}.$$

Step 3: Compute the contribution of each informative coordinate.

For any informative k ,

$$1 - \bar{\phi}_k(\lambda) = 1 - \frac{\lambda}{\sigma_k^2 + \lambda} = \frac{\sigma_k^2}{\sigma_k^2 + \lambda},$$

so

$$\sigma_k^2 \left[1 - \bar{\phi}_k(\lambda) \right] = \frac{\sigma_k^4}{\sigma_k^2 + \lambda}.$$

In our problem, $\lambda = \lambda^*(n, \mathcal{T}) = \frac{1 - S(t)}{n}$, which does not depend on k . Hence all informative coordinates in $\mathcal{P}_p \cap \mathcal{T}_t$ contribute the same amount, and

$$V(n, \mathcal{T}, \mathcal{P}_p) = M \cdot \frac{s^4}{s^2 + \lambda^*(n, \mathcal{T})}.$$

Substitute $M = S(\min\{t, p\})/s^2$:

$$V(n, \mathcal{T}, \mathcal{P}_p) = \frac{S(\min\{t, p\})}{s^2} \cdot \frac{s^4}{s^2 + \lambda^*(n, \mathcal{T})} = S(\min\{t, p\}) \cdot \frac{s^2}{s^2 + \lambda^*(n, \mathcal{T})}.$$

Step 4: Substitute λ^ and simplify.*

We have

$$\lambda^*(n, \mathcal{T}) = \frac{1 - S(t)}{n},$$

so

$$\frac{s^2}{s^2 + \lambda^*(n, \mathcal{T})} = \frac{s^2}{s^2 + \frac{1 - S(t)}{n}} = \frac{1}{1 + \frac{1 - S(t)}{ns^2}}.$$

Therefore

$$A(n, t, p) = V(n, \mathcal{T}, \mathcal{P}_p) = \frac{S(\min\{t, p\})}{1 + \frac{1 - S(t)}{ns^2}},$$

which is the desired expression.

If we reparametrize the sample size in effective units $n' := ns^2$ (so that $s^2 = 1$ in the formula), this becomes

$$A(n', t, p) = \frac{S(\min\{t, p\})}{1 + \frac{1 - S(t)}{n'}},$$

as claimed. □

Proposition 15 (Harmful Acquisition). *The incumbent I always acquires the entrant E .*

?proofname? If the acquisition does not occur, each firm $i \in \{I, E\}$ sets $p_i = A_1$ and sells to its unit mass of captive buyers, so profits are $\Pi_I^{\text{no}} = \Pi_E^{\text{no}} = A_1$ and

$$W^{\text{no}} = \Pi_I^{\text{no}} + \Pi_E^{\text{no}} = 2A_1.$$

If the acquisition occurs, the merged firm sets $p_2 = A_2$ and earns profit $\Pi_2 = A_2$, while the planner incurs the dynamic loss ξ . Hence

$$W^{\text{acq}} = \Pi_2 - \xi = A_2 - \xi.$$

The acquisition is (weakly) harmful whenever $W^{\text{acq}} \leq W^{\text{no}}$, that is,

$$A_2 - \xi \leq 2A_1 \iff A_2 - 2A_1 \leq \xi.$$

Under the symmetry and equal-informativeness assumptions, Proposition 12 implies that

$$A_2 - 2A_1 = \Delta(n, k) = \frac{n}{\left(\frac{1+n}{k} - 1\right) \left(\frac{2(n+1)}{k} - 1\right)},$$

which is increasing in k . Define $\tilde{k}(n, \xi)$ as the (unique) solution to $\Delta(n, k) = \xi$. By monotonicity, $A_2 - 2A_1 \leq \xi$ holds if and only if $k \leq \tilde{k}(n, \xi)$, and solving $\Delta(n, k) = \xi$ for k yields

$$\tilde{k}(n, \xi) = \frac{n+1}{2} \frac{\sqrt{\xi(8n+\xi)} - 3\xi}{n-\xi}.$$

This proves the claim. □

Lemma 4 (Optimal Pool Price). *The optimal pool price is*

$$P^* = \frac{A_2}{\alpha + 1}.$$

?proofname? The pool price is

$$P^* = \arg \max_P \{PD(P - A_2)\},$$

which can be solved using the first-order condition

$$D(P - A_2) + PD'(P - A_2) = 0.$$

□

Lemma 6 (Sample Fragmentation Price). *If the brokers have different observations on the same covariates,*

$$p_i = \min \left\{ A_2 - A_1, \frac{A_2}{2 + \alpha} \right\},$$

and the buyers will buy from both brokers.

?proofname? We follow Lerner and Tirole (2004). **Demand Margin Binds.** Suppose that the brokers offer prices $\mathcal{P} \equiv (p_1, p_2)$, and wlog $p_1 \leq p_2$. Prediction Sellers decide how many

datasets to buy.

$$\mathcal{V}(\mathcal{P}) = \max_{q \in \{1,2\}} \{A_q - p_1 - p_2 \mathbf{1}(\{q = 2\})\}$$

Second, the user adopts the technology if and only if

$$\mathcal{V}(\mathcal{P}) \geq \theta.$$

Lerner and Tirole (2004) demonstrate the existence of a symmetric equilibrium. Individual data sellers solve

$$\hat{p} = \arg \max_{p_i} \{p_i D(p_i + \hat{p} - A_2)\}$$

which has FOC

$$\hat{p} D'(2\hat{p} - A_2) + D(2\hat{p} - A_2) = 0$$

which has a unique solution by hazard-rate monotonicity. It can be seen as selling the whole pool setting total price P and keeping $p_i = P - \hat{p}$ for itself. Therefore

$$\hat{P} = \arg \max_P \{(P - \hat{p}) D(P - A_2)\}.$$

The term \hat{p} can be seen as a marginal cost $\hat{c} = \hat{p}$. In this interpretation when there is the pool $c^* = 0$ so by revealed preference

$$\hat{P} \geq P^*.$$

If demand margin binds in the absence of a pool then the pool reduces price paid by data buyers. This means that if all datasets can increase the price marginally without being excluded, the pool is pro-competitive.

With our CDF G ,

$$p_{\text{dem}} = \frac{A_2}{2 + \alpha}.$$

Competition Margin Binds. Define the price when the competition margin binds will be p_{comp} defined by

$$A_2 - 2p_{\text{comp}} = \max_{q \in \{0,1\}} \{A_q - qp_{\text{comp}}\}.$$

If $A_1 - p_{\text{comp}} \geq 0$, then $p_{\text{comp}} = A_2 - A_1$. This is consistent because $A_1 - (A_2 - A_1) > 0$ by concavity of $V(n, t)$ in n . Otherwise, if $A_1 - p_{\text{comp}} < 0$, then $p_{\text{comp}} = A_2/2$. This is not consistent because $A_1 - A_2/2 > 0$ by concavity of $V(n, t)$ in n . \square

Proposition 18 (Welfare of Observation Pooling). *A pool of observations is procompetitive if and only if*

$$k < \hat{k}(n, \alpha) \equiv 1 - \frac{n}{2\alpha}.$$

?proofname? Pools are procompetitive if the demand margin binds i.e.

$$p_{\text{comp}} > p_{\text{dem}} \iff A_2 - A_1 > \frac{A_2}{2 + \alpha} \iff \alpha > \alpha_{\text{marg}} \equiv \frac{2A_1 - A_2}{A_2 - A_1}.$$

In this case

$$p_i = \frac{A_2}{2 + \alpha}.$$

This implies that observation pooling can be procompetitive when n is not too large and k is not too small, meaning data is relatively abundant and models are relatively complex. Furthermore, as the RHS is increasing in Q , data pools are more likely to be competitive if Q is small meaning if data fragmentation is limited.

Otherwise if the competition margin binds,

$$p_i = A_2 - A_1.$$

the pool is procompetitive if the pool price is lower than the competition price, i.e.

$$P^* < Qz(Q) \iff \frac{A_2}{\alpha + 1} < 2(A_2 - A_1) \iff \alpha > \alpha_{\text{comp}} \equiv \frac{A_1 - \frac{A_2}{2}}{A_2 - A_1}.$$

As $\alpha_{\text{marg}} > \alpha_{\text{comp}}$, the relevant threshold is α_{comp} and so a pool of observations is procompetitive if and only if

$$\alpha > \alpha_{\text{comp}} = \frac{A_1 - \frac{A_2}{2}}{A_2 - A_1}.$$

This implies that as k and n increase it becomes less likely that the pool is procompetitive. When data is abundant and demand is inelastic observation pools are anticompetitive. Direct application of Proposition 5 in Lerner and Tirole (2004) implies that the pool is strongly unstable, therefore enforcing independent licensing of datasets will prevent pooling if and only if the pool is welfare-reducing. In this case each data broker will charge $p_{\text{comp}} = A_2 - A_1$. \square

Lemma 5 (Covariate Fragmentation Price). *If the brokers have distinct covariates each B_i prices at*

$$p_i = \frac{A_2}{2 + \alpha},$$

and the buyers will buy from both brokers.

?proofname? If $A_1 - p_{\text{comp}} \geq 0$, then $p_{\text{comp}} = A_2 - A_1$. This is not consistent because $A_1 - (A_2 - A_1) < 0$ by convexity of $V(n, t)$ in t . Otherwise, if $A_1 - p_{\text{comp}} < 0$, then $p_{\text{comp}} = A_2/2$. This is consistent because $A_1 - A_2/2 < 0$ by convexity of $V(n, t)$ in t . Demand margin always binds as

$$p_{\text{comp}} > p_{\text{dem}} \iff \frac{A_2}{2} > \frac{A_2}{2 + \alpha}.$$

\square

Lemma 7 (Shopper-price equilibrium). *Fix (A_1, A_2) . The pricing subgame in shopp[Shopper-price equilibrium]er prices admits a unique equilibrium under trembling-hand perfection:*

$$p_i^s = (A_i - A_j)^+, \quad i \neq j, \tag{5}$$

so that (i) if $A_i > A_j$, firm i sets $p_i^s = A_i - A_j$, firm j sets $p_j^s = 0$, and all shoppers buy from i ; (ii) if $A_i = A_j$, both set $p_i^s = p_j^s = 0$ and shoppers can be split arbitrarily.

?proofname? We first show that (5) is an equilibrium. Consider $A_i > A_j$. If i sets $p_i^s = A_i - A_j$ and j sets $p_j^s = 0$, shoppers are indifferent and (by the refinement) all go to i .¹⁹ Any deviation by i :

If $p_i^s > A_i - A_j$, then $A_i - p_i^s < A_j - p_j^s = A_j$ and i loses all shoppers, strictly reducing its shopper revenue to 0.

If $p_i^s < A_i - A_j$, i still serves all shoppers but leaves revenue on the table; since demand is inelastic at one (all shoppers) at the margin, profit increases by raising p_i^s up to $A_i - A_j$.

Any deviation by j :

If $p_j^s > 0$, then $A_j - p_j^s < A_j = A_i - p_i^s$ and j serves no shoppers with the same zero shopper revenue; under trembling hand, the weakly dominated positive price is eliminated in the limit, selecting $p_j^s = 0$.

If $p_j^s < 0$ is infeasible; if $p_j^s = 0$ already, no profitable deviation exists.

Thus (5) is an equilibrium when $A_i > A_j$. The case $A_i = A_j$: for any $p_i^s = p_j^s$, shoppers are indifferent and each firm earns $\sigma p_i^s/2$ from shoppers; any unilateral increase loses all shoppers, any decrease reduces price with the same demand, so $p_i^s = p_j^s = 0$ is the unique trembling-hand limit (positive common prices are not robust to small payoff perturbations). Symmetry covers $A_j > A_i$.

Uniqueness under trembling-hand perfection follows from the standard Bertrand-vertical-differentiation logic: if $A_i > A_j$, any equilibrium must have the higher-quality firm serving all shoppers; then the highest sustainable price for i that keeps all shoppers is $A_i - A_j$, and the lower-quality firm's best response is any price with zero demand and zero revenue, refined to 0. □

Proposition 21. *If $c < \bar{c}$, the planner prefers real non-exclusivity to de facto exclusivity.*

?proofname?

$$W^e \equiv \frac{1+s}{2}\bar{A} - c < W^{ne} \equiv \bar{A} - 2c \iff c < \frac{1-s}{2}\bar{A}.$$

which contradicts $c < \bar{c} \equiv \frac{1-s}{2}\bar{A}$. □

Proposition 22. *We distinguish two subcases:*

- **De facto exclusivity:** *Without loss of generality, let $\ell_1 = 1$, $\ell_2 = 0$, the unique NE is*

$$r_1 = 1, \quad r_2 = 0.$$

¹⁹Formally, with tiny perturbations (trembles) that give i an ε quality advantage or j an ε higher price with positive probability, the unique limit assigns the shoppers to i .

- **Both firms license:** $\ell_1 = \ell_2 = 1$. The unique symmetric mixed strategy equilibrium has each firm investing with probability

$$\xi^* = \frac{A(1+s) - 2c}{2As}, \quad \xi^* \in [0, 1].$$

In equilibrium, profits are zero: $\Pi_i^{\text{mix}} = 0$.

?proofname? We condition throughout on the licensing profile (ℓ_1, ℓ_2) and characterize the equilibrium in proprietary investment $r_i \in \{0, 1\}$.

De facto exclusivity. Suppose without loss of generality that $\ell_1 = 1$ and $\ell_2 = 0$. Firm 2 cannot profitably invest in proprietary covariates, so its best response is

$$r_2 = 0.$$

If firm 1 invests, it obtains monopoly access to shoppers. By construction of the threshold \bar{c} , monopoly investment yields strictly positive net profit when $c > \bar{c}$, whereas not investing yields zero. Hence firm 1 strictly prefers to invest:

$$r_1 = 1.$$

Thus $(r_1, r_2) = (1, 0)$ is the unique Nash equilibrium in proprietary investment when $(\ell_1, \ell_2) = (1, 0)$.

Both firms license. Now suppose $\ell_1 = \ell_2 = 1$. Let ξ denote the probability with which firm $j \neq i$ invests, so that firm i 's expected payoff from choosing $r_i = 1$ is

$$\begin{aligned} \Pi_i(r_i = 1) &= (1 - \xi) \left[\frac{1-s}{2} \bar{A} + s\bar{A} \right] + \xi \left[\frac{1-s}{2} \bar{A} \right] - c \\ &= \frac{1-s}{2} \bar{A} + s\bar{A}(1 - \xi) - c, \end{aligned}$$

while the payoff from not investing is

$$\Pi_i(r_i = 0) = 0.$$

First, no symmetric pure-strategy equilibrium exists. If $(r_1, r_2) = (0, 0)$, then $\xi = 0$ and

$$\Pi_i(1) = \frac{1-s}{2} \bar{A} + s\bar{A} - c > 0$$

by the definition of \bar{c} , so each firm has a profitable deviation to $r_i = 1$. If $(r_1, r_2) = (1, 1)$, then $\xi = 1$ and

$$\Pi_i(1) = \frac{1-s}{2} \bar{A} - c < 0$$

for $c > \bar{c}$, so each firm has a profitable deviation to $r_i = 0$. Hence neither $(0, 0)$ nor $(1, 1)$ is a symmetric equilibrium.

Consider now a symmetric mixed-strategy equilibrium in which each firm invests with probability $\xi^* \in (0, 1)$. Symmetry requires that firm i be indifferent between $r_i = 1$ and $r_i = 0$, so

$$\Pi_i(r_i = 1) = \Pi_i(r_i = 0) = 0.$$

Substituting the expression for $\Pi_i(r_i = 1)$ and solving for ξ yields

$$\frac{1-s}{2}\bar{A} + s\bar{A}(1-\xi^*) - c = 0 \implies \xi^* = \frac{\bar{A}(1+s) - 2c}{2\bar{A}s}.$$

Under Assumption ??, this mixing probability satisfies $\xi^* \in [0, 1]$. Moreover, because $\Pi_i(r_i = 1)$ is strictly decreasing in ξ , the indifference condition can hold for at most one value of ξ , so the symmetric mixed strategy equilibrium is unique.

Finally, since firm i is indifferent between $r_i = 1$ and $r_i = 0$ at ξ^* , its expected equilibrium profit (before paying the license fee) equals the payoff from not investing:

$$\Pi_i^{\text{mix}} = 0.$$

This proves the two cases stated in the proposition. □

Proposition 23. *If $c \in (\bar{c}, \bar{\bar{c}})$, the planner prefers de facto exclusivity to real non-exclusivity.*

?proofname? Under exclusivity:

$$W^e = \frac{1+s}{2}\bar{A} - c.$$

Under non-exclusivity (mixed equilibrium):

$$W^{ne} = \xi^2(\bar{A} - 2c) + 2\xi(1-\xi) \left(\frac{1+s}{2}\bar{A} - c \right) = \frac{c(\frac{c}{\bar{A}} - s - 1) + \bar{A} \left(\frac{s+1}{2} \right)^2}{s}.$$

One can verify that

$$W^e > W^{ne} \quad \text{for } c \in (\bar{c}, \bar{\bar{c}}).$$

□

B Extensions

B.1 Scope as Model Complexity and LLMs

Scope as Model Complexity. Instead, make no restriction on Σ . Furthermore suppose the firm observes all covariates for all individuals but faces constraints on the number of covari-

| House i | y^i | x_{size}^i | x_{year}^i | x_{dist}^i | x_{sun}^i | |
|-----------|----------|---------------------|---------------------|---------------------|--------------------|---|
| 0 | ? | x_{size}^0 | NA | x_{dist}^0 | NA | } Prediction Vector: \mathbf{x}'_p |
| 1 | y^1 | x_{size}^1 | x_{year}^1 | NA | x_{sun}^1 | |
| 2 | y^2 | x_{size}^2 | x_{year}^2 | NA | x_{sun}^2 | } Training Matrix: $\mathbf{M}_{\mathcal{T}}^{(n)}$ |
| \vdots | \vdots | \vdots | \vdots | \vdots | \vdots | |
| n | y^n | x_{size}^n | x_{year}^n | NA | x_{sun}^n | |

Table 1: Example of Zillow dataset with prediction covariates $\mathcal{P} = \{\text{size}, \text{dist}\}$ and training covariates $\mathcal{T} = \{\text{size}, \text{year}, \text{sun}\}$, where *size* denotes square meters, *dist* the distance to the nearest supermarket, *year* the construction year, and *sun* the daily sunlight exposure.

ates it can effectively use in the learning and targeting steps. The scope of learning, ℓ , is the number of principal components the firm can use in learning. The scope of targeting, t , is the number of principal components that can be used in targeting. This interpretation captures the *model complexity*, which reflects the higher computing cost deriving from analyzing more covariates.

To reduce the dimensionality whilst extracting the maximum information in the constraints, Jolliffe (2002) shows that the optimal procedure is Principal Component Analysis (PCA). Let the eigendecomposition of the variance/covariance matrix be

$$\Sigma = \mathbf{U}\mathbf{S}\mathbf{U}', \quad \mathbf{S} = \text{diag}(s_1 \geq \dots \geq s_{\bar{\ell}} \geq 0), \quad \mathbf{U} \text{ orthonormal.}$$

Define principal components $\mathbf{z}^i \equiv \mathbf{x}^i \mathbf{U}$. Then

$$\mathbf{z}^i \sim \mathcal{N}(0, \mathbf{\Lambda}), \quad \mathbf{z}_j^i \text{ are uncorrelated with variances } s_j.$$

Remark 1 (Application to Large Language Models (LLMs)). Although LLMs are trained with cross-entropy loss, near a trained solution their behavior can be well approximated by a linear predictor under squared loss in a suitable linear transformation of the covariates (MacKay (1992); Jacot, Gabriel, and Hongler (2018)). In this local view, our primitives map directly: the scale of learning n corresponds to the amount of training information (e.g., the number of training observations/tokens), the scope of learning ℓ captures the effective number of informative directions used at the learning stage, and the scope of targeting t captures the amount of information observed at the targeting stage for specific instances. Under this mapping, comparative statics in (n, ℓ, t) align with empirical scaling laws for language models (Kaplan et al. (2020)). Supplying richer information at prediction time corresponds to increasing t via retrieval-augmented inputs (Lewis et al. (2020)), with benefits contingent on relevance and known long-context effects (Liu et al. (2023)).

B.2 Shrinkage Interpretation

We express the Bayes estimator in terms of a generalization of the ordinary least-squares (OLS) estimator — the minimum-norm least-squares (MNLS) estimator, defined as

$$\hat{\boldsymbol{\beta}}_{\mathcal{T}}^{\text{MNLS}} \equiv (\mathbf{X}'_{\mathcal{T}} \mathbf{X}_{\mathcal{T}})^+ \mathbf{X}'_{\mathcal{T}} \mathbf{y} = \begin{cases} \hat{\boldsymbol{\beta}}_{\mathcal{T}}^{\text{OLS}}, & \text{if } |\mathcal{T}| \leq n, \\ \min_{\mathbf{b}_{\mathcal{T}}} \{\|\mathbf{b}_{\mathcal{T}}\|_2 : \mathbf{X}_{\mathcal{T}} \mathbf{b}_{\mathcal{T}} = \mathbf{y}\}, & \text{if } |\mathcal{T}| > n, \end{cases}$$

where $(\cdot)^+$ denotes the Moore–Penrose pseudo-inverse.²⁰ The MNLS is the estimator that the firm would adopt if the residual variance were approximately zero (i.e., the cumulative signal $S(\mathcal{T}) \approx 1$). It comes in two flavors, depending on whether the number of parameters is greater than the sample size:

- Underparametrized regime ($n \geq |\mathcal{T}|$): the MNLS estimator coincides with the OLS estimator, which is uniquely defined because $\mathbf{X}'_{\mathcal{T}} \mathbf{X}_{\mathcal{T}}$ is invertible.
- Overparametrized regime ($n < |\mathcal{T}|$): the OLS estimator is not defined because the system $\mathbf{X}_{\mathcal{T}} \mathbf{b}_{\mathcal{T}} = \mathbf{y}$ has infinitely many solutions; the MNLS chooses the solution with the smallest Euclidean norm.

The MNLS is useful because it is well-defined in both regimes and coincides with the maximum-likelihood estimator. The Bayes estimator is a shrinkage of the MNLS estimator towards the prior mean $\mathbf{0}_{|\mathcal{T}|}$

Corollary 9. *The Bayes Estimator is the MNLS estimator with shrinkage:*

$$\mathbb{E}[\boldsymbol{\beta}_{\mathcal{T}} \mid \mathbf{M}_{\mathcal{T}}] = \left(\underbrace{(1 - S(\mathcal{T}))}_{\text{Shrinkage Factor}} \cdot (\mathbf{X}'_{\mathcal{T}} \mathbf{X}_{\mathcal{T}})^+ + \mathbf{I}_{\mathcal{T}} \right)^{-1} \hat{\boldsymbol{\beta}}_{\mathcal{T}}^{\text{MNLS}}.$$

Because it is the maximum likelihood estimator, the MNLS estimator attributes all the variation in the learning matrix $\mathbf{M}_{\mathcal{T}}$ to the parameters $\boldsymbol{\beta}_{\mathcal{T}}$. In reality, a fraction $1 - S(\mathcal{T})$ of the variation in \mathbf{y} is residual variance and not due to $\boldsymbol{\beta}_{\mathcal{T}}$. The posterior mean corrects for this by shrinking $\hat{\boldsymbol{\beta}}_{\mathcal{T}}^{\text{MNLS}}$ towards the prior mean $\mathbf{0}_{|\mathcal{T}|}$ with a shrinkage factor equal to the residual variance $1 - S(\mathcal{T})$. Adding a new covariate $j \notin \mathcal{T}$ reduces the residual variance by s_j , the variance of x_j , thereby lowering the shrinkage factor and the weight of the prior mean. Hence, the posterior mean moves closer to the MNLS estimator. Hence, covariates lend precision to each other: observing a new variable improves the accuracy of the estimated parameters of the others.

²⁰For a matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$, the Moore–Penrose pseudo-inverse is the unique matrix $\mathbf{A}^+ \in \mathbb{R}^{m \times n}$ satisfying

$$\mathbf{A} \mathbf{A}^+ \mathbf{A} = \mathbf{A}, \quad \mathbf{A}^+ \mathbf{A} \mathbf{A}^+ = \mathbf{A}^+, \quad (\mathbf{A} \mathbf{A}^+)' = \mathbf{A} \mathbf{A}^+, \quad (\mathbf{A}^+ \mathbf{A})' = \mathbf{A}^+ \mathbf{A}.$$

B.3 Double Descent

Corollary 10. *If covariates in \mathcal{L} are highly informative, the Bayes Estimator is equivalent to the ridgeless estimator and the MNLS estimator*

$$\lim_{S(\mathcal{L}) \rightarrow 1^-} \mathbb{E}[\boldsymbol{\beta}_{\mathcal{L}} | \mathbf{M}_{\mathcal{L}}] = \lim_{\lambda \rightarrow 0^+} \hat{\boldsymbol{\beta}}_{\mathcal{L}}^{\text{ridge}}(\lambda) = \hat{\boldsymbol{\beta}}_{\mathcal{L}}^{\text{MNLS}}.$$

In general, sophisticated algorithms are needed to compute or approximate the posterior mean $\mathbb{E}[\boldsymbol{\beta}_{\mathcal{L}} | \mathbf{M}_{\mathcal{L}}]$. Instead, the MNLS can be obtained by a simple machine learning algorithm, *gradient descent*. This equivalence therefore shows that once the data is sufficiently rich, even such a rudimentary algorithm approximates the Bayes estimator arbitrarily well. When data is linear-separable, prediction accuracy is driven almost entirely by data, not by algorithms.

Remark 2. The result also sheds light on a central puzzle in modern statistics and machine learning: the double descent phenomenon first discussed in Belkin et al. (2019). Classical statistics tells us the prediction error of gradient descent is U-shaped in the number of parameters $|\mathcal{L}|$: with too few parameters the model underfits, while beyond the optimum $|\mathcal{L}|^* \in (0, n)$ prediction error increases due to overfitting, as residual variation $\boldsymbol{\varepsilon}$ is mistakenly attributed to $\boldsymbol{\beta}_{\mathcal{L}}$. However, empirical work shows that expanding \mathcal{L} further can reduce the error again—the second descent in the error. Double descent is not yet fully understood: the dominant explanations rely on intricate properties of high-dimensional geometry (see Hastie et al. (2020)). Our model offers a simpler account that also applies to low-dimensions. As the learning set \mathcal{L} expands, the residual variance $1 - S(\mathcal{L})$ decreases, and the shrinkage operator in the Bayes estimator vanishes. When $S(\mathcal{L}) \approx 1$, the Bayes estimator is arbitrarily close to the MNLS even in finite samples, so gradient descent is approximately optimal.

B.4 Connection with Shannon’s Information Theory

Remark 3. Let a real-valued additive white Gaussian residual variance (AWGN) channel be given by

$$y = w + z, \quad z \sim \mathcal{N}(0, \sigma^2),$$

with an input power constraint $\mathbb{E}[w^2] \leq P$. Classical results due to Shannon (1948) show that the mutual information between w and y is²¹

$$I(w; y) = \frac{1}{2} \log_2 \left(1 + \frac{P}{\sigma^2} \right) \quad \text{nats.} \quad (\text{R.1})$$

If the channel is decomposed into independent “frequency” slices indexed by $j \in \mathcal{T}$ that

²¹See C. E. Shannon, *Bell System Technical Journal*, 1948, eq. (26); or T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed., §9.1.

each carry an SNR of

$$\text{SNR}_j = \frac{s_j}{\lambda^*},$$

then (R.1) adds up across slices by orthogonality. The total mutual information revealed by a learning sample of *strength* t is therefore²²

$$I_{\mathcal{T}}(\lambda^*) = \frac{1}{2} \sum_{j \in \mathcal{T}} \log_2 \left(1 + \frac{s_j}{\lambda^*} \right). \quad (\text{R.2})$$

Equation (R.2) is exactly the functional that appears in our model. Thus the economic value function I study,

$$v(t) = \sum_{j \in \mathcal{T}} \frac{t \lambda_j}{1 + t \lambda_j},$$

equals

$$v(\mathcal{L}, \mathcal{T}) = 2 \left(\frac{I'_{\mathcal{T}}(\lambda^*(\mathcal{L}))}{\lambda^*(\mathcal{L})} - I_{\mathcal{T}}(\lambda^*(\mathcal{L})) \right),$$

linking our “value of accuracy” directly to the canonical Shannon measure of information. Two substantive insights follow:

1. **Capacity-driven diminishing returns.** Because $I''(t) < 0$ by Shannon’s law, marginal economic value $v'(t) = 2I'(t)$ must also fall. No additional curvature assumption is needed; the concavity of v is pinned down by fundamental information limits. In policy terms, data economies of scale saturate exactly when further capacity gains are information-theoretically expensive.

Table 2: Types of predictions and policy implications

| Type of prediction | Data abundant? | Tails thick? | Monopoly Remedy |
|---|----------------|--------------|--------------------|
| Genomic risk prediction (health) | No | Yes | Access regulation |
| Clinical decision support for rare diseases | No | Yes | Access regulation |
| Credit scoring / SME default probability | No | Yes | Access regulation |
| Fraud / AML detection | No | Yes | Access regulation |
| Industrial predictive maintenance (OEM IoT) | No | Yes | Access regulation |
| Smart grid anomaly detection (critical infra) | No | Yes | Access regulation |
| Autonomous driving safety edge cases | Yes | Yes | Hybrid |
| Weather nowcasting for extremes | Yes | Yes | Hybrid |
| E-commerce CTR / product recommendation | Yes | No | Competition policy |
| Targeted Ads | Yes | No | Competition policy |
| Media streaming recommendation | Yes | No | Competition policy |
| Web search ranking | Yes | No | Competition policy |

²²This integral form follows immediately from Gallager, *Information Theory and Reliable Communication*, 1968, Ch. 8, where parallel Gaussian sub-channels are treated.