

# Experiential and Social Learning in Firms: The Case of Hydraulic Fracturing in the Bakken Shale

## JOB MARKET PAPER

Thomas R. Covert\*

November 15, 2013

### Abstract

Learning how to utilize new technologies is a key step in innovation, yet little is known about how firms actually learn. This paper examines firms' learning behavior using data on their operational choices, profits, and information sets. I study companies using hydraulic fracturing in North Dakota's Bakken Shale formation, where firms must learn the relationship between fracking input use and oil production. Using a new dataset that covers every well since the introduction of fracking to this formation, I find that firms made more profitable input choices over time, but did so slowly and incompletely, only capturing 67% of possible profits from fracking at the end of 2011. To understand what factors may have limited learning, I estimate a model of fracking input use in the presence of technology uncertainty. Firms are more likely to make fracking input choices with higher expected profits and lower standard deviation of profits, consistent with passive learning but not active experimentation. Most firms over-weight their own information relative to observable information generated by others. I use these model estimates to measure the impact of information availability regulations on learning.

## 1 Introduction

New technologies are important contributors to economic growth<sup>1</sup>, but little is known about how firms learn to profitably use them. While there is longstanding evidence that firms learn from their own experiences (learning-by-doing), and from others (social learning), the specific actions that firms actually take in learning are not well understood. Models of learning predict that firms efficiently analyze information about new technologies, invest in experiments to create new information, and incorporate information generated by other firms.<sup>2</sup> However, to test these models, it is necessary to observe data on the information that firms have, which is difficult to acquire. This paper empirically tests predictions of learning behavior for the first time, using data on oil companies that employ hydraulic fracturing (fracking) in the North Dakota Bakken Shale. The data covers operational choices, profits, and measures of the information firms had when making choices. The oil companies in this data learn to use fracking more profitably over time, but are slow to respond to new information, rarely experiment and ignore much of the data provided by their competitors.

---

\*Harvard Business School and Harvard University Department of Economics; tcovert@hbs.edu. I am grateful to Paul Asquith, Bharat Anand, Greg Lewis, Ariel Pakes and Parag Pathak for their guidance and encouragement. I thank Michael Luca, Thomas G. Wollman, Richard Sweeney, Bryce Millet-Steinberg, Stephanie Hurder, Alex Peysakhovich, Evan Herrnstatt, Joseph Shapiro, Hugh Daigle, Heath Flowers and the Harvard Industrial Organization, Environment Economics, and Work-in-progress lunches for their helpful comments and discussions. Funding from the Harvard Business School Doctoral Programs, a Harvard University Dissertation Completion Fellowship and the Sandra Ohrn Family Foundation is gratefully acknowledged.

<sup>1</sup>See, for example, Arrow (1962), Romer (1986) and Kogan et al. (2012)

<sup>2</sup>See Aghion et al. (1991) in the single agent context and Bolton and Harris (1999) in the multi-agent context.

Fracking is a useful context to study learning behavior in firms. The profit maximizing choice of fracking inputs may vary across drilling locations in unpredictable ways, so firms must empirically learn this relationship over time and change their behavior accordingly. In North Dakota, firms have a wealth of information to learn from because regulators collect and publicly disseminate unusually detailed, well-specific information about oil production and fracking input choices. Moreover, regulators delay dissemination until 6 months after a well is fracked, making it possible to measure differences in knowledge about fracking across firms. The industry is not concentrated, which motivates studying learning as a single agent problem. During the time period I study, there are 70 active firms, the market share of the largest firm is only 13% and the combined share of the five largest firms is under 50%. The two main inputs to fracking, sand and water, are commodities, as is the output of fracking, crude oil. The unique data and environment make fracking in the Bakken shale an unusually compelling setting for studying learning in firms. The stakes in fracking are large. Using a production function, I estimate that the average NPV of profits per well for actual fracking choices is about \$12.7 million, while the average profit for each well's most profitable choice is \$23.6 million. Since the regulator in North Dakota expects that 40,000 wells will eventually be fracked over the next 18 years, the potential for lost profits from inefficient learning is substantial.<sup>3</sup>

Learning-by-doing and social learning are both important in this context. Between 2005 and 2006, the average well is fracked by a firm that had fracked only a single well before. By 2011, the average well is fracked by a firm that had previously fracked 117 wells. Thus, firms have an increasing amount of learning-by-doing experience to draw from. However, North Dakota's disclosure laws make it possible for firms to access their competitors' data. Between 2005 and 2006, the average well is fracked by a firm that can observe 10 wells previously fracked by other firms, a number which rises to 1,783 in 2011. As a result, most of the information firms have comes from others, and firms have the ability to socially learn.

The data I collect from the regulator in North Dakota is well suited to estimate the relationship between location, fracking, and oil production. I observe the complete operating history of every firm and every well they frack in the Bakken Shale between January 2005 and December 2011 (70 firms and 2,699 wells). The data contains precise measurements of a well's production, location and fracking inputs. Because I observe information on inputs and outputs for all wells that were ever fracked, I can avoid the omitted variable and survivorship biases that are common in production function estimation.

Using the data I collect, I semi-parametrically estimate a production function for fracking to establish what firms need to learn. The amount of oil in the ground and the sensitivity of its production to fracking both vary over space, a result that is consistent with geological theory and data. Estimates made using subsets of the data that were available to firms when they were fracking have quantitatively similar results, suggesting that firms could have used this data to make informed fracking decisions. The estimated production function fits the data well and is stable across cross-validation robustness tests.

I use this production function to measure how quickly firms learn. Wells fracked in 2005 capture only 21% of the profits that optimally fracked wells would have produced. However, profit capture grows almost monotonically over time, with firms capturing 67% of maximal profits in 2011. This growth is driven by improved fracking input choices, with firms gradually increasing their use of sand and water towards optimal levels over time. I interpret this upward trend in the profitability of fracking input choices as evidence for learning.

Existing research measures learning from upward trends in *productivity*, or residual production that is not explained by input choices. I test for productivity based learning by analyzing the

---

<sup>3</sup>See <https://www.dmr.nd.gov/oilgas/presentations/NDOGCPC091013.pdf>

growth of estimated production function residuals over time. Wells fracked in 2011 are 34% more productive than wells fracked in 2005, suggesting some role for productivity-driven learning. However, the majority of the growth in productivity occurs by 2008, and there is no statistically significant difference in productivity between 2008 and 2011. This contrasts with the fraction of profits captured, which increases monotonically over time, and from 44% to 67% between 2008 and 2011. Thus, during 2008-2011, when 95% of wells in my data are fracked, there is little productivity growth, even though there is substantial growth in the fraction of profits captured. These results help clarify the difference between models of learning in which knowledge is a direct input in the production function, and a model of learning about the production function itself.

To see if firms are using their information to make better fracking choices over time, I estimate an *ex ante* production function, using the subset of the data that firms had when they were making choices. I use these estimates to compute *ex ante* profits. Though firms capture nearly 80% of *ex ante* optimal profits in 2007, they capture only 67% in 2011. The fraction of *ex ante* profits falls because initial fracking input choices are close to the (then) estimated optimal levels, but optimal levels subsequently change more quickly than choices do.

Theory predicts that firms may sacrifice estimated profits in the current period by experimenting in order to generate information for the future. To test if experimenting behavior can rationalize the decline in the fraction of estimated *ex ante* optimal profits captured, I estimate a simple model of fracking input choice under technology uncertainty. In this model, firms have preferences over the expectation and standard deviation of their *ex ante* estimates of profits for a fracking input choice. If firms are experimenting, they should be empirically more likely to choose inputs with higher standard deviations of profit. I do not find support for this theory. Firms are more likely to select fracking designs with higher expected profits and *lower* standard deviation of profits. Firms are indifferent between a \$0.62-\$1.11 increase in expectation of profits and a \$1 reduction in the standard deviation of profits.

My calculation of the expectation and standard deviation of profits assumes that firms equally learn from their own and others' experiences. However, firms may treat the social portion of their data differently than the data they directly experience, and in the process form different estimates of profits than what I calculate. To account for this possibility, I modify my fracking input choice model to allow for weighted estimates. I use this model and data on firms' choices to estimate the weight they place on their own experiences relative to their competitors' experiences. Most firms place more weight on their own experiences than their competitors' experiences. Even after controlling for weighted estimates, firms still prefer fracking choices with lower standard deviations and higher means.

I am currently using the fracking input choice model to study the role of information disclosure rules in learning. Regulators use "well confidentiality" laws to incentivize exploration of new drilling areas. In North Dakota, firms that discover productive formations have 6 months to acquire nearby mineral rights before their discovery is made public. However, this delay in public disclosure of may affect learning by reducing the information firms can learn from. To measure the impact that delays have on learning, I compute the choices firms would have made under alternative well confidentiality rules. First, I consider the repeal of North Dakota's well confidentiality law, which would cause all firms to have the same information at the same time. Second, I consider an increase in well confidentiality from 6 months to 1 year, the level of well confidentiality available to firms in neighboring Montana. These counterfactual analyses will be of interest to policy makers in states and countries that are beginning to regulate fracking for the first time.

This paper finds that firms are reluctant to experiment and ignore valuable data generated by their competitors. These firms are not unsophisticated or under-incentivized. They have access to capital markets, are managed by executives with engineering and business education and are

the primary equity holders in the wells they frack. These findings stand in contrast to theories of efficient learning behavior by rational agents, which predict that firms will take experimental risk and learn from all the information they have.

Aside from its usefulness as a laboratory to study learning, fracking is interesting from a public policy perspective due to its contribution to national oil production and the current debate about its impact on local economies and the environment. The US EIA reports that fracking has caused national oil production to grow nearly 30% since 2008, reversing almost two decades of declines.<sup>4</sup> There is early evidence that these resource booms have affected housing prices<sup>5</sup> and local banking markets.<sup>6</sup> However, there are growing concerns about the potential for fracking to negatively affect the quantity and quality of local ground water supplies<sup>7</sup>, which the US EPA is currently studying.<sup>8</sup> In response to these concerns, federal regulators have proposed significant increases to disclosure requirements for fracking operations.<sup>9</sup> Though this push for increased transparency around fracking is driven by environmental concerns, new disclosure regulations may also have an impact on learning by increasing the availability of data.

Finally, the Bakken Shale unlikely to be the last oil and gas formation where fracking and the learning it requires play an important role. Fracking is currently in use in the Eagle Ford and Barnett Shales in Texas, the Woodford Shale in Oklahoma, and several locations in Canada. International oil companies are now developing shale resources in Argentina, Poland and China. The results of this paper may be useful to both policy makers and oil & gas companies alike in regulating access to information and understanding the benefits of more efficient learning behavior.

## 1.1 Related literature

Firms in many industries and time periods have become more productive by learning from their own experiences. Researchers studying the manufacturing of World War II ships (Thornton and Thompson (2001)), aircraft (Benkard (2000)) and automobiles (Levitt et al. (2012)) have documented an important empirical regularity: with the same inputs, firms are able to produce more output as they accumulate experience in production.<sup>10</sup> That is, they learn by doing (LBD). The LBD literature establishes that productivity is correlated with experience, suggesting that the knowledge embedded in this experience is a direct input to the production function. Changes over time in capital, labor and materials are interpreted as a profit-maximizing response to increases in experience. In this paper, I instead assume that the production technology itself is initially unknown and that experience has no direct impact on production. As firms accumulate experience in fracking, they acquire more data about the fracking production function, perform inference on this data, and make more profitable input choices on the basis of their inference. This is similar to the approach taken by Conley and Udry (2010).

Economic theory predicts that when firms are learning about a new technology, they face a tradeoff between “exploration” and “exploitation”. Firms may actively learn by exploring, or experimenting, with fracking input choices that have highly uncertain profits or passively learn by exploiting choices with high expected profits. Except in the simplest theory models, the optimal amount of exploration and exploitation is a challenging problem to solve. However, most models of learning predict that forward-looking firms will always do some exploring. In the single agent

---

<sup>4</sup><http://www.eia.gov/todayinenergy/detail.cfm?id=13251>

<sup>5</sup>Muehlenbachs et al. (2012) find that housing prices increase after the introduction of fracking to a community, except for houses that depend on groundwater.

<sup>6</sup>See Gilje (2012)

<sup>7</sup>See Vidic et al. (2013) for an overview

<sup>8</sup>See <http://www2.epa.gov/hfstudy>

<sup>9</sup>See Deutsch (2011).

<sup>10</sup>This phenomenon has also been observed by Anand and Khanna (2000) in the corporate strategy setting.

context, Aghion et al. (1991) show that in many settings, forward-looking firms will always do some exploration. Bolton and Harris (1999) find a similar result in the multi-agent context. Wieland (2000) employs computational methods to characterize the costs and benefits of exploration, finding that firms who only exploit can get stuck, and repeatedly choose suboptimal actions. To my knowledge, this paper is the first to empirically measure the amount of exploration that firms do in learning situations.

This paper adds to a wide literature documenting the existence and importance of social learning between firms. Most of this evidence is in agricultural settings. Ryan and Gross (1943), Griliches (1957) and Foster and Rosenzweig (1995) demonstrate that farmers learn about the benefits of adopting new technologies from the experiences of their neighbors. Conley and Udry (2010) show that farmers in Ghana learn about the efficient use of fertilizer from other farmers in their social networks, demonstrating that social learning in agriculture is not limited to the adoption decision. Social learning has also been observed in manufacturing. During the construction of WWII ships, Thornton and Thompson (2001) find that firms benefited from accumulated experience by other firms. Similarly, Stoyanov and Zubanov (2012) find evidence that firms in Denmark experience increased productivity after hiring workers away from more productive competitors.

Finally, this paper is complementary to the existing literature on learning behavior by oil and gas companies. Levitt (2011) shows that the observed temporal and spatial patterns of the oil exploration process match the predictions of a forward-looking learning model. Kellogg (2011) studies repeated interactions between oil companies and their drilling service providers, finding that they jointly learn to be more productive in drilling as they accumulate shared operating experience.

The remainder of the paper is as follows. In section 2, I provide institutional background on fracking in North Dakota and describe the data I have on operations, results and information sets. Next, in section 3, I estimate a production function model of fracking and evaluate its ability to predict oil production. In section 4, I use the production function estimates to test if firms learned to make more profitable fracking choices over time. In section 5, I specify and estimate the model of fracking input choice under technology uncertainty. Then, in section 6, I use this model to evaluate the impact of learning strategies and information availability regulations on oil production. Finally, I conclude in section 7.

## 2 Institutional Background and Data

### 2.1 Fracking and US oil production

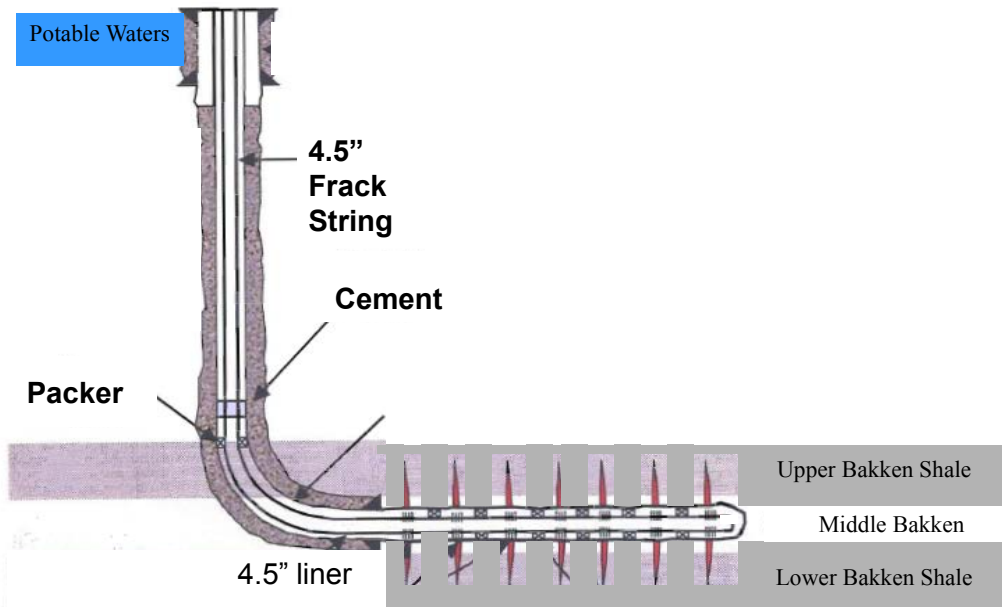
The hydraulic fracturing of shale formations, like the Bakken, has had a profound impact on the fortunes of energy producing states and the US as a whole. In 2009, the US Energy Information Administration reported that national oil production grew 6.8% year-over-year, the first increase in over two decades.<sup>11</sup> This trend has continued and between 2009 and 2012, national oil production has increased 21.7%. Three states represent the majority of this reversal of fortunes: Texas, Oklahoma and North Dakota. This paper focuses on what has happened in North Dakota.

In March 2012, North Dakota surpassed Alaska to become the second most prolific oil producing state in the US, after Texas. Between January 2005 and July 2013, oil production in North Dakota increased from 93,000 barrels (bbl) per day to 874,000 bbl per day. During the same time period, total US oil production increased from 5.63 million bbl per day to 7.48 million bbl per day, meaning that increased production in North Dakota amounted to 42% of the net increase in total production. Though production increased in Texas and Oklahoma as well, it is striking that North Dakota went

---

<sup>11</sup>See the EIA Annual Energy Review, 2009. <http://www.eia.gov/totalenergy/data/annual/archive/038409.pdf>

Figure 1: Diagram of a Hydraulically Fractured Bakken Shale well. Adapted from Hicks (2012)



from producing less than 2% of national oil production to almost 12% in the span of 8 years.<sup>12</sup> This vast expansion in North Dakotan oil production coincides with the introduction of fracking to the Bakken Shale formation.

## 2.2 The Bakken Shale and Hydraulic Fracturing

The Bakken Shale spans 200,000 square miles in North Dakota, Montana and Saskatchewan.<sup>13</sup> It lies 10,000 feet underground and contains 3 distinct layers: the Upper Bakken Shale, the Middle Bakken (which is not a shale), and the Lower Bakken Shale. The US Geological Survey estimates that the Upper and Lower Shales together contain 4.6 billion bbl of recoverable oil.<sup>14</sup> Though the Middle Bakken does not contain any oil of its own, firms typically drill horizontally through it and use hydraulic fracturing, or “fracking”, to make contact with the oil bearing shales above and below, as shown in Figure 1.

Fracking is the process of pumping a mix of water, sand and chemicals into a well at high pressures. The high pressure of the mix fractures the surrounding rock and the sand in the mix props those fractures open.<sup>15</sup> The fractures created by fracking the Middle Bakken radiate outwards into the Upper and Lower Bakken Shales, as shown in Figure 1. These fractures both serve as a conduit between the Middle Bakken and the Upper and Lower Bakken Shales, and also increase the permeability of the Upper and Lower Bakken Shales.

<sup>12</sup>Texas also experienced production significant production increases during that same time period, though from a much higher base level (from 1.08 million bbl per day to 2.62 million bbl per day, a 143% increase). Much of this increase can also be attributed to the technology changes described here. Operators applied fracking technology successfully to the Eagle Ford, Permian and Barnett shales.

<sup>13</sup>See Gaswirth (2013)

<sup>14</sup>See Gaswirth (2013)

<sup>15</sup>Chemicals reduce mineral scaling, inhibit bacterial growth and increase the buoyancy of sand in the mixture. See <http://www.fracfocus.org> for an overview.

Permeability is a geological measure of the ease at which oil naturally flows through rock. The Upper and Lower Bakken Shales are unusually impermeable, making it impossible for the oil they contain to naturally reach a wellbore drilled through the Middle Bakken. Without fracking, wells drilled into the Middle Bakken will not produce profitable quantities of oil.<sup>16</sup> After fracking, oil inside the Lower and Upper Bakken Shales can travel through the new fractures into the wellbore in the Middle Bakken.

Firms choose how much water and sand to use in fracking and this choice can have a large impact on the profits generated by a well. Wells fracked with more sand and water may produce more oil than wells fracked with less, but fracking is expensive, and water and sand represent the bulk of this expense. The reported costs of fracking in 2013 range from \$2-5 million per well, out of a total well costs of \$9 million.<sup>17</sup> To maximize profits, firms must balance the benefits of sand and water use in fracking with their costs. This requires firms to understand the relationship between oil production and fracking inputs, and it is unlikely that firms initially knew this relationship. The first Bakken wells to be developed with fracking were not drilled until 2005, and at the time, the firms developing those wells had limited experience in fracking shale formations.<sup>18</sup> Without prior experience, firms had to learn how to use fracking by doing it.

Today, there is a growing literature about best practices in fracking. Petroleum engineers have found that in many cases, wells fracked with more water and sand are more productive than similar wells with less aggressive fracking treatments.<sup>19</sup> However, there is evidence that the relationship between oil production and fracking inputs is not necessarily monotonic and that it varies over drilling locations.<sup>20</sup> While neither of these results were publicly available to firms during the time period I study, they suggest that firms faced a complicated learning problem.

### 2.3 The information environment in North Dakota

Firms in North Dakota can learn about the relationship between oil production, location, and fracking inputs from the past experiences of other firms. After a firm fracks a well, the oil and gas regulator in North Dakota requires the firm to submit a well completion report, detailing the well's horizontal length, location and fracking inputs. Additionally, the regulator and tax authorities require the firm to submit audited production records on a monthly basis. The regulator publishes this information on the internet. This public disclosure infrastructure makes it easy for firms to learn detailed information about every previously fracked well in the state, including information about wells that they took no part in developing.

There is a delay between when firms submit well completion reports and when the regulator makes them public, due to North Dakota's well confidentiality laws. If a firm requests confidential treatment for a well, the regulator delays the publication of that well's completion report for 6 months. Production information is not subject to well confidentiality requests, and is publicly available within 2 months. Well confidentiality creates differences across firms in what wells they can learn from at each point in time, as the operating firm of a well has a temporary knowledge

---

<sup>16</sup>See Hicks (2012)

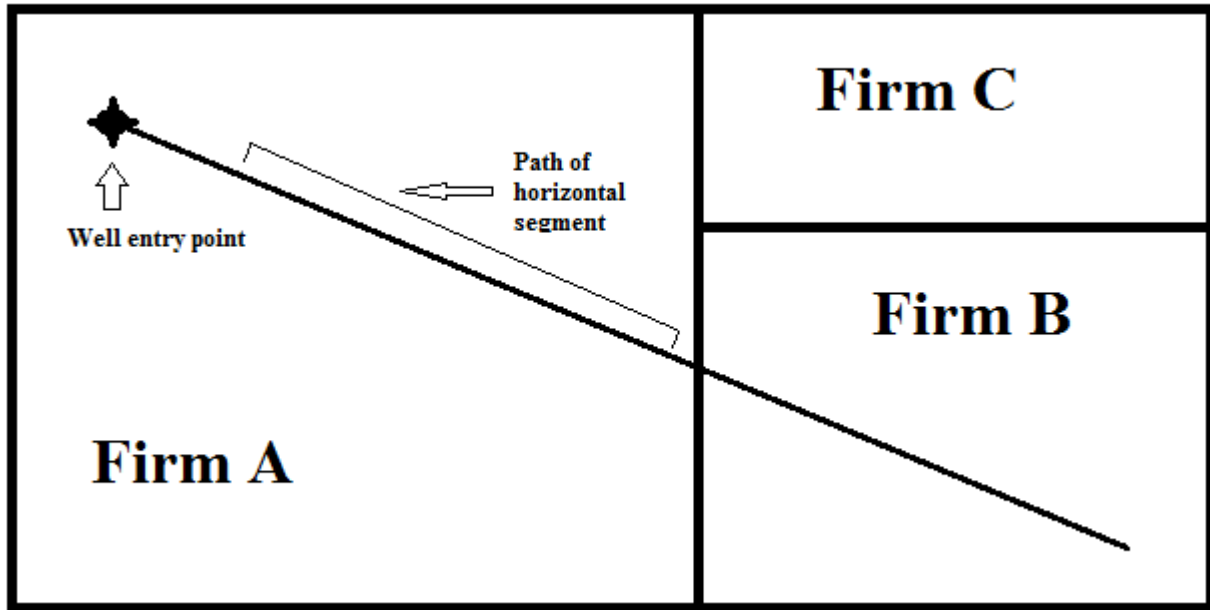
<sup>17</sup>See Hicks (2012)

<sup>18</sup>Fracking was first successfully used in shale formations in the 1990s. Under the hunch that permeability issues could eventually be resolved through the use of fracking, Mitchell Energy worked for years on its own and with the help of the US Department of Energy to determine how to apply fracking technology to the Barnett shale in Texas. They succeeded in 1997. See Michael Shellenberger and Jenkins (2012). Two firms active in North Dakota, EOG and XTO, were active in the Barnett as well. However, the Barnett Shale is different from the Bakken. Barnett wells are drilled directly into the shale layer, and produce natural gas instead of oil. It is unlikely that any knowledge that these firms may have had about fracking in the Barnett was useful in the Bakken.

<sup>19</sup>See Shelley et al. (2012)

<sup>20</sup>See Baihly et al. (2012)

Figure 2: Diagram of a hypothetical spacing unit



advantage over other firms.

While well confidentiality generates differences in knowledge between firms, the complicated ownership structure of mineral rights in a well reduces these differences. The mineral rights for a well are often owned by many separate firms. Every firm that owns mineral rights in the area spanned by a well is entitled to pay a share of the capital expenditures needed to develop the well in exchange for a share of the revenue generated by the well. The firm with the largest mineral rights claim in a well is called the “operator”, and it retains all control rights, including the choice of the well’s fracking inputs. The remaining owners of mineral rights are called “non-operating participants”. Figure 2 depicts a hypothetical ownership situation for a well in the Bakken. The land spanned by the well is a 2 mile by 1 mile rectangle, called a “spacing unit”. Within this spacing unit, Firm A has the largest mineral rights claim, followed by firms B and C. The wellhead enters the ground in A’s claim and the horizontal segment passes through B’s claim. Though the well does not directly pass through C’s claim, it is close enough to C’s claim that it may be drawing oil from the claim. While A retains control rights, B and C must pay their respective share of capital expenditures.<sup>21</sup>

Non-operating participants have immediate access to a well’s completion report.<sup>22</sup> This means that non-operating participants in a well are not subject to well confidentiality rules and thus observe information regarding a well before the public does.

<sup>21</sup>Firms can choose to opt out of a spacing unit, but that does not allow them to operate another well within the spacing unit, so opt outs are rare.

<sup>22</sup>See Larsen (2011)



Table 1: Summary Statistics

Variable	Mean	Std. Dev	P25	P50	P75	N
lbs sand per foot	265.02	138.68	158.27	264.53	378.66	2,699
gals water per foot	188.87	110.73	100.31	181.70	249.52	2,699
horizontal feet in length	8,040	2,138	5,600	9,135	9,518	2,699
avg producing days per month	26.80	2.99	25.90	27.56	28.67	2,699
oil production per foot in first year	10.86	8.95	5.38	8.39	12.99	2,699
# non-operating participants	3.00	2.50	1.00	3.00	4.00	2,699
# past wells fracked by operator	80	82	16	49	125	2,699
# past wells fracked by others	1,089	658	511	1,062	1,698	2,699

## 2.4 Data

### 2.4.1 Well characteristics and production history

I have collected operating and production data for every well targeting the Bakken shale formation in North Dakota that was fracked between January 1, 2005 and December 31, 2011. This data is reported by oil companies to the North Dakota Industrial Commission (NDIC), and the NDIC publishes their submissions on the internet. For each well  $i$ , I observe the location of its wellhead in latitude  $lat_i$  and longitude  $lon_i$  coordinates, its horizontal length  $H_i$ , the mass of sand  $S_i$  and volume of water  $W_i$  per foot of horizontal length used in fracking operations and the identity of the operating firm  $f_i$ . Additionally, I observe oil production  $Y_{it}$  for well  $i$  in its  $t$ -th month of existence and the number of days  $D_{it}$  during that month that the well was actually producing. Let  $X_{it}$  denote the set  $(H_i, f_i, D_{it})$  and let  $Z_i$  denote the set  $(lat_i, lon_i, S_i, W_i)$ . Then the dataset  $(Y_{it}, X_{it}, Z_i)$  has a panel structure, where  $i$  indexes wells and  $t$  indexes well-specific timing. Though I only study wells fracked during 2005-2011, I have production data through February 2013, making it possible to study the performance of all wells for at least a year. While the production history is reported electronically on the NDIC website, the static well characteristics are stored in PDF format, so much of this dataset was entered into the computer manually. Using a well's latitude and longitude, I also compute the "township" the wellhead lies in. Townships are 6 mile by 6 mile squares, defined by the US Geological Survey and are a standard measure of location in the oil & gas business. There are 272 townships with Bakken wells during 2005-2011. I have also collected the geographic boundaries of the spacing units for every well. This data comes from various portions of the NDIC website.

Table 1 shows the distribution of well characteristics and production statistics. There is substantial variation across wells in both fracking input use and oil production. The 75th percentiles of sand, water and oil production are more than double their respective 25th percentiles. This variation will be important later on in estimating the relationship between sand and water use in fracking and oil production. Most wells have horizontal segments that are 9,000 feet or longer. The length of a well's horizontal segment is determined by the size of its spacing unit. Though not shown in the table, approximately 75% of wells have rectangular spacing units that are two miles wide and one mile tall. The remaining 25% have 1 mile square spacing units. The average well produces almost 11 bbl per foot of horizontal length in its first year. Since the price of oil averaged \$76 per bbl during 2005-2011, the value of production in the first year for the average well is worth \$6.6 million. Most wells tend to produce on the majority of days during a month, and though not shown in the table, only 93 wells have fewer than 20 average producing days.

The bottom rows of Table 1 show the distribution of non-operating participants and past experience across wells. In the average well, 3 other firms obtain knowledge about a well at the

Table 2: Summary statistics by year

		2005	2006	2007	2008	2009	2010	2011
	# wells fracked	10	20	94	352	463	691	1,069
	# active townships	9	17	37	102	132	179	231
	# active firms	5	11	17	28	34	47	49
Sand	Average	94.50	136.88	134.64	180.00	212.75	308.82	302.79
	Std. Dev	22.01	152.43	143.15	146.79	145.32	121.68	110.85
Water	Average	49.53	64.29	95.67	108.28	137.08	215.14	232.68
	Std. Dev	25.03	61.87	83.72	59.90	88.36	99.88	111.28
Length	Average	6,883	6,062	7,017	7,283	7,238	8,006	8,795
	Std. Dev	1,679	2,001	2,048	2,233	2,316	2,144	1,715
Oil	Average	3.08	4.85	10.76	13.41	11.55	11.15	9.73
	Std. Dev	1.94	7.59	13.72	15.16	9.83	6.72	5.78

same time as the well's operator. The average well is fracked by a firm that has previously fracked 80 of its own wells, and can observe the data on 1,089 wells fracked by others.

Table 2 shows the distribution of fracking activity over time. The number of wells fracked and the number of active townships and firms all increase over time. More than 65% of all wells are fracked during the last two years, and in 2011, wells are fracked in 85% of townships by 70% of all firms. Over time, firms frack longer wells, using more sand and more water. Wells fracked in 2011 use more than three times as much sand and four times as much water, on average, as wells fracked in 2005. However, average oil production does not rise monotonically, reaching its peak in 2008 and then falling thereafter. The set of drilling locations expands considerably over time, which may confound measurement of the relationship between input use and oil production.

To control for this, I compare relative oil production within a township across wells with differing levels of sand and water use. In each township, I subtract the average levels of oil production and input use from actual production and input use. Then, I add back the overall average levels, creating normalized measures of production and input use, net of township fixed effects. Table 3 reports average normalized oil production per foot by quintiles of normalized sand and water use. Across both sand and water use, the highest input levels are associated with higher oil production. For every quintile of water use (columns), the top quintile of sand use has higher production than the bottom quintile. For all but the second quintile of sand use (rows), the top quintile of water use has higher production than the bottom quintile. Thus the data supports the idea that sand and water use affect oil production.

Current petroleum engineering research suggests that there is spatial heterogeneity in the relationship between fracking inputs and oil production. I check for this by estimating a simple Cobb-Douglas production function for fracking, with and without township fixed effects. I regress the log of oil production per foot of horizontal length on the well's log sand use and log water use:

$$\log \text{ oil per foot}_i = \alpha_0 + \alpha_S \log S_i + \alpha_W \log W_i + \text{township}_i + \epsilon_i$$

Table 4 reports coefficient estimates for this regression. The first column shows estimates without fixed effects, and the second column shows estimates with fixed effects. Consistent with the results in Table 3, higher sand and water use are associated with higher production. This is true with and without fixed effects. However, the inclusion of township fixed effects decreases the coefficient on sand use and increases the coefficient on water use, suggesting the existence of spatial heterogeneity in oil production and the possibility that firms make different input choices in different locations.

Table 3: Average first year’s oil production per foot of horizontal length by quintiles of sand and water use. Net of township fixed effects. Standard errors in parentheses.

		Quintiles of Water Use				
		First	Second	Third	Fourth	Fifth
Quintiles of Sand Use	First	8.09 (0.33)	8.27 (0.44)	6.77 (0.73)	8.82 (2.20)	10.16 (0.81)
	Second	9.53 (0.37)	10.50 (0.35)	9.27 (0.33)	10.50 (0.56)	9.37 (1.37)
	Third	10.25 (0.52)	11.51 (0.36)	10.91 (0.29)	10.56 (0.35)	10.81 (0.76)
	Fourth	10.71 (0.52)	10.48 (0.58)	13.24 (0.55)	11.46 (0.40)	11.87 (0.48)
	Fifth	10.80 (1.06)	12.24 (0.83)	13.37 (0.98)	13.19 (0.52)	13.85 (0.37)

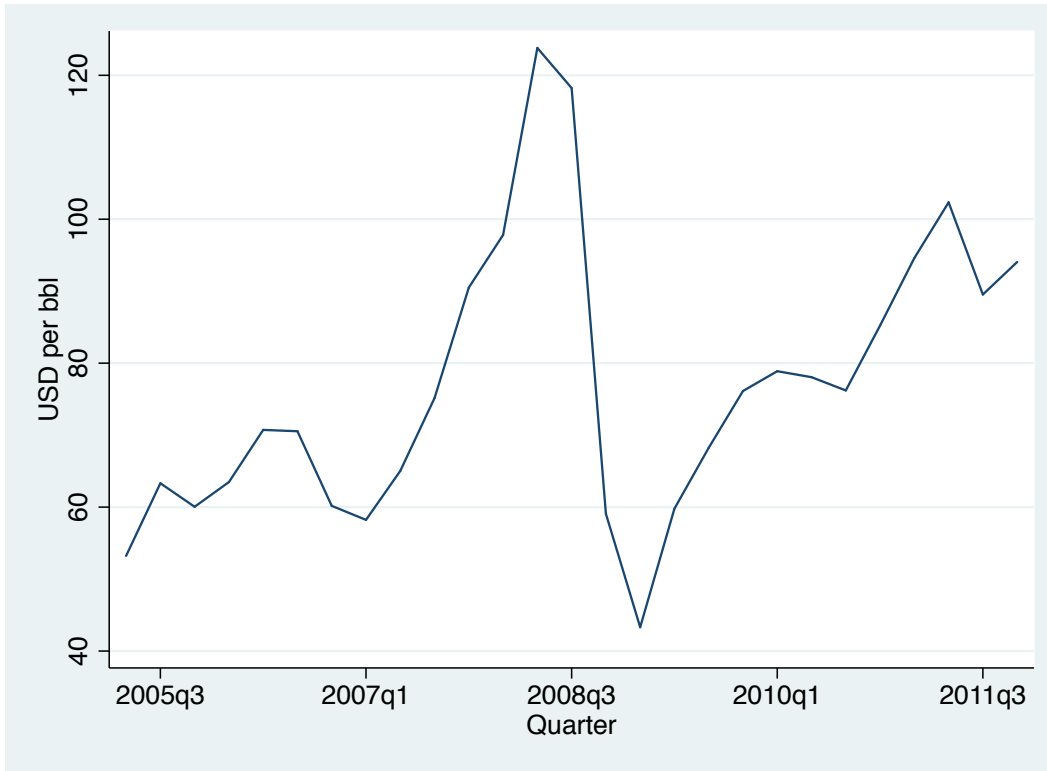
Table 4: Spatial Heterogeneity in the Relationship Between Sand, Water and Oil Production

	(1) Log Oil per foot	(2) Log Oil per foot
$\alpha_S$	0.352 (0.0211)	0.208 (0.0183)
$\alpha_W$	0.0512 (0.0228)	0.137 (0.0185)
$\alpha_0$	-0.0280 (0.104)	0.319 (0.0948)
Township FE		X
$N$	2,698	2,698
$R^2$	0.159	0.618

Standard errors in parentheses. OLS estimates of

$$\log \text{ oil per foot}_i = \alpha_0 + \alpha_S \log S_i + \alpha_W \log W_i + \text{township}_i + \epsilon_i$$

Figure 3: Quarterly Average Cushing Oil Prices



#### 2.4.2 Oil prices and fracking costs

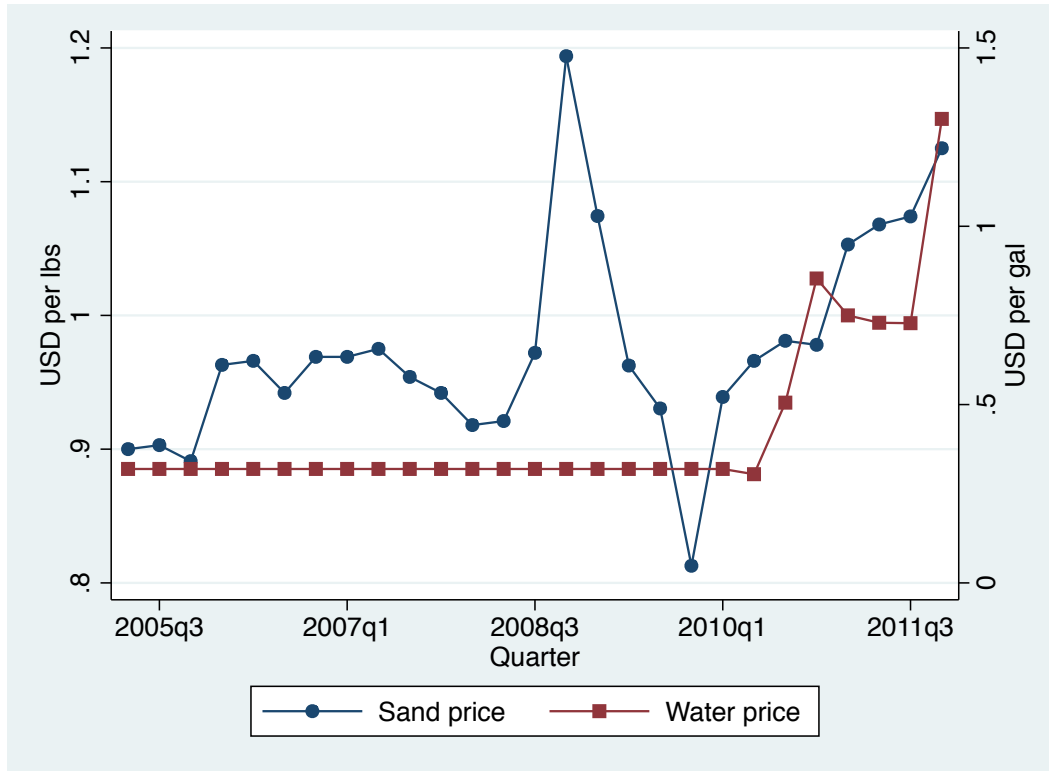
I collect the daily spot prices for West Texas Intermediate crude oil at the Cushing, Oklahoma oil trading hub from the US Energy Information Administration. The Cushing price is the reference price for oil futures traded on the NYMEX commodity exchange, and the Cushing hub is connected to North Dakota through the Keystone and Enbridge pipeline systems. Figure 3 plots quarterly average oil prices at the Cushing hub. Between 2005-2011, there was a boom and bust in oil prices, with prices climbing from approximately \$60 per bbl in early 2007, reaching more than \$120 per bbl in mid 2008 and falling to \$45 per bbl in early 2009. In 2010-2011, when more than 65% of the wells are fracked, oil prices average \$87 per bbl.

I use data from Spears & Associates to compute the costs of drilling and fracking in North Dakota. Spears & Associates surveys independent engineers in North Dakota quarterly, asking them to estimate the cost of a reference well. The cost estimates are divided into 14 categories, of which 4 are fracking related and 10 are drilling related. The data is separately available for a vertical reference well design, which begins in the first quarter of 2008 and a horizontal reference well design, which begins in the first quarter of 2010. The vertical reference design does not include a fracking treatment. The characteristics of the reference wells stay constant over time, so the changes in estimated costs are due to changes in prices, not quantities.

I assume that the drilling related costs for the vertical reference well represent fixed costs of drilling and that the drilling costs for the horizontal reference well in excess of drilling costs of the vertical well represent variable costs of drilling. I divide the total variable cost of drilling by 9,000 feet (the horizontal length of the horizontal reference well) to compute the per-foot cost of horizontal drilling.

I assume that the fracking related costs for the horizontal reference well are equal to the cost

Figure 4: Estimated Sand and Water Prices



of sand times sand use plus the cost of water times water use. Unfortunately, I do not observe the characteristics of the fracking design for the horizontal reference well, so I assume that it has the average characteristics of the wells fracked in the first quarter of 2010 with horizontal lengths of 7,500 feet or more: 238 lbs sand per foot and 173 gals water per foot.<sup>23</sup> This assumption by itself is not enough to infer the marginal costs of both water and sand, since there are two prices to infer from one cost observation, so I also assume that the cost of using an additional pound of sand in fracking is equal to 3 times the wholesale cost of ceramic beads, a substitute for sand in fracking. I collect the wholesale cost of ceramic beads from the SEC filings of CARBO Ceramics, the world’s largest manufacturer of ceramic beads. Finally, I compute the cost of using an additional gallon of water in fracking during quarter  $t$  as:

$$p_{W,t} = \frac{\text{total fracking cost}_t - 9,000 \times 238 \times 3 \times p_{\text{ceramic},t}}{9,000 \times 173}$$

Since the horizontal reference well design costs are only available starting in the first quarter of 2010, I assume that variable drilling costs and water costs remain constant at the level in the first quarter of 2010 between the first quarter 2005 and the fourth quarter 2009. Though I have assumed away any variation in these prices between 2005-2009, recall that 65% of wells are fracked in 2010-2011. Thus the biases caused by this assumption affect a fewer than half the wells.

Figure 4 plots my estimates of sand and water prices by quarter.

<sup>23</sup>Broadly speaking, there are two kinds of horizontal wells: those with a horizontal length of approximately 4,500 feet, situated on 1 mile by 1 mile spacing units, and those with a horizontal length of approximately 9,000 feet, situated on 2 mile by 1 mile spacing units. The distribution of horizontal lengths in my data has modes at approximately 4,500 and 9,500 feet, with the least likely horizontal length in between the modes at 7,500 feet.

Table 5: Wells completed by the eight largest firms, by location, time and well characteristics

Firm	North Dakota		Outside North Dakota		
	2005-2011	1995-2004		2005-2011	
	Bakken Shale	Vertical	Horizontal	Vertical	Horizontal
Brigham	113	161	0	93	0
Burlington	105	3,826	26	2,792	532
Continental Resources	313	597	3	657	167
EOG	354	4,659	91	6,566	2,914
Hess	165	639	2	219	15
Marathon	223	2,221	4	813	87
Whiting	247	131	0	1,150	11
XTO	101	2,349	53	7,749	2,801
Rest of industry	1,078				

### 2.4.3 Information Sets

At time  $t$ , firm  $f$  can learn about fracking from the union of three sets of wells. First,  $f$  can observe all wells that the regulator has made public by time  $t$ . This public knowledge includes wells that  $f$  operated and wells that other firms operated. Second,  $f$  can observe its own wells which are not yet public knowledge. Third,  $f$  can observe other firms' wells in which it is a non-operating participant. I can compute the first two sets of information from well completion reports alone. To compute the third set, I must identify the mineral rights owners in each well's spacing unit.

I collect mineral rights lease transaction data from DrillingInfo.com, a firm which digitally records the universe of mineral rights transactions filed in county deed registries. These leases are often between a surface owner, such as a farmer, and an intermediary lease broker operating on behalf of an oil company. Once the broker acquires a lease, it assigns this lease back to its client, a transaction which is not recorded by DrillingInfo.com. To capture the information in the lease assignment process, I also scrape the website of the North Dakota Registry Information Network ([www.ndrin.com](http://www.ndrin.com)), which electronically records all lease assignments. I combine this lease and lease assignment data into a single dataset identifying the names of any firm that has mineral rights in a spacing unit. I assume that all firms with mineral rights in a well's spacing unit that are not the well's operator are non-operating participants.

### 2.4.4 Outside Experience

Throughout the paper, I assume that the only knowledge firms have about fracking comes from the wells fracked in North Dakota during 2005-2011. To assess the validity of this assumption, I collect firm-specific drilling history from IHS International for the 8 largest firms in my data, which I report in Table 5. In the first column, I list the number of wells each firm completed in the Bakken during 2005-2011. These 8 firms frack 60% of the wells in the dataset. During the time period I study, these firms are all publicly held, either as independent firms (Brigham, Continental Resources, EOG, Hess, Marathon and Whiting) or as subsidiaries of larger oil companies (Burlington is owned by Conoco Phillips, XTO is owned by Exxon Mobil).

On the right hand side of Table 5, I list the US operating history of these firms outside of North Dakota. In the ten years prior to the period I study, these firms collectively completed tens of thousands of vertical wells, which are typically drilled into conventional formations. However, they only completed 179 horizontal wells, suggesting that they had very little experience with the technology necessary to develop wells in the Bakken Shale. Only three firms had previously

completed more than ten horizontal wells, and two had done none. During 2005-2011, all eight firms are active outside North Dakota, with four firms completing more than a thousand wells each. Except for EOG and XTO, the vast majority of contemporaneous operational experience outside North Dakota is in vertical wells, though seven of the eight firms do complete horizontal wells. Thus, there appears to be limited scope for learning about fracking from experience outside of the Bakken.

### 3 The fracking production function

To quantify what knowledge firms learn about fracking, it is necessary to measure the empirical relationship between oil production, location and fracking input choices. I do this by estimating a production function for fracking. This production function accounts for variation in oil production across a well’s life and variation between wells in average production levels.

A well’s production changes over time due to age and maintenance-driven downtime. I measure these factors using a simple model common in the petroleum engineering literature. Because a well’s age is outside the firm’s control and because maintenance needs are both similar across wells and hard to predict, I argue that the time-varying error in production is plausibly exogenous.

Wells have different average production levels due to differences in their horizontal lengths, locations and fracking inputs. Location and fracking inputs may nonlinearly affect production, so I measure their impact non-parametrically, using Gaussian process regression (GPR), which I describe in detail below. The well-specific error in average production includes the effects of unobserved inputs, such as chemicals, variation across locations in the amount of oil that can be recovered and its sensitivity to fracking. I argue that chemical choices are independent of sand and water choices, and that the unobserved information firms observe about the well’s specific geological properties while drilling is unlikely to be correlated with the production error.

In the next two sections, I explain the production function model in further detail.

#### 3.1 The time series of oil production

Per unit of time, wells of all kinds (including non-fracked wells in conventional formations) tend to produce more oil when they are younger and less oil when they are older. This decline in performance over time is not surprising, because the amount of oil that can be recovered is finite and as more of it is pumped out of the ground, the rest becomes more difficult to recover. For nearly 70 years, petroleum engineers have used the simple "Arps" model to illustrate this basic phenomenon (see Fetkovich (1980)). The Arps model states that oil production in the  $t$ -th month of a well  $i$ ’s life is:

$$Y_{it} = Q_i t^\beta \exp(\nu_{it})$$

where  $Q_i$  is the *baseline* level of production,  $\beta < 0$  is a constant governing the production decline of the well and  $\nu_{it}$  is a mean-zero production shock. In log terms, this is

$$\log Y_{it} = \log Q_i + \beta \log t + \nu_{it}$$

meaning that a 1% increase in a well’s age should decrease per period production by  $-\beta\%$ , on average.

The operator of a well chooses  $D_{it}$ , the number of days during month  $t$  that well  $i$  is producing. Unless the well needs maintenance, there is no reason the operator would choose to produce for fewer than the full number of days during a month. All wells experience two routine maintenance events: the installation of external pumping hardware, and the connection of the well to a gas pipeline network. During maintenance, the operator must shut the well down, reducing  $D_{it}$ . My

data does not indicate whether maintenance occurs in a month, but it does report the number of producing days  $D_{it}$ , which I incorporate in the model:

$$\log Y_{it} = \log Q_i + \beta \log t + \delta \log D_{it} + \nu_{it}$$

The time-varying shock to log production,  $\nu_{it}$ , is the result of approximation error in the model. Firms cannot control  $t$ , the age of a well, and it is unlikely that firms observe anything correlated with  $\nu$  before choosing to do maintenance. Even if they did, firms would rather have the well producing on more days than fewer days, independent of  $\nu$ . Moreover, firms cannot predict  $\nu$  when fracking the well, which happens before production starts. For these reasons, I assume that:

$$\mathbb{E}[\nu_{it} \mid t, H_i, D_{it}, S_i, W_i, lat_i, lon_i] = 0$$

### 3.2 The cross section of oil production

I specify a semi-parametric model for log  $Q$ , the log of baseline production:

$$\log Q_i = \alpha + \eta \log H_i + f(S_i, W_i, lat_i, lon_i) + \epsilon_i$$

The parametric part of this model,  $\alpha + \eta \log H_i$ , is a Cobb-Douglas production function relating the horizontal length of a well to its baseline production. Though it may seem intuitive that  $\eta$  should equal one, there are practical reasons why this may not be true. Conversations with petroleum engineers indicate that fracking applied to the furthest away points of the horizontal segment of a well do not always perform as well as fracking applied to the closest points. If this decline in effectiveness is nonlinear, wells with longer horizontal segments may not proportionally outperform wells with shorter horizontal segments. The Hicks-neutral productivity  $\alpha$  measures the average log baseline production across wells. I discuss the well-specific productivity shock  $\epsilon_i$  in more detail below.

The function  $f(S_i, W_i, lat_i, lon_i) = f(Z_i)$  captures the relationship between baseline production, location and fracking choices. Table 4 in the data section suggests that this relationship differs across locations, and current petroleum engineering suggests that it may be nonlinear. For this reason, I estimate  $f(Z_i)$  non-parametrically, using Gaussian process regression, or GPR. GPR combines kernel regression techniques with a probability distribution over bandwidth parameters. Because there are few examples of GPR in applied economic settings, I provide a basic overview of its application here.

#### 3.2.1 Gaussian process regression

A *Gaussian process*  $G$  is a probability distribution over continuous real functions. Gaussian processes are defined by two functions: a mean function  $m(Z)$  and a positive definite covariance function  $k(Z, Z')$ . The mean function is the expectation of the value of a function  $f$  drawn at random from  $G$  at the point  $Z$ . The covariance function is the covariance between  $f(Z)$  and  $f(Z')$ . In mathematical terms, the mean and covariance functions satisfy:

$$m(Z) = \int f(Z) dG(f)$$

$$k(Z, Z') = \int (f(Z) - m(Z))(f(Z') - m(Z')) dG(f)$$



A Gaussian process is “Gaussian” because the joint distribution of the values  $f(Z_1)\dots f(Z_N)$  is multivariate normal, with a mean vector  $\mu$  and covariance matrix  $\Sigma$  given by:

$$\begin{aligned}\mu &= (m(Z_1)\dots m(Z_N))^\top \\ \Sigma_{i,j} &= k(Z_i, Z_j)\end{aligned}$$

This implies that the distribution of  $f(Z)$  is also normal with mean  $m(Z)$  and variance  $k(Z, Z)$ . The normality property makes it easy to compute the likelihood that a dataset  $(g_i, Z_i)_{i=1}^N$  is generated by the relationship  $g = f(Z)$  for a function drawn from a Gaussian process with mean  $m(Z)$  and covariance  $k(Z, Z')$ . By selecting mean and covariance functions from parametric families, the parameters that best fit the dataset can be estimated using maximum likelihood.

To estimate the function  $f(Z_i)$  above, I assume  $m(Z) = 0$  due to the presence of the constant term,  $\alpha$ , in the parametric portion of the production function. I assume that  $k(Z, Z')$  takes the form of a multivariate normal kernel:

$$k(Z_i, Z_j | \gamma) = \exp(2\gamma_0) \exp\left(-\frac{1}{2} \sum_{d \in S, W, lat, lon} \frac{(Z_{i,d} - Z_{j,d})^2}{\exp(2\gamma_d)}\right)$$

The first parameter,  $\gamma_0$ , measures the variance of the unknown function  $f(Z)$ . As points  $(Z_i, Z_j)$  become arbitrarily close to each other, the covariance function approaches the variance of  $f$ , and its formula collapses to  $\exp(2\gamma_0)$ . The remaining parameters ( $\gamma_S, \gamma_W, \gamma_{lat}, \gamma_{lon}$ ) are log-bandwidths, which measure how smooth  $f$  is in each dimension.

If the mean function is 0 and the covariance function parameters are  $\gamma = (\gamma_0, \gamma_S, \gamma_W, \gamma_{lat}, \gamma_{lon})$ , then the log likelihood of the data  $(g_i, Z_i)_{i=1}^N$  is:

$$\log \mathcal{L}(\gamma) = -\frac{1}{2} g^\top K(\gamma)^{-1} g - \log |K(\gamma)| - \frac{N}{2} \log(2\pi)$$

where  $g = (g_1 \dots g_N)^\top$  and  $K(\gamma)_{i,j} = k(Z_i, Z_j | \gamma)$ . The process of maximizing this likelihood over  $\gamma$  is called *Gaussian process regression*, or GPR. Conditional on  $\gamma$  and the data  $(g, \mathbf{Z})$ , the distribution of  $f$  evaluated at an out-of-sample point  $\tilde{Z}$  is normal, with mean and variance given by:

$$\begin{aligned}\mathbb{E} [f(\tilde{Z}) | g, \mathbf{Z}, \gamma] &= k(\tilde{Z} | \gamma)^\top K(\gamma)^{-1} g \\ \mathbb{V} [f(\tilde{Z}) | g, \mathbf{Z}, \gamma] &= k(\tilde{Z} | \gamma)^\top K(\gamma)^{-1} k(\tilde{Z} | \gamma)\end{aligned}$$

where  $k(\tilde{Z} | \gamma) = (k(Z_1, \tilde{Z} | \gamma) \dots k(Z_N, \tilde{Z} | \gamma))^\top$ . Note that the formula for the mean of  $f(\tilde{Z})$  is identical to the formula for the estimated regression function in kernel regression. The only difference between GPR and standard kernel regression techniques is the additional assumption that the unknown function  $f$  is drawn from a Gaussian process, making it possible to select bandwidths  $\gamma$  using likelihood techniques.

Gaussian processes are commonly used in the artificial intelligence and operations research literatures, though their application in economics is so far limited to econometric theory.<sup>24</sup> For a detailed treatment of Gaussian processes, see Rasmussen and Williams (2005).

---

<sup>24</sup>See Kasy (2013) for a recent example.

### 3.2.2 The well-specific shock $\epsilon_i$

The well-specific shock to log baseline production,  $\epsilon_i$ , contains approximation error, unobserved inputs to the fracking process and unobservable variation in geology. Fracking chemicals are the primary unobserved input. Firms use chemicals to inhibit bacterial growth in the fracking mixture, to provide lubrication for the pumping units used in fracking and to prevent corrosion and mineral scaling.<sup>25</sup> Because chemicals do not affect the efficiency of observed fracking inputs, I assume that sand and water choices are independent of them.

The petroleum engineering literature predicts that different parts of the Bakken contain different amounts of oil and respond to fracking inputs differently.<sup>26</sup> Aside from the latitude and longitude coordinates of a well, I do not have data to control for these differences. If firms observe signals about how much oil a well contains or how amenable it is to fracking, they may adjust their fracking inputs in response and  $\epsilon_i$  will not be independent of these choices. Unfortunately, I do not have instruments for fracking input choices, so it is important to consider what additional information firms have.

The most relevant information firms observe before fracking a well is contained in the samples of rock that firms collect during the drilling process. As the drill bit passes through the Upper Bakken Shale on its way into the Middle Bakken, firms have an opportunity to analyze the returned cuttings. These may be indicative of the amount of the oil in the Upper Bakken Shale at the location where the horizontal segment starts. However, since the goal in horizontal drilling is to stay inside the Middle Bakken, firms receive no additional information about the Upper Bakken Shale and receive no information at all about the Lower Bakken Shale. Moreover, the characteristics of the Upper Bakken Shale can change over the length of the horizontal segment, and there is no guarantee that the Lower Bakken Shale has the same characteristics at a point as the Upper Bakken Shale. Thus, the information firms can acquire during drilling is unlikely to be helpful in choosing fracking inputs.

For this reason, I argue that  $\epsilon_i$  is exogenous to firm choices and other well characteristics:

$$\mathbb{E}[\epsilon_i \mid t, H_i, D_{it}, S_i, W_i, lat_i, lon_i] = 0$$

Combining everything together, the whole production function model is:

$$\log Y_{it} = \alpha + \beta \log t + \delta \log D_{it} + \eta \log H_i + f(Z_i) + \epsilon_i + \nu_{it}$$

Since Gaussian process regression generates a normal likelihood for  $f(Z_i)$ , I assume that  $\nu_{it}$  and  $\epsilon_i$  are both normal, with zero mean and variances  $\sigma_\nu^2$  and  $\sigma_\epsilon^2$ , respectively.

### 3.3 Likelihood

I compute the likelihood function in two steps. In the first step, I treat the unobserved effect of fracking and location  $f(Z_i)$  as observed and compute the likelihood of  $(Y_{it}, X_{it})$  conditional on  $f(Z_i)$  and the parameters. In the second step, I integrate out the unobserved values of  $f(Z_i)$  using the likelihood function for  $f(Z_i)$  generated by GPR. I describe the likelihood calculation in detail in the appendix.

### 3.4 Production function estimates

Table 6 shows maximum likelihood estimates of the semi-parametric production function described above in addition to a simpler parametric specification. The parametric specification replaces

<sup>25</sup>See <http://www.fracfocus.org> for further details on the chemicals used in fracking.

<sup>26</sup>See Baihly et al. (2012)

$f(S_i, W_i, lat_i, lon_i)$  with township fixed effects,  $\tau(i)$ , and a Cobb-Douglas production technology in sand and water,  $\kappa_S \log S_i + \kappa_W \log W_i$ .

All of the parametric model coefficients are statistically significantly different from zero in both specifications and the coefficients common to both have similar estimates. As expected, wells produce less oil per month as they age, with an estimated log decline rate of  $-0.56$ .<sup>27</sup> The coefficient on days producing is 1.75, suggesting that when wells undergo maintenance, production per day is lower than when wells do not have maintenance issues. Wells with longer horizontal segments produce more than wells with shorter segments, but the effect is not linear. Doubling the horizontal length of a well increases production by 80% in the Cobb-Douglas specification and 85% in the Gaussian process. The variance of  $\epsilon$  is larger in the Cobb-Douglas specification than in the Gaussian process, suggesting that the flexibility of the Gaussian process explains more of the variation in oil production than Cobb-Douglas and location fixed effects do. The estimated Cobb-Douglas marginal productivities of sand and water are precisely estimated and are smaller than the preliminary estimates in Table 4. Sand and water both increase oil production, with decreasing returns to scale.

The estimated GPR bandwidth coefficients do not have an intuitive interpretation, so I illustrate the estimated production relationships graphically in Figure 5. The top panel is a contour plot of the non-parametrically estimated function  $f(S_i, W_i, lat_i, lon_i)$ , evaluated at the geographic centroid of the most active township. The lines are iso-production curves, which are combinations of sand and water choices with the same estimated value of  $f$ . Across all levels of water use, greater sand use is associated with higher oil production, while greater water use is only associated with higher production at the highest level of sand use, and only in a limited range. The middle panel shows contour lines for the Cobb-Douglas specification. The Gaussian process and Cobb-Douglas specifications make starkly different predictions about the impact of fracking inputs and location on oil production. At the average sand and water choices for this township, 266 lbs and 131 gals per foot, respectively, the Gaussian process predicts -3.5 log points of production, while Cobb-Douglas predicts -3.1. Additionally, the non-parametric specification makes different predictions in different locations. The bottom panel shows contour lines for the production function evaluated at the centroid of a nearby township. The location of the most productive sand and water choices differ across the two townships. In the top panel, the maximal choice is approximately 600 lbs sand and 200 gals water, per foot, while in the bottom panel it is 400 lbs sand and 500 gals water, per foot. This variation across townships in the relationship between oil production and inputs is not possible with the Cobb-Douglas specification, so for the rest of the paper, I focus on the Gaussian process specification.

The fit of both models is high, with  $R^2$ 's of 78% for the Cobb-Douglas model and 81% for the Gaussian process model. The “between”  $R^2$ 's, which measure the correlation of predicted baseline production and actual baseline production, are higher, at 81% and 88%, respectively. The production function models fit the data well for several reasons. Both the inputs to fracking, sand and water, and the single output of fracking, crude oil production, are precisely measured. The one unobserved input, fracking chemicals, does not directly affect production or observed input choices. Moreover, the production function for fracking is an approximation to a “true” physical relationship between sand, water, location and oil production. However, since I estimate this approximation non-parametrically, there is the possibility that the estimated bandwidths are too narrow, leading to over-fitting.

To check for this, I perform a cross-validation test of the model estimates. For each of 25 test runs, I randomly split the wells into two separate datasets: a training dataset containing 90% of the

---

<sup>27</sup>Note, current geophysics research on the Bakken has found similar decline rates. Hough and McClurg (2011), for example, estimates the decline rate to be  $-0.5$ .

Table 6: Production function model estimates

Coefficient	Cobb-Douglas		Gaussian Process	
	Estimate	Std. Error	Estimate	Std. Error
$\alpha$			-4.4152	(0.3278)
$\beta$	-0.5576	(0.0024)	-0.5570	(0.0024)
$\delta$	1.7543	(0.0035)	1.7549	(0.0035)
$\eta$	0.7977	(0.0363)	0.8479	(0.0357)
$\gamma_0$			-0.3945	(0.0572)
$\gamma_S$			6.1757	(0.1343)
$\gamma_W$			5.9467	(0.1232)
$\gamma_{lat}$			-2.4702	(0.0539)
$\gamma_{lon}$			-2.2376	(0.0609)
$\kappa_S$	0.1582	(0.0157)		
$\kappa_W$	0.1148	(0.0159)		
$\log \sigma_\epsilon$	-0.9086	(0.0147)	-1.0591	(0.0187)
$\log \sigma_\nu$	-0.4898	(0.0024)	-0.4897	(0.0024)
Township Fixed-effects	X			
Overall $R^2$	0.783		.811	
Between $R^2$	0.813		.882	
Within $R^2$	0.764		.764	
# Wells	2,699			
# Well-months	91,783			

Maximum likelihood estimates of the Cobb-Douglas production function model:

$$\log Y_{it} = \alpha + \beta \log t + \delta \log D_{it} + \eta \log H_i + \kappa_S \log S_i + \kappa_W \log W_i + \tau(i) + \epsilon_i + \nu_{it}$$

and the Gaussian process production function model:

$$\log Y_{it} = \alpha + \beta \log t + \delta \log D_{it} + \eta \log H_i + f(Z_i | \gamma) + \epsilon_i + \nu_{it}$$

$Y_{it}$  is oil production for well  $i$  when it is  $t$  months old,  $D_{it}$  is the number of days producing,  $H_i$  is the horizontal length, and  $Z_i$  is the vector of sand use  $S_i$ , water use  $W_i$ , latitude  $lat_i$  and longitude  $lon_i$ .  $\tau(i)$  is a set of township fixed effects. “Between”  $R^2$  is the  $R^2$  for the average predicted log baseline production. “Within”  $R^2$  is the  $R^2$  for the predicted time series of production.

Figure 5: Contour Plots of Production Function Estimates

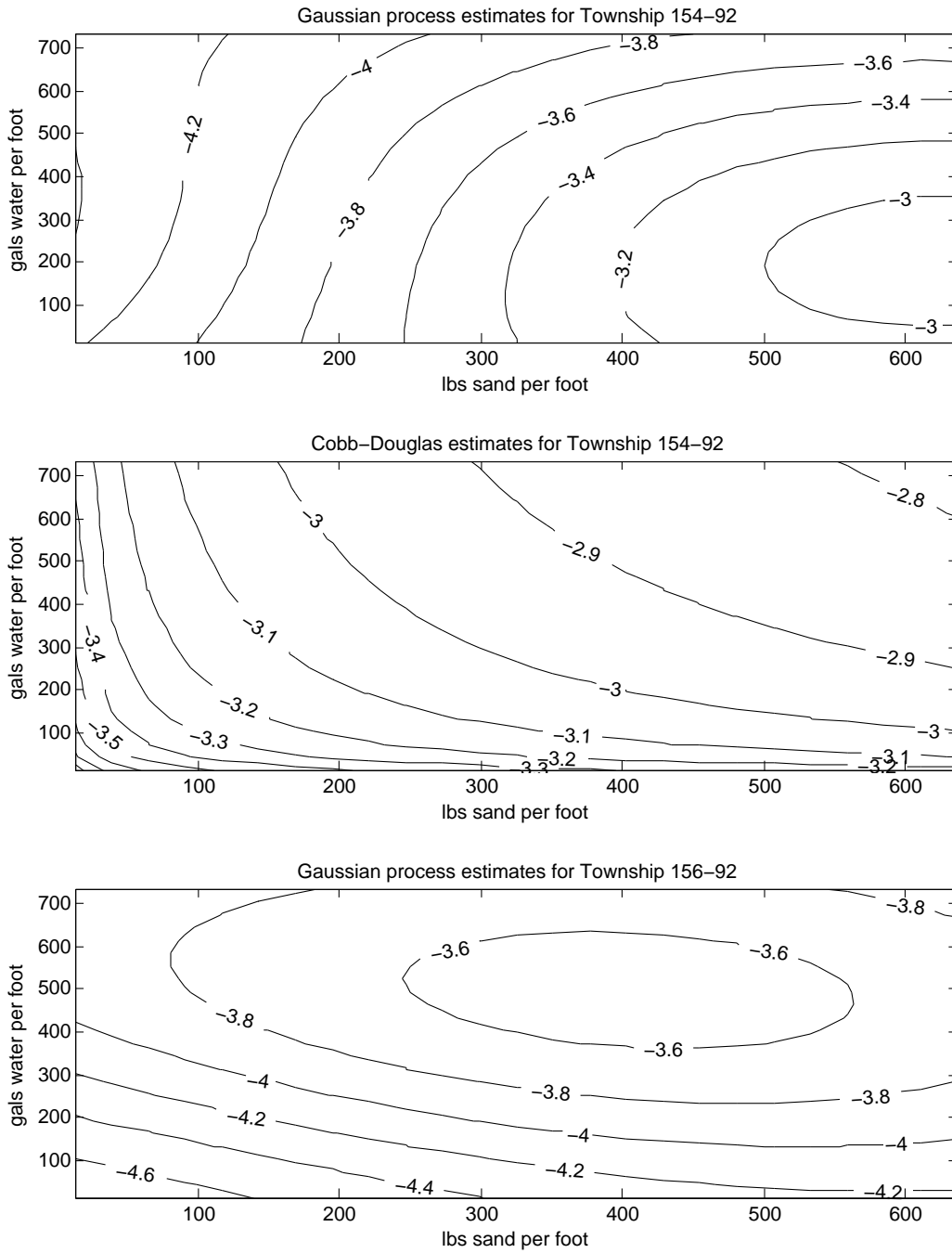


Table 7: Production function model cross-validation statistics

Coefficient	Average Estimate	Std. Dev. of Estimate
$\alpha$	-4.3388	0.1295
$\beta$	-0.5570	0.0016
$\delta$	1.7533	0.0055
$\eta$	0.8404	0.0136
$\gamma_0$	-0.4046	0.0290
$\gamma_S$	6.1659	0.0532
$\gamma_W$	5.9278	0.0616
$\gamma_{lat}$	-2.4454	0.0236
$\gamma_{lon}$	-2.2211	0.0392
$\log \sigma_\epsilon$	-1.0545	0.0109
$\log \sigma_\nu$	-0.4916	0.0063
$R^2$ comparisons		
$R^2$ type	Avg. in training	Avg. in validation
Overall $R^2$	0.8116	0.7835
Between $R^2$	0.8826	0.8098
Within $R^2$	0.7652	0.7594
# Wells	2,699	
# Well-months	91,783	
# Cross validation samples	25	

Maximum likelihood estimates of the production function model:

$$\log Y_{it} = \alpha + \beta \log t + \delta \log D_{it} + \eta \log H_i + f(Z_i | \gamma) + \epsilon_i + \nu_{it}$$

$Y_{it}$  is oil production for well  $i$  when it is  $t$  months old,  $D_{it}$  is the number of days producing,  $H_i$  is the horizontal length, and  $Z_i$  is the vector of sand use  $S_i$ , water use  $W_i$ , latitude  $lat_i$  and longitude  $lon_i$ .

wells, and a validation dataset containing the remaining 10%. I re-estimate the production function on the training dataset and use the estimates to predict production in the validation dataset. I save the estimated production function coefficients, the  $R^2$  values applied to the training data and the  $R^2$  values applied to the validation data, and report their distribution across test runs in Table 7. The parametric components of the production function model are quite stable across runs, with the average model estimates being similar to the full dataset maximum likelihood estimates. The standard deviations across runs are smaller than the maximum likelihood standard errors for the full dataset. Though the  $R^2$  values for validation samples are generally lower than for training samples, they are still quite high, with the average overall  $R^2$  for validation samples at approximately 78%, compared to 81% in the training samples. The stability of the coefficient estimates and the consistently high goodness-of-fit measures in validation samples suggest that the maximum likelihood estimates in Table 6 do not suffer from over-fitting and represent a stable and causal relationship between inputs and production.

## 4 Evidence for Learning

As firms learn to use fracking technology more efficiently, they should make more profitable fracking design choices. If oil prices, input costs and the quality and size of drilling locations were constant over time, I could test this prediction by extrapolating future production from current production and simply checking if average expected discounted profits per well increased over time. However, oil prices, input costs and locations do vary over time, so I control for this variation by examining trends in the ratio of actual profits to counterfactual maximal profits. That is, I compute a profitability measure which compares the profits firms earned with the highest amount of profits they could have earned with the best fracking design.

I use the fracking production function to compute these profits. The profits to well  $i$  fracked using design  $j$  are

$$\Pi_{ij} = \lambda P_i \mathbb{E} \left[ \sum_{t=1}^T \rho^t \tilde{Y}_{ijt} \right] - c_i(S_j, W_j)$$

where  $\lambda$  is the fraction of oil production the firm keeps for itself,  $P_i$  is the price the firm will receive for its oil production,  $T$  is the number of periods the well is expected to produce for,  $\rho$  is the per-period discount rate,  $\tilde{Y}_{ijt}$  is the realization of the level of oil production for well  $i$  under fracking design  $j$  at age  $t$ , and  $c_i(S_j, W_j)$  is the total cost of drilling and fracking that design.<sup>28</sup> The main empirical object needed in the calculation of  $\Pi_{ij}$  is the expected present value of discounted oil production,  $\mathbb{E}[DOP_{ij}]$ :

$$\begin{aligned} \mathbb{E}[DOP_{ij}] &= \mathbb{E} \left[ \sum_{t=1}^T \rho^t \tilde{Y}_{ijt} \right] \\ &= \sum_{t=1}^T \rho^t \mathbb{E} \left[ \tilde{Y}_{ijt} \right] \end{aligned}$$

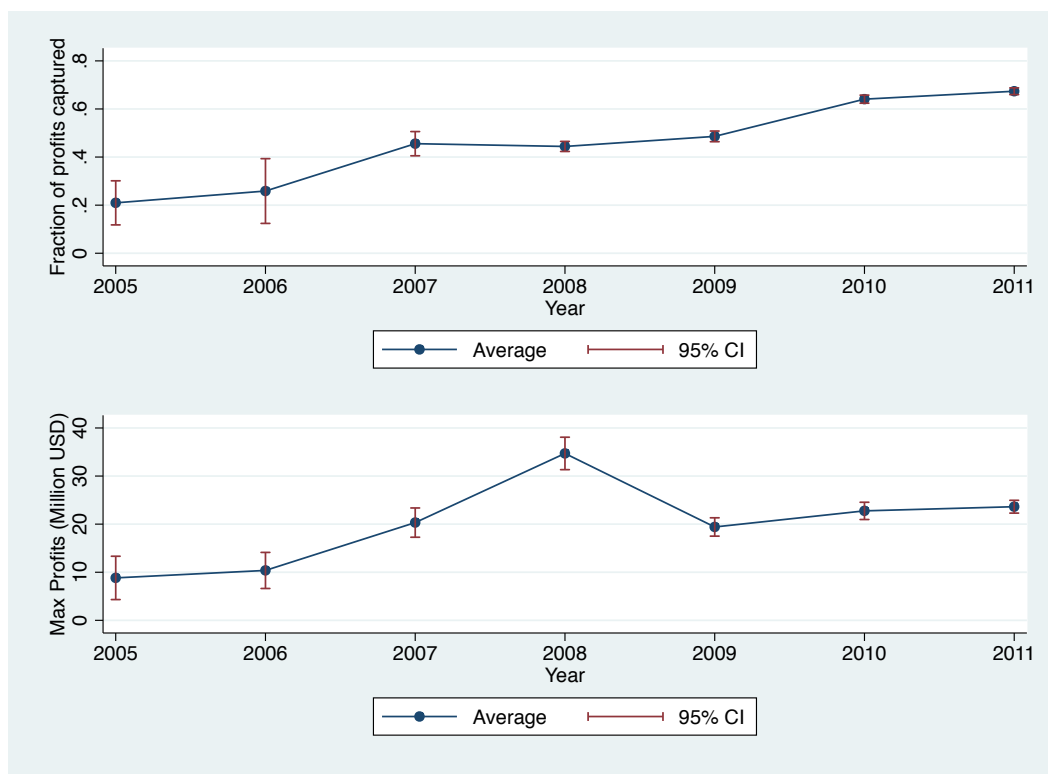
I compute this expectation conditional on two different information sets: the full data that I have, and the data each firm had when it made a fracking design decision. The first case represents an *ex post* expectation, and provides a way of asking whether firms made better fracking design decisions over time, given today's knowledge. The second case represents an *ex ante* expectation, and provides a way of asking whether firms' choices were consistent with static profit maximization, given my measures of their information sets.

In both cases, I combine the production function parameter estimates in Table 6 with the normality assumptions on the unobserved terms to compute a probability distribution over oil production. Since the production function estimates depend on the full dataset, this means that I am computing *ex ante* expectations under the assumption that firms had the same beliefs about the production function parameters as I do now. This is a strong assumption. The *ex ante* calculation of expected oil production will be biased if firms had different beliefs than I do about the decline rate  $\beta$ , the productivity of producing days  $\delta$  and horizontal length  $\eta$ , the bandwidth parameters  $\gamma$  and the variances  $\sigma$  of the unobservable production shocks. I assume that these biases are small, as decline rates and productivity parameters can be predicted using geophysical models<sup>29</sup>,

<sup>28</sup>I assume that the fraction of oil revenue that accrues to the firms is 70%, based on typical royalty rates of 16.5%, state taxes of 11.5% and ongoing operating costs of 2%. I assume  $T = 240$  months, though the NDIC expects Bakken wells to produce for 540 months, making these profit calculations an underestimate. I set  $\rho = .9$ , which is the standard discount rate use in oil & gas accounting. At this rate, the difference between 540 months and 240 months is only 2.6% in present value terms.

<sup>29</sup>See Fetkovich (1980).

Figure 6: Fraction of Positive Profits Captured and Maximal Profits by Year, ex post



and bandwidth and variance parameters do not affect the asymptotic properties the production function estimate.<sup>30</sup> Moreover, the impact of fracking design and location  $f(Z)$  is computed nonparametrically from both the bandwidth parameters  $\gamma$  and the information set. Thus firms with different information sets will have different beliefs about  $f(Z)$ , and these beliefs will differ from the *ex post* beliefs as well.

I present the full calculation of expected discounted oil production in the appendix.

#### 4.1 *ex post* Comparisons

Over time, firms choose fracking designs with higher *ex post* expected profits. The top half of Figure 6 plots the *ex post* ratio of actual profits to maximal profits per well.<sup>31</sup> The average fraction of profits captured increases nearly monotonically over time, from 21% in 2005 to 67% in 2011. Much of this growth happens in two phases. Between 2006 and 2007, the fraction increases from 25.9% to 45.6%, and between 2009 and 2010, the fraction increases from 48.6% to 64.1%. By 2011, firms earn an average of 67.3% of the maximum profits they could have earned with optimal fracking input choices.

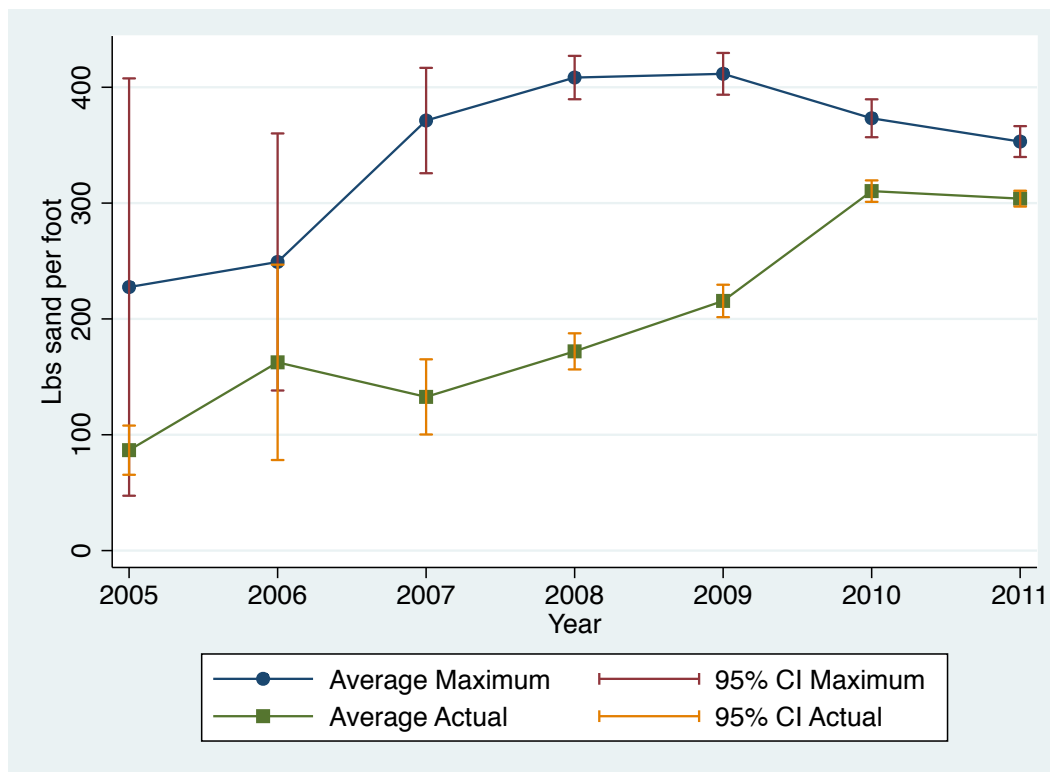
The bottom half of Figure 6 shows how these maximal profits evolve over time. When oil prices were at their peak in 2008, the profit maximizing input choice for the average would have generated \$34.7 million in profits, meaning that in 2008, foregone profits from inefficient fracking choices averaged \$19.3 million per well. By 2011, lower oil prices reduced these maximal profits

<sup>30</sup>See section 7.1 in Rasmussen and Williams (2005).

<sup>31</sup>I only include wells in this calculation that have both positive actual profits and positive maximal profits. Over the entire sample, 5.2% of wells have either negative actual profits or negative maximal profits. This selection is not uniform over time, as 29% of excluded wells occur in 2009 and 35% occur in 2011.



Figure 7: Average Profit Maximizing Sand Use and Actual Sand Use Per Well, ex post



to \$23.6 million per well. Combined with the higher fraction of profits captured, firms in 2011 left only \$7.7 million on the table.

Firms captured more profits by selecting more profitable fracking designs over time. In Figures 7 and 8, I plot average profit maximizing and actual input use per well over time. Though firms use less sand in fracking than the estimated profit maximizing levels, actual choices approach optimal choices, starting in 2009. In 2005 and 2006, the average well was fracked with approximately 100 lbs sand per foot less than the profit maximizing level. This difference in sand use actually increases in 2007 and 2008 to 215 lbs per foot, before starting to fall in 2009. By 2011, the difference between optimal sand use and actual sand use is only 39 lbs per foot.

Though the differences in actual and optimal water use start out considerably larger than the differences in sand use, actual water choices get closer to optimal water choices in every year. In 2005, firms fracked the average well with 468 gals per foot less water than the water use in the optimal well. By 2011, the difference is only 56 gals per foot. These trends in actual input use towards optimal input use are consistent with the idea that firms are learning about the efficient use of fracking inputs as they observe more data, and with this knowledge they make more profitable choices.

## 4.2 Profitability vs. Productivity

The existing literature on learning in firms focuses on *productivity* instead of *profitability*. Researchers in this literature measure learning by comparing estimates of Hicks-neutral productivity

Figure 8: Average Profit Maximizing Water Use and Actual Water Use Per Well, ex post

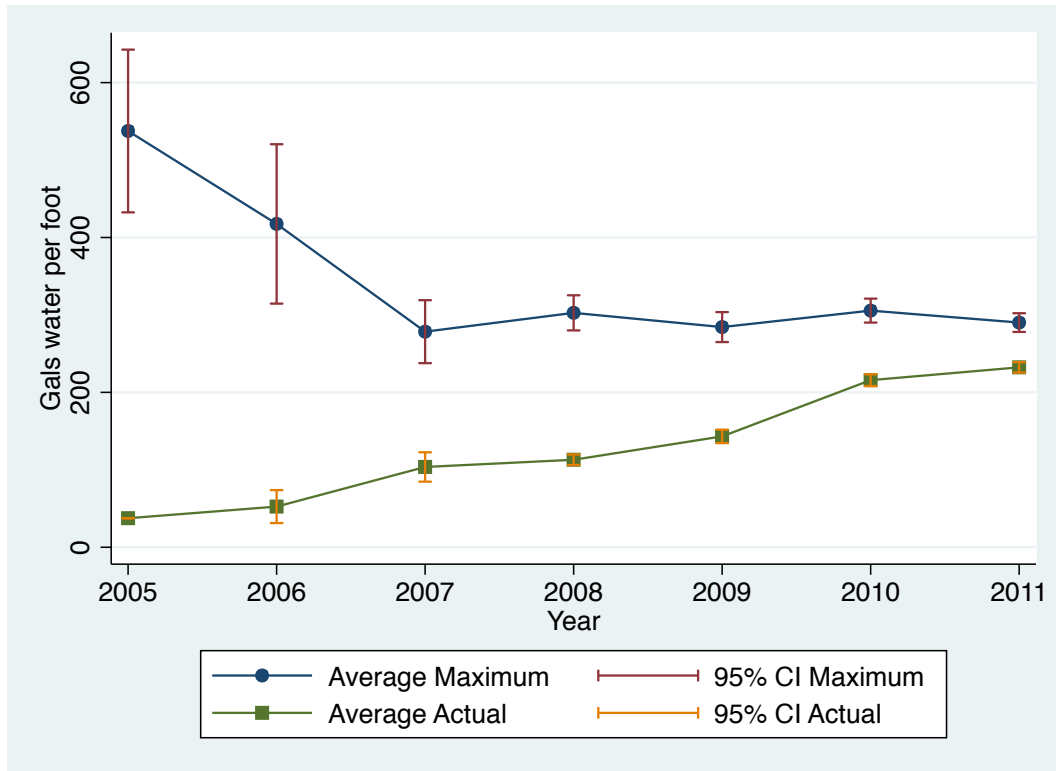
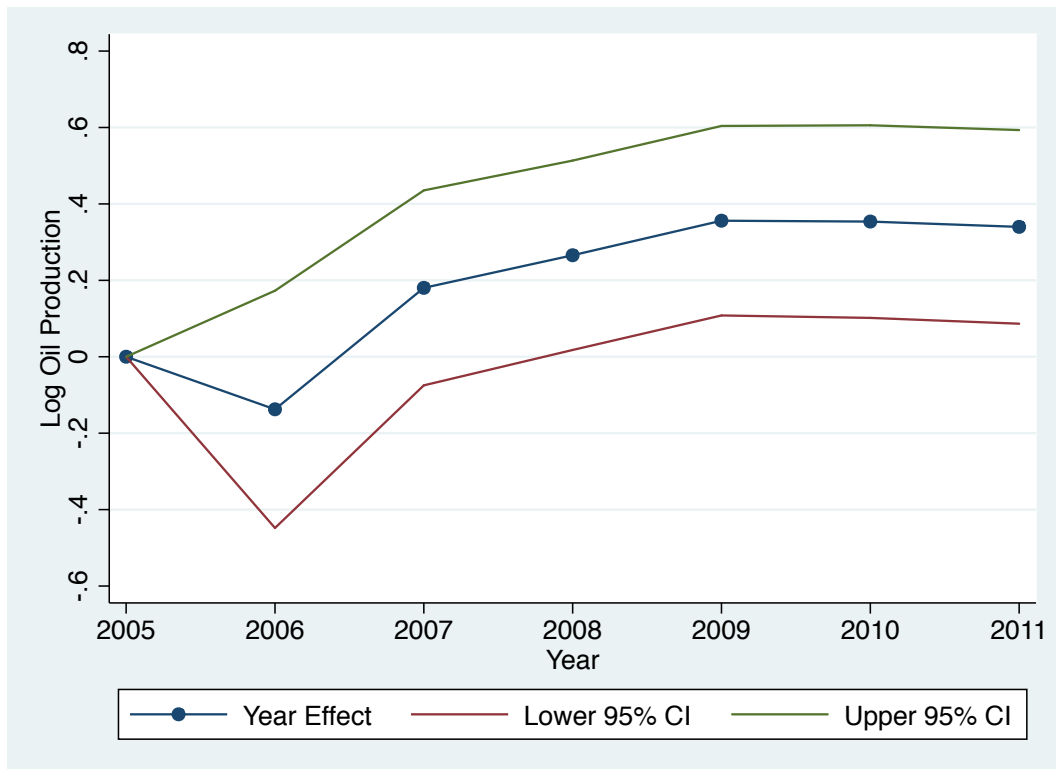


Figure 9: Gaussian Process Year Effects



with the amount of experience a firm has in producing.<sup>32</sup> This approach to studying learning does not treat the production function as an object for firms to learn. Rather, the knowledge from accumulated experience serves as an *input* to the firm’s production function, in the same way that labor, capital and materials do. To determine if firms in this dataset became more productive, in addition to more profitable, I add year fixed effects to the Gaussian process production function specification, and plot their estimated values and confidence intervals in Figure 9.

Wells fracked in 2005 are actually 13.7% *more* productive than wells fracked in 2006. However, the confidence interval around this estimate is wide enough to include zero, as there are only 10 wells in 2005 and 20 wells in 2006. Wells fracked in later years are significantly more productive than wells fracked in 2005 or 2006. For example, wells fracked in 2009 are 35.6% more productive than those in 2005, and 49.3% more productive than those in 2006. Again, the confidence intervals around these estimates are wide, and I cannot reject the hypothesis that there is no change in productivity between 2006 and 2009. In each of the next 2 years, productivity falls slightly, though the differences are not statistically significant. Overall, wells fracked between 2008-2011 cohorts are more productive than the earliest wells, but there is no productivity growth during 2008-2011. Since this time period covers 95% of the wells studied in this paper, I interpret this as evidence that firms learned to be more productive only in the earliest years. In contrast, the results in the previous section show that firms learned to be more profitable in all years, and especially during 2008-2011.

### 4.3 *ex ante* Comparisons

Though firms make choices which approach the *ex post* estimates of optimal choices over time, those choices do not always maximize the *ex ante* estimates of expected profits. The top half of Figure 10 plots the ratio of actual profits to maximal profits per well using *ex ante* expectations.<sup>33</sup> Firms initially make fracking input choices with expected profits that are close to the optimal choices, capturing 79.5% of potential *ex ante* profits in 2007. However, profit capture actually falls over time, reaching 67.4% in 2011, approximately the same level as the *ex post* case in 2011.

While the fraction of profits captured falls, *ex ante* expectations of maximal profits rise over time, as show in the bottom half of Figure 6. Unlike the *ex post* case, where the highest level of maximal profits coincides with the 2008 peak in oil prices, *ex ante* maximal profits are highest in 2011, reaching \$26.3 million per well. Though average oil prices are similar in 2008 (\$100 per bbl) and 2011 (\$95 per bbl), firms have much more information about fracking in 2011 and this information generates more optimistic expectations. The combined effect of falling *ex ante* profit capture and rising maximal profits increases foregone *ex ante* profits from \$2.8 million in 2007 to \$9.6 million in 2011.

Firms capture a shrinking fraction of *ex ante* profits over time because their actual sand use diverges from the expected profit maximizing sand use. Figure 11 plots average profit maximizing and actual sand use per well over time. In 2007 and 2008, actual sand use is quite similar to *ex ante* optimal sand use, and in 2008, the average well is actually fracked with 10 lbs per foot more sand than the optimal level. However, as the data firms have to learn from accumulates, optimal sand use increases faster than actual sand use, and by 2011, the difference between optimal and

---

<sup>32</sup>For example, Benkard (2000) correlates log labor requirements per unit of production with measures of experience (and forgetting), and Thornton and Thompson (2001) estimate a semi-parametric production function model in which various measures of experience are direct inputs to production.

<sup>33</sup>As in the *ex post* case, I only include wells in this calculation that have both positive actual profits and positive maximal profits. Over the entire sample, 4.8% of wells have either negative actual profits or negative maximal profits. Half of these wells are fracked in 2009. Moreover, I further limit the set of wells by computing expected profits for the subset of wells that are fracked by firms which can observe 50 wells and 300 well-months of production history. The first wells that satisfy this criteria are not fracked until 2007.

Figure 10: Fraction of Positive Profits Captured and Maximal Profits by Year, ex ante

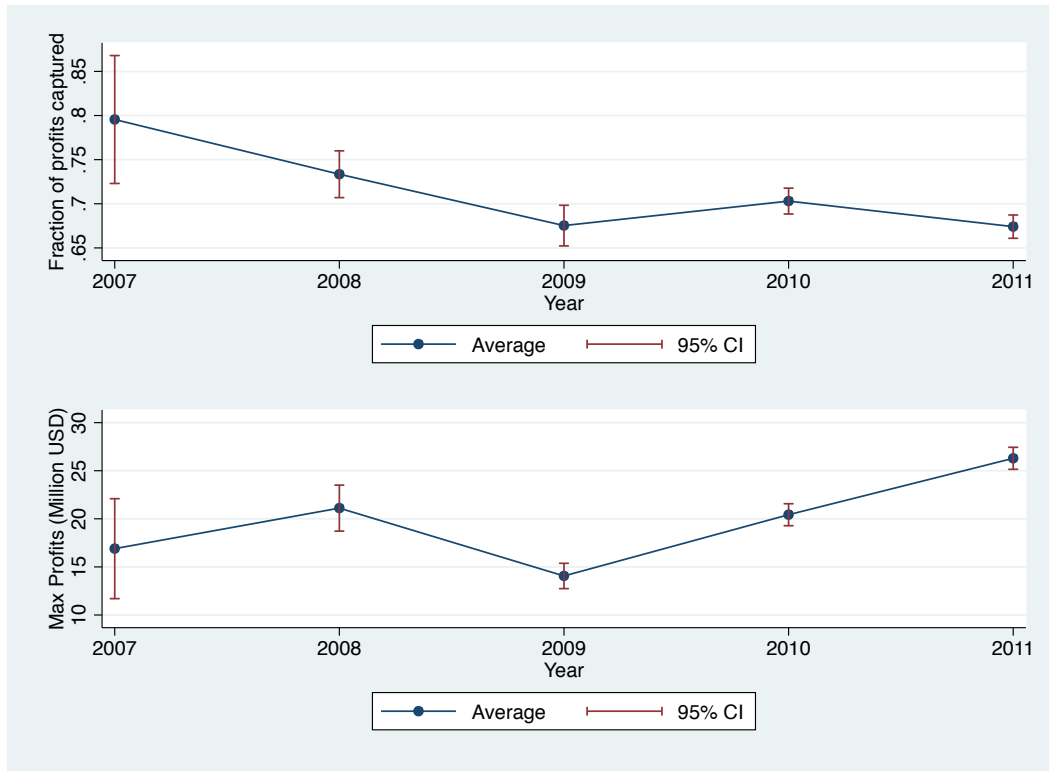


Figure 11: Average Profit Maximizing Sand Use and Actual Sand Use Per Well, ex ante

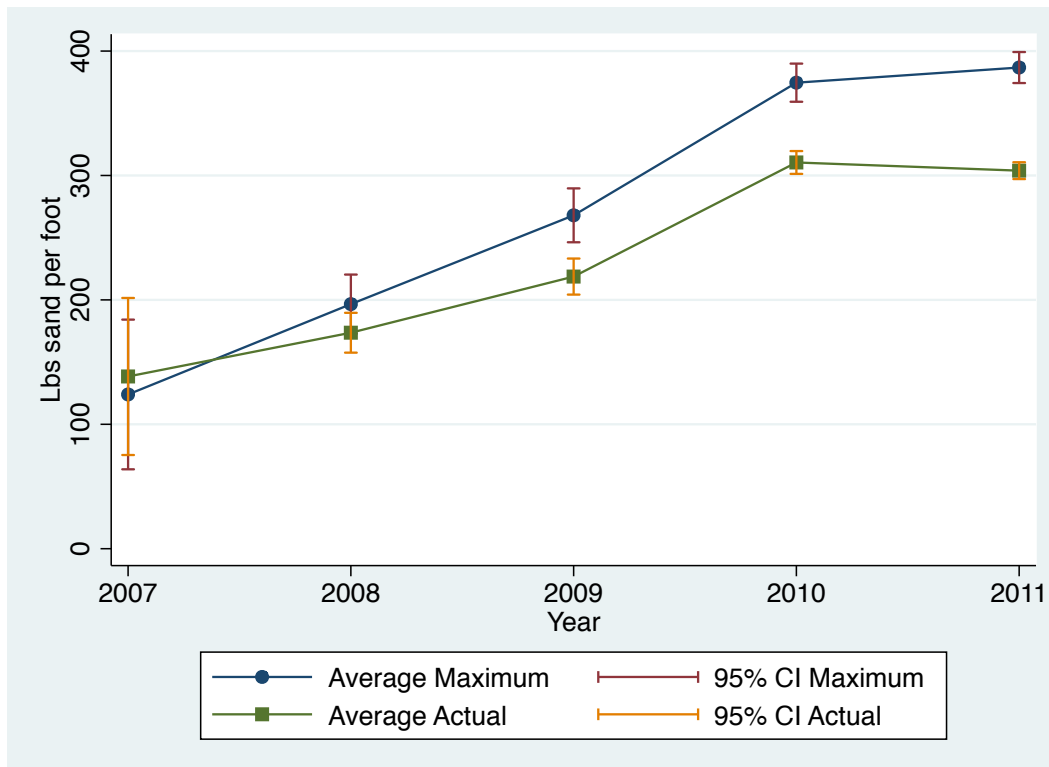
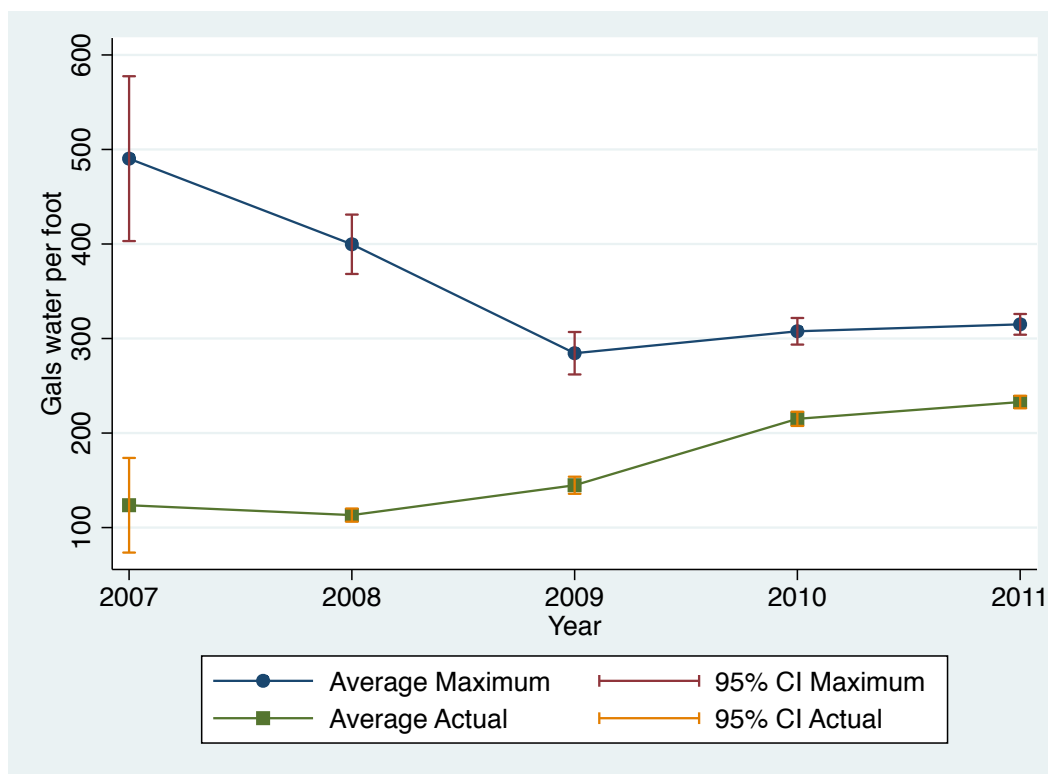


Figure 12: Average Profit Maximizing Water Use and Actual Water Use Per Well, *ex ante*



actual sand use reaches 77 lbs per foot. Though this difference is similar to the difference in the *ex post* case during 2011, it is striking that the differences in actual and optimal sand use increase over time in the *ex ante* case while decreasing in the *ex post* case.

Figure 12 plots average *ex ante* optimal and actual water use per well is similar to the *ex post* case in Figure 8. On average, firms use less than the *ex ante* optimal amount of water in fracking, but make improved water choices over time. In 2007, firms use 395 gals per foot less water than the optimal level. This difference shrinks in each year, and by 2011, it is only 139 gals per foot.

## 5 Fracking input choice model

Though firms do learn over time, many of their choices do not coincide with the predicted optimal choices, even on an *ex ante* basis. I consider two possible explanations for this phenomenon based firm preferences. First, firms may care about the uncertainty in their estimates of the profits of a fracking design. Second, in estimating the profits of a fracking design, firms may prefer to weigh their own data differently than the data generated by their competitors.

### 5.1 Preferences over uncertainty

In comparing the expected profits a firm earned to the maximal expected profits a firm could have earned, I have implicitly assumed that the *correct* learning strategy is for firms to select fracking designs solely on the basis of expected profits, without regard to the extent of profit variance across designs. There are two potential problems with this assumption. First, viewing fracking design as an investment project selection problem, there may be financial or organizational factors that

cause firms to have preference over variance. Second, when learning about the performance of different fracking designs, firms may care about variance through the *explore vs. exploit* tradeoff that exists in all learning problems.

In a simple, frictionless model of investment project selection, it is appropriate for firms to ignore variance. In practice, there may be reasons why this would not be optimal. Firms raise outside capital to finance operations and the presence of debt capital can lead firms to select fracking designs with higher variance, as bond holders bear some of the downside risk. On the other hand, capital constrained firms may not necessarily have the option of selecting fracking designs with higher variance if they are more expensive. Financial considerations can thus push firms towards or away from fracking designs with higher variance in profits. Firms must also hire and incentivize potentially risk averse engineers, who select fracking designs. Depending on the extent of their career concerns and the structure of their compensation, engineers may prefer fracking designs with more or less variance.

The prescribed learning strategies in most theoretical models of learning involve variance seeking behavior. Analyses of the *explore vs. exploit* tradeoff in learning predict that agents should always do some amount of exploration, by occasionally selecting actions with high variance. In most of the settings studied by Aghion et al. (1991), a fully rational, expected present discounted value maximizing agent will do some amount of exploring forever and a similar result obtains in the multi-agent context studied by Bolton and Harris (1999). The implied preferences for variance in both of these models arise out of the natural dynamics of learning problems. Agents are still risk neutral over their payoffs, but because there is present value to better information in the future, between two actions with the same expected flow payoff but different variances, they prefer the action with higher variance.

Empirically, oil companies exhibit both risk seeking and risk averse behavior. The process of acquiring mineral rights for new drilling prospects and establishing the existence of oil within those prospects is an especially risky one (see, for example Walls and Dyer (1996) and Reiss (1989)). However, oil companies are price takers in the world market for oil, and many use financial markets to hedge some or all of their future oil production, suggesting that firms may wish to avoid risks associated with future price fluctuations (see Haushalter (2000)).

Whether the companies I study here prefer fracking input choices with more or less variance is an empirical question. I estimate firm preferences over expectations and variance of fracking designs by analyzing realized choices. To do this, I fit a multinomial logit preference model of fracking design choice in which the “utility” a firm has for fracking design  $j$  applied to well  $i$  is:

$$\begin{aligned} u_{ij} &= \lambda P_i \left( \xi_m \mathbb{E}[DOP_{ij}] + \xi_s (\mathbb{V}[DOP_{ij}])^{\frac{1}{2}} \right) - c_i(S_j, W_j) + \epsilon_{ij} \\ &= \tilde{u}_{ij}(\xi_m, \xi_s) + \epsilon_{ij} \end{aligned}$$

where  $\lambda$  is the fraction of oil revenues firms keep,  $P_i$  is the price of oil for well  $i$ ,  $c_i(S_j, W_j)$  is the cost of fracking design  $j$  for well  $i$ , and  $\epsilon_{ij}$  is an iid logit error. The parameters  $(\xi_m, \xi_s)$  represent the firm’s preference over expected present discounted revenues and the standard deviation of present discounted revenues, conditional on the data they have. Under this preference specification, the probability that a firm selects design  $j$  for well  $i$  is

$$p_{ij} = \frac{\exp(\tilde{u}_{ij})}{\sum_k \exp(\tilde{u}_{ik})}$$

The mean utilities in this preference model are linear in the expectation and standard deviation of profits to a fracking design. Preferences of this type have precedence in the theoretical learning literature. Brezzi and Lai (2002) show that a linear combination of the expectation and standard

deviation of the payoff to a choice can represent a simple and efficient approximation to the Gittins index value for the choice, if the choices have independently distributed payoffs. Since Gittins and Jones (1979) show that ordinal preferences over Gittins indices result in dynamically efficient learning behavior, agents that utilize these linear approximations attain near-optimal learning. Though the profits to fracking input choices are not distributed independently, authors in the computer science and operations research literatures have found that these learning strategies also perform well in the general case. In those literatures, learning strategies which select the choice with the highest value of a linear combination of the expectation and standard deviation of payoffs are called “upper confidence bound”, or UCB strategies. Rusmevichientong and Tsitsiklis (2010) and Srinivas et al. (2012) have established that UCB strategies quickly identify the highest performing choice, and do so in a way which minimizes an agent’s *ex post* cumulative regret over its past choices. UCB strategies are also reported to be in use at major technology companies, like Yahoo, Microsoft and Google (see Chapelle and Li (2011), Graepel et al. (2010) and Scott (2010)). In all of the existing literature which utilizes UCB learning strategies, the weight on the standard deviation of the payoffs to a choice is positive, hence the “upper” in upper confidence bound strategies. This paper is not the first in economics to utilize UCB learning strategies in an empirical context. Dickstein (2013) estimates the parameters of a UCB learning strategy in a study of learning behavior by physicians.

With data on the choices firms made, expectation and standard deviation calculations made using their information sets, and oil price and fracking cost data, I estimate the parameters  $(\xi_m, \xi_s)$  using maximum likelihood. I estimate separate values of  $(\xi_m, \xi_s)$  for each of the 8 largest firms, and also estimate a pooled value of  $(\xi_m, \xi_s)$  the industry as a whole. Table 8 reports these coefficient estimates, standard errors, and several measures of goodness-of-fit. All firms and the pooled industry have positive “taste” for the expectation of profits of a fracking design and negative “taste” for the standard deviation. That is, every firm appears to avoid fracking input choices with high variance. I can reject variance-neutrality for all firms and for the pooled industry. In dollar terms, firms make choices as if they are willing to accept a reduction in expected profits of \$0.62 to \$1.11 for a reduction of \$1 in the standard deviation of profits.

I report three goodness-of-fit statistics. The likelihood based pseudo- $R^2$ , which I refer to as *LLPR*, is defined as 1 minus the ratio of the optimized log-likelihood over the log-likelihood evaluated at the null hypothesis:

$$LLPR = 1 - \frac{\log \mathcal{L}(\hat{\xi}_m, \hat{\xi}_s)}{\log \mathcal{L}(0)}$$

This statistic is similar to a real  $R^2$  in that it varies between 0 and 1, with 0 indicating that the model does not fit any better than no model and 1 indicating that the model fits the data perfectly (see Train (2009)). This measure of fit indicates how far from “perfect” the fit actually is, but it does not have a “fraction of variance explained” interpretation the way a true  $R^2$  does. I also compute the correlation between the expected input use implied by the model’s estimated choice probabilities and actual input use, for both sand and water. If expected input use is similar to what is observed in the data, these correlations should be positive and (ideally) close to 1.

The fit of this model varies a fair amount across firms, but is generally modest. The pseudo- $R^2$  measures are small for all firms and for the pooled industry, suggesting that the best fitting values of the model’s parameters still require a lot of support from the logit errors to rationalize firm behavior. For six of the eight firms, the correlation between predicted sand use and realized sand use is positive, and for four firms it is at least 50%. The correlations between predicted and realized water use are much smaller, with only 4 firms having correlations at or above 20%. Though the coefficient estimates are all significantly different from zero, the low fit statistics suggest that

Table 8: Risk preference model estimates

Firm	$\widehat{\xi}_m$	$se(\widehat{\xi}_m)$	$\widehat{\xi}_s$	$se(\widehat{\xi}_s)$	# Wells	<i>LLPR</i>	$\rho_S$	$\rho_W$
Brigham	6.22	0.67	-5.59	0.68	111	0.15	0.03	0.05
Burlington	8.39	0.89	-9.95	0.99	102	0.30	0.51	0.35
Continental	7.59	0.47	-9.09	0.54	313	0.26	0.50	0.20
EOG	4.92	0.33	-6.08	0.47	339	0.15	-0.10	0.34
Hess	6.68	0.57	-8.30	0.66	143	0.27	0.57	0.25
Marathon	10.63	0.82	-13.88	0.97	209	0.41	0.61	0.12
Whiting	6.23	0.43	-10.09	0.68	247	0.32	-0.06	0.02
XTO	5.35	0.58	-7.02	0.68	101	0.23	0.37	-0.27
All	4.90	0.11	-6.62	0.14	2,604	0.19	0.47	0.24

Maximum likelihood estimates of the risk preference model model:

$$u_{ij} = P_i \left( \xi_m \mathbb{E}[DOP_{ij}] + \xi_s (\mathbb{V}[DOP_{ij}])^{\frac{1}{2}} \right) - c_i(S_j, W_j) + \epsilon_{ij}$$

$P_i$  is the price of oil for well  $i$ ,  $\mathbb{E}[DOP_{ij}]$  is the expectation of the present discounted value of oil production for well  $i$  when it is fracked using design  $j$ ,  $\mathbb{V}[DOP_{ij}]$  is the variance of the present discounted value of oil production for  $i$  under design  $j$ ,  $c_i(S_j, W_j)$  is the cost of implementing design  $j$  on well  $i$ , and  $\epsilon_{ij}$  is an iid logit shock. *LLPR* is a likelihood-based pseudo- $R^2$ :

$$LLPR = 1 - \frac{\log \mathcal{L}(\widehat{\xi}_m, \widehat{\xi}_s)}{\log \mathcal{L}(0, 0)}$$

where  $\mathcal{L}(\widehat{\xi}_m, \widehat{\xi}_s)$  is the likelihood of the model evaluated at the MLE and  $\mathcal{L}(0, 0)$  is the likelihood of the model evaluated at the null hypothesis.  $\rho_S$  and  $\rho_W$  are the correlations of actual sand and water use decisions with their predicted values from the model.



Table 9: Risk preferences model estimates, with interaction

Firm	$\widehat{\xi}_m$	$se(\widehat{\xi}_m)$	$\widehat{\xi}_s$	$se(\widehat{\xi}_s)$	$\widehat{\xi}_I$	$se(\widehat{\xi}_I)$	# Wells	<i>LLPR</i>	$\rho_S$	$\rho_W$
Brigham	9.79	1.06	-3.92	0.78	-1.19	0.25	111	0.18	0.01	0.02
Burlington	9.42	1.22	-9.71	1.00	-0.32	0.26	102	0.30	0.52	0.35
Continental	11.22	0.73	-8.51	0.57	-1.10	0.17	313	0.28	0.59	0.09
EOG	5.27	0.37	-6.03	0.48	-0.07	0.03	339	0.15	-0.10	0.35
Hess	7.03	0.78	-8.27	0.67	-0.09	0.14	143	0.27	0.57	0.25
Marathon	13.14	1.17	-13.26	1.01	-0.88	0.29	209	0.42	0.64	0.11
Whiting	7.31	0.53	-9.95	0.71	-0.22	0.06	247	0.32	0.00	0.01
XTO	4.06	0.69	-8.01	0.79	0.50	0.16	101	0.24	0.46	0.05
All	6.00	0.14	-6.68	0.15	-0.19	0.02	2,604	0.20	0.50	0.26

Maximum likelihood estimates of the risk preference model model:

$$u_{ij} = P_i \left( \xi_m \mathbb{E}[DOP_{ij}] + \xi_s (\mathbb{V}[DOP_{ij}])^{\frac{1}{2}} + \xi_I \mathbb{E}[DOP_{ij}] (\mathbb{V}[DOP_{ij}])^{\frac{1}{2}} \right) - c_i(S_j, W_j) + \epsilon_{ij}$$

$P_i$  is the price of oil for well  $i$ ,  $\mathbb{E}[DOP_{ij}]$  is the expectation of the present discounted value of oil production for well  $i$  when it is fracked using design  $j$ ,  $\mathbb{V}[DOP_{ij}]$  is the variance of the present discounted value of oil production for  $i$  under design  $j$ ,  $c_i(S_j, W_j)$  is the cost of implementing design  $j$  on well  $i$ , and  $\epsilon_{ij}$  is an iid logit shock. *LLPR* is a likelihood-based pseudo- $R^2$ :

$$LLPR = 1 - \frac{\log \mathcal{L}(\widehat{\xi}_m, \widehat{\xi}_s, \widehat{\xi}_I)}{\log \mathcal{L}(0, 0, 0)}$$

where  $\mathcal{L}(\widehat{\xi}_m, \widehat{\xi}_s, \widehat{\xi}_I)$  is the likelihood of the model evaluated at the MLE and  $\mathcal{L}(0, 0, 0)$  is the likelihood of the model evaluated at the null hypothesis.  $\rho_S$  and  $\rho_W$  are the correlations of actual sand and water use decisions with the predicted values from the model.

preferences that are linear in the mean and standard deviation of profits only explain a small portion of observed behavior.

I also estimate a version of this model which includes an interaction term between expected profits and the standard deviation of profits. While learning rules which are nonlinear in the mean and standard deviation do not appear in the existing learning literature, it is possible that true firm preferences over risk and reward are more complicated than a linear model can capture. By including an interaction between expected profits and the standard deviation of profits, I allow for risk preferences that may vary with the mean. Table 9 reports estimates of these models. The results are qualitatively the same as Table 8, with all firms showing risk aversion and all but one firm showing increasingly negative taste for risk as reward increases. In dollar terms, firms make choices as if they are willing to accept a reduction in expected profits of \$0.51 to \$2.50 for a reduction of \$1 in the standard deviation of profits. Goodness-of-fit measures are slightly better for these models than for the standard mean/variance models, though this is to be expected from the inclusion of an additional covariate.

Overall, Tables 8 and 9 provide evidence that firms tend to select fracking designs with higher expected profits and avoid fracking designs with higher standard deviation of profit. This behavior is not consistent with the notion that firms are actively exploring uncertain fracking designs, but it is consistent with passively learning firms that are constrained by organizational or financially motivated variance aversion.

## 5.2 Own-data bias

A different explanation for firms' apparent unwillingness to select the fracking design with the largest expected profits is that I am computing expectations with respect to different beliefs than those held by firms. There are many ways that a firm's beliefs may be different than the ones I calculate: firms may have biased prior beliefs about the role of fracking design and location, they may have simpler beliefs about the functional form relating fracking design and location to production, or my fracking cost and oil price data could be different from the costs and prices firms experience. However, using the data that I have, I am only able to test a simpler explanation. I assume that firms do have the belief structure I have described here, but do not necessarily treat all of the data available to them equally. In particular, firms may weigh data from their own experiences differently than data from the experiences of other firms that they observe through the public disclosure process. I refer to this explanation as "own-data bias".

To test for this phenomenon, I introduce a new parameter,  $\lambda \in (0, 1)$ , which represents the firm's relative weighting scheme. If  $\lambda = 0$ , the firm places no weight on the data generated by other firms and if  $\lambda = 1$ , the firm places no weight on its own data, relying entirely on outside data to learn. At  $\lambda = \frac{1}{2}$ , the firm puts equal weight on its own data and the data generated by others, which gives the preference model described in the previous section. For each value of  $\lambda$ , I can compute the expectation and standard deviation of *weighted* discounted profits for well  $i$  with fracking design  $j$ , for which I provide a calculation in the appendix. I then use these weighted profits in the same multinomial logit choice model described in the previous section, and refer to the choice model with weighted estimates as the weighted preference model.

In Table 10, I report maximum likelihood estimates of  $\lambda$ , as well as the other preference model coefficients, for the same specification in Table 8. The estimated value of  $\lambda$  is less than  $\frac{1}{2}$  for 6 of the 8 individual firms, and for 5 of those 6, the 95% confidence intervals do not include  $\frac{1}{2}$ . The pooled estimate is also less than  $\frac{1}{2}$  and its 95% confidence interval does not include  $\frac{1}{2}$ . Comparing Tables 8 and 10, the preference model coefficients do change slightly, but allowing for weighted beliefs does not affect the previous conclusion that all firms dislike uncertainty in the profits of a fracking input choice. Firms are willing to trade \$0.62 to \$0.96 in expected profits for a reduction of \$1 in the standard deviation of profits, which is a similar range to the model estimated in Table 8. The fit of the model in Table 10 is somewhat better than the model in Table 8, but it is still modest.

## 6 Information Disclosure Rules and the Speed of Learning

Firms in North Dakota can learn from each other because the data they file with the NDIC is made public. However, firms have the ability to keep secrets about the wells they operate for short periods of time. Under section 43-02-03-31 of the North Dakota Administrative code, firms may request that the NDIC delay its disclosure of their regulatory submissions for up to six months.<sup>34</sup> Every well in this dataset received confidential treatment from the NDIC, so firms' revealed preferences suggest that they value this temporary period of secrecy.

Well confidentiality exists to provide incentives to firms who drill in unexplored areas.<sup>35</sup> Upon the discovery of a previously unknown oil deposit, a firm can use its temporary period of confidentiality to acquire additional mineral rights in the area around the discovery. In the absence of a confidentiality law, it is likely that the firm's competitors would also try to acquire these mineral rights, and in the process make it impossible for the firm to earn a profit commensurate with the

<sup>34</sup>See <https://www.dmr.nd.gov/oilgas/webhelpfaq.asp>

<sup>35</sup>For an overview of well confidentiality, see Larsen (2011)).

Table 10: Weighted risk preference model estimates

Firm	$\widehat{\xi}_m$	$se(\widehat{\xi}_m)$	$\widehat{\xi}_s$	$se(\widehat{\xi}_s)$	$\widehat{\lambda}$	$se(\widehat{\lambda})$	# Wells	<i>LLPR</i>	$\rho_S$	$\rho_W$
Brigham	7.64	0.80	-7.37	0.85	0.30	0.05	111	0.17	0.08	0.11
Burlington	8.72	0.92	-10.47	1.02	0.42	0.06	102	0.30	0.53	0.33
Continental	8.95	0.58	-9.70	0.58	0.37	0.02	313	0.27	0.56	0.19
EOG	5.18	0.35	-7.25	0.55	0.17	0.04	339	0.17	-0.11	0.38
Hess	7.81	0.71	-9.26	0.74	0.52	0.06	146	0.27	0.58	0.20
Marathon	13.43	1.05	-15.40	1.09	0.36	0.04	212	0.43	0.64	0.10
Whiting	6.15	0.46	-9.64	0.69	0.05	0.06	247	0.32	-0.01	0.02
XTO	6.21	0.69	-7.86	0.77	0.58	0.07	101	0.25	0.39	-0.00
All	5.22	0.12	-7.02	0.15	0.36	0.02	2,605	0.19	0.50	0.28

Maximum likelihood estimates of the risk preference model model:

$$u_{ij} = P_i \left( \xi_m \mathbb{E}[DOP_{ij} | \lambda] + \xi_s (\mathbb{V}[DOP_{ij} | \lambda])^{\frac{1}{2}} \right) - c_i(S_j, W_j) + \epsilon_{ij}$$

$P_i$  is the price of oil for well  $i$ ,  $\mathbb{E}[DOP_{ij} | \lambda]$  is the expectation of the present discounted value of oil production for well  $i$  when it is fracked using design  $j$ ,  $\mathbb{V}[DOP_{ij} | \lambda]$  is the variance of the present discounted value of oil production for  $i$  under design  $j$ ,  $\lambda$  is the weighting parameter,  $c_i(S_j, W_j)$  is the cost of implementing design  $j$  on well  $i$ , and  $\epsilon_{ij}$  is an iid logit shock. *LLPR* is a likelihood-based pseudo- $R^2$ :

$$LLPR = 1 - \frac{\log \mathcal{L}(\widehat{\xi}_m, \widehat{\xi}_s, \widehat{\lambda})}{\log \mathcal{L}(0, 0, \frac{1}{2})}$$

where  $\mathcal{L}(\widehat{\xi}_m, \widehat{\xi}_s, \widehat{\lambda})$  is the likelihood of the model evaluated at the MLE and  $\mathcal{L}(0, 0, \frac{1}{2})$  is the likelihood of the model evaluated at the null hypothesis.  $\rho_S$  and  $\rho_W$  are the correlations of actual sand and water use decisions with the predicted values from the model.

risk it took in doing exploration in the first place.

In the context of shale exploration using fracking technology, well confidentiality can also reduce the quantity of information that firms have to learn from. Thus, in evaluating the impact of well confidentiality rules, the potential benefits of enhanced exploration activity must be weighed against the potential learning costs. In ongoing work, I am using the fracking design choice model developed in this paper to estimate the effects of well confidentiality rules on the speed of learning.

## 7 Conclusion

This paper provides one of the first empirical analyses of learning behavior in firms using operational choices, realized profits, and information sets. Oil companies in the North Dakota Bakken Shale learned to more efficiently use fracking technology between 2005-2011, increasing their capture of possible profits from 21% to 67%. These gains are mostly driven by improved fracking design choices over time. Contrary to the predictions of most theoretical models of learning, I do not find evidence that firms actively experiment in order to learn. Instead, I find evidence that firms prefer fracking input choices with lower variance, and are willing to give up \$0.62-1.11 in expected profits for a reduction of \$1 in the standard deviation of profits. Finally, several firms in my data appear to overweight data from their own operations relative to the data they observe from their competitors.

From a neoclassical economics perspective, it is surprising that these firms learn slowly and do not experiment, even though it is valuable to do so. The firms in this data operate in an industry known for its appetite for risk and use of advanced technology. However, they leave money on the table. Across the 2,699 wells in this data, the average well appears to forego \$11.2 million in profits on an *ex post* basis and \$5.8 million on an *ex ante* basis, resulting in \$14-30 billion in lost profits.

In future work, I plan to collect better data on the costs firms face, and the informational they acquire while drilling. This will allow me to relax my identification assumptions and construct better estimates of the profits firms did earn and could have earned. Additionally, I intend to collect data on the organizational structure of firms in the Bakken and use it to analyze the role that agency factors play in encouraging or inhibiting learning.

## References

- Aghion, Philippe, Patrick Bolton, Christopher Harris, and Bruno Jullien**, “Optimal Learning by Experimentation,” *The Review of Economic Studies*, 1991, 58 (4), pp. 621–654.
- Anand, Bharat N and Tarun Khanna**, “Do firms learn to create value? The case of alliances,” *Strategic management journal*, 2000, 21 (3), 295–315.
- Arrow, Kenneth J.**, “The Economic Implications of Learning by Doing,” *The Review of Economic Studies*, 1962, 29 (3), pp. 155–173.
- Baihly, Jason, Raphael Altman, and Isaac Aviles**, “Has the Economic Stage Count Been Reached in the Bakken Shale?,” in “SPE Hydrocarbon Economics and Evaluation Symposium” 2012.
- Benkard, C.L.**, “Learning and Forgetting: The Dynamics of Aircraft Production,” *The American Economic Review*, 2000.

- Bolton, Patrick and Christopher Harris**, “Strategic Experimentation,” *Econometrica*, 1999, 67 (2), pp. 349–374.
- Brezzi, Monica and Tze Leung Lai**, “Optimal learning and experimentation in bandit problems,” *Journal of Economic Dynamics and Control*, 2002, 27 (1), 87–108.
- Chapelle, Olivier and Lihong Li**, “An empirical evaluation of thompson sampling,” in “Advances in Neural Information Processing Systems” 2011, pp. 2249–2257.
- Conley, T.G. and C.R. Udry**, “Learning about a new technology: Pineapple in Ghana,” *The American Economic Review*, 2010, 100 (1), 35–69.
- Deutsch, John**, “Secretary of Energy Advisory Board Shale Gas Production Subcommittee Second Ninety Day Report,” Technical Report November 2011.
- Dickstein, M.J.**, “Efficient provision of experience goods: Evidence from antidepressant choice,” Working Paper 2013.
- Fetkovich, MJ**, “Decline curve analysis using type curves,” *Journal of Petroleum Technology*, 1980, 32 (6), 1065–1077.
- Foster, A.D. and M.R. Rosenzweig**, “Learning by doing and learning from others: Human capital and technical change in agriculture,” *Journal of Political Economy*, 1995, pp. 1176–1209.
- Gaswirth, Stephanie B.**, “Assessment of Undiscovered Oil Resources in the Bakken and Three Forks Formations, Williston Basin Province, Montana, North Dakota, and South Dakota, 2013,” Technical Report, U.S. Geological Survey April 2013.
- Gilje, E.**, “Does Local Access To Finance Matter?: Evidence from US Oil and Natural Gas Shale Booms,” Working Paper November 2012.
- Gittins, John C and David M Jones**, “A dynamic allocation index for the discounted multi-armed bandit problem,” *Biometrika*, 1979, 66 (3), 561–565.
- Graepel, Thore, Joaquin Q Candela, Thomas Borchert, and Ralf Herbrich**, “Web-scale bayesian click-through rate prediction for sponsored search advertising in microsoft’s bing search engine,” in “Proceedings of the 27th International Conference on Machine Learning (ICML-10)” 2010, pp. 13–20.
- Griliches, Zvi**, “Hybrid corn: An exploration in the economics of technological change,” *Econometrica, Journal of the Econometric Society*, 1957, pp. 501–522.
- Haushalter, G. David**, “Financing Policy, Basis Risk, and Corporate Hedging: Evidence from Oil and Gas Producers,” *The Journal of Finance*, 2000, 55 (1), 107–152.
- Hicks, Bruce E.**, “North Dakota Oil & Gas Update,” in “Presented for Dunn County Oil Day in Killdeer, ND on February 21, 2012” North Dakota Industrial Commission 2012.
- Hough, E and T McClurg**, “Impact of Geological Variation and Completion Type in the U.S. Bakken Oil Shale Play Using Decline Curve Analysis and Transient Flow Character,” in “AAPG International Conference and Exhibition, Milan, Italy, October 23-26, 2011” 2011.
- Kasy, Maximilian**, “Why experimenters should not randomize, and what they should do instead,” *Working Paper*, 2013.

- Kellogg, Ryan**, “Learning by Drilling: Interfirm Learning and Relationship Persistence in the Texas Oilpatch,” *The Quarterly Journal of Economics*, 2011, 126 (4), 1961–2004.
- Kogan, Leonid, Dimitris Papanikolaou, Amit Seru, and Noah Stoffman**, “Technological Innovation, Resource Allocation, and Growth,” Working Paper 17769, National Bureau of Economic Research January 2012.
- Larsen, Lamont C**, “Horizontal Drafting: Why Your Form JOA May Not Be Adequate for Your Companys Horizontal Drilling Program,” *Rocky Mtn. Min. L. Found. J.*, 2011, 48, 51.
- Levitt, Clinton J.**, “Learning through Oil and Gas Exploration,” *Working Paper*, November 2011.
- Levitt, Steven D., John A. List, and Chad Syverson**, “Toward an Understanding of Learning by Doing: Evidence from an Automobile Assembly Plant,” Working Paper 18017, National Bureau of Economic Research April 2012.
- Muehlenbachs, Lucija, Elisheba Spiller, and Christopher Timmins**, “Shale Gas Development and Property Values: Differences across Drinking Water Sources,” Working Paper 18390, National Bureau of Economic Research September 2012.
- Nordhaus, Alex Trembath Michael Shellenberger Ted and Jesse Jenkins**, “Where the Shale Gas Revolution Came From,” Technical Report, The Breakthrough Institute May 2012.
- Rasmussen, Carl Edward and Christopher KI Williams**, *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*, The MIT Press, 2005.
- Reiss, Peter C.**, “Economic and Financial Determinants of Oil and Gas Exploration Activity,” Working Paper 3077, National Bureau of Economic Research August 1989.
- Romer, Paul M.**, “Increasing Returns and Long-Run Growth,” *Journal of Political Economy*, 1986, 94 (5), pp. 1002–1037.
- Rusmevichientong, Paat and John N Tsitsiklis**, “Linearly parameterized bandits,” *Mathematics of Operations Research*, 2010, 35 (2), 395–411.
- Ryan, Bryce and Neal C Gross**, “The diffusion of hybrid seed corn in two Iowa communities,” *Rural sociology*, 1943, 8 (1), 15–24.
- Scott, Steven L**, “A modern Bayesian look at the multi-armed bandit,” *Applied Stochastic Models in Business and Industry*, 2010, 26 (6), 639–658.
- Shelley, Robert, Nijat Guliyev, and Amir Nejad**, “A Novel Method to Optimize Horizontal Bakken Completions in a Factory Mode Development Program,” in “SPE Annual Technical Conference and Exhibition” 2012.
- Srinivas, Niranjana, Andreas Krause, Sham M Kakade, and Matthias Seeger**, “Information-theoretic regret bounds for gaussian process optimization in the bandit setting,” *Information Theory, IEEE Transactions on*, 2012, 58 (5), 3250–3265.
- Stoyanov, Andrey and Nikolay Zubanov**, “Productivity spillovers across firms through worker mobility,” *American Economic Journal: Applied Economics*, 2012, 4 (2), 168–198.

**Thornton, R.A. and P. Thompson**, “Learning from experience and learning from others: An exploration of learning and spillovers in wartime shipbuilding,” *American Economic Review*, 2001, pp. 1350–1368.

**Train, Kenneth**, *Discrete choice methods with simulation*, Cambridge university press, 2009.

**Vidic, RD, SL Brantley, JM Vandenbossche, D Yoxthimer, and JD Abad**, “Impact of Shale Gas Development on Regional Water Quality,” *Science*, 2013, *340* (6134).

**Walls, Michael R. and James S. Dyer**, “Risk Propensity and Firm Performance: A Study of the Petroleum Exploration Industry,” *Management Science*, 1996, *42* (7), 1004–1021.

**Wieland, Volker**, “Learning by doing and the value of optimal experimentation,” *Journal of Economic Dynamics and Control*, 2000, *24* (4), 501–534.

## A Likelihood Calculation

### A.1 Step 1

Let  $\theta = (\alpha, \beta, \delta, \eta)$  represent the vector of the non-fracking parameters and let  $\phi = (\sigma_\epsilon, \sigma_\nu)$  represent the vector of the variance parameters. I compute the pseudo-observation  $g_i$  from  $(Y_{it}, X_{it})$ , conditional on  $\theta$  as

$$\begin{aligned} g_i &= \frac{1}{N_i} \sum_{t=1}^{N_i} (\log Y_{it} - X_{it}\theta) \\ &= \frac{1}{N_i} \sum_{t=1}^{N_i} (g(Z_i) + \epsilon_i + \nu_{it}) \\ &= f(Z_i) + \epsilon_i + \frac{1}{N_i} \sum_{t=1}^{N_i} \nu_{it} \end{aligned}$$

$g_i$  is the sum of the “true” effect of fracking and location on oil production and a normally distributed error with zero mean and variance  $\sigma_\epsilon^2 + \frac{1}{N_i}\sigma_\nu^2$ .

### A.2 Step 2

Conditional on the pseudo-observations  $g_i$ , the likelihood of  $(Y_{it}, X_{it})$  follows the standard formula for panel data with a random effect on each well. Let  $\psi(\cdot | \mu, \sigma)$  denote the normal likelihood with mean  $\mu$  and standard deviation  $\sigma$  and let  $e_{it} = \log Y_{it} - X_{it}\theta$ . Finally, let bolded capital letters represent vectors of the time series of a variable. The likelihood of observing  $(\mathbf{Y}_i, \mathbf{X}_i)$  conditional

on the parameters  $(\theta, \phi)$  and the unobserved impact of fracking  $g_i$  is

$$\begin{aligned}
\mathcal{L}(\mathbf{Y}_i, \mathbf{X}_i \mid g_i, \theta, \phi) &= \int \psi(\epsilon_i \mid 0, \sigma_\epsilon) \prod_{t=1}^{T_i} \psi(e_{it} - g_i - \epsilon_i \mid 0, \sigma_\nu) d\epsilon_i \\
&= \exp \left( -\frac{1}{2} \left[ \frac{1}{\sigma_\nu^2} \left( \sum_{t=1}^{T_i} (e_{it} - g_i)^2 - \frac{\sigma_\epsilon^2}{T_i \sigma_\epsilon^2 + \sigma_\nu^2} \left( \sum_{t=1}^{T_i} e_{it} - g_i \right)^2 \right) \right] \right) \\
&\quad \times \left( \left( T_i \frac{\sigma_\epsilon^2}{\sigma_\nu^2} + 1 \right) (2\pi \sigma_\nu^2)^{T_i} \right)^{-\frac{1}{2}} \\
&= \psi \left( g_i \mid \frac{1}{T_i} \sum_{t=1}^{T_i} e_{it}, \sigma_\epsilon^2 + \frac{1}{T_i} \sigma_\nu^2 \right) \\
&\quad \times \exp \left( \frac{1}{\sigma_\nu^2} \left( \left( \sum_{t=1}^{T_i} e_{it} \right)^2 \left( \frac{\sigma_\epsilon^2}{T_i \sigma_\epsilon^2 + \sigma_\nu^2} - \frac{1}{2T_i} \right) - \frac{1}{2} \sum_{t=1}^{T_i} e_{it}^2 \right) \right)^{-\frac{1}{2}} \\
&= \psi \left( g_i \mid \frac{1}{T_i} \sum_{t=1}^{T_i} e_{it}, \sigma_\epsilon^2 + \frac{1}{T_i} \sigma_\nu^2 \right) J(\mathbf{Y}_i, \mathbf{X}_i, T_i \mid \theta, \phi)
\end{aligned}$$

The first term in this final expression is simply a normal likelihood, evaluated at  $g_i$ , the effect of fracking and location for well  $i$ . The second term does not depend on  $g_i$ . Though  $g_i$  is unobserved, by the properties of GPR, the vector  $\mathbf{g}$  of  $g_i$ 's for all  $N$  wells is distributed multivariate normal with mean zero and variance  $K(\mathbf{Z} \mid \gamma)$ . Thus, I can integrate over the values of  $g_i$  to obtain the likelihood in terms of observable data and parameters. Let  $\mathbf{T}$  denote the vector of values of  $T_i$ ,  $\Sigma(\mathbf{T}, \phi)$  be a  $N$  by  $N$  matrix with  $\sigma_\epsilon^2 + \frac{1}{T_i} \sigma_\nu^2$  in the  $i$ -th diagonal position and zeros elsewhere and let  $\mu(\mathbf{Y}, \mathbf{X}, \mathbf{T}, \theta)$  be a vector with  $\frac{1}{T_i} \sum_{t=1}^{T_i} e_{it}$  in the  $i$ -th position. Then the full likelihood is:

$$\begin{aligned}
\mathcal{L}(\mathbf{Y}, \mathbf{X}, \mathbf{Z}) &= \int \psi(\mathbf{g} \mid \mathbf{0}, K(\mathbf{Z} \mid \gamma)) \prod_{i=1}^N \mathcal{L}(\mathbf{Y}_i, \mathbf{X}_i \mid g_i, \theta, \phi) dg_i \\
&= \left[ \prod_{i=1}^N J(\mathbf{Y}_i, \mathbf{X}_i, T_i \mid \theta, \phi) \right] \int \psi(\mathbf{g} \mid \mathbf{0}, K(\mathbf{Z} \mid \gamma)) \psi(\mathbf{g} \mid \mu(\mathbf{Y}, \mathbf{X}, \mathbf{T}, \theta), \Sigma(\mathbf{T}, \phi)) d\mathbf{g} \\
&= \left[ \prod_{i=1}^N J(\mathbf{Y}_i, \mathbf{X}_i, T_i \mid \theta, \phi) \right] \psi(\mu(\mathbf{Y}, \mathbf{X}, \mathbf{T}, \theta) \mid \mathbf{0}, \Sigma(\mathbf{T}, \phi) + K(\mathbf{Z} \mid \gamma))
\end{aligned}$$

where the last line comes as a result of equations A.7 and A.8 from Rasmussen and Williams (2005). Having integrated out the unobserved values  $g_i$ , the full likelihood is completely in terms of the observed data  $(\mathbf{Y}, \mathbf{X}, \mathbf{T})$ , the parameter vectors  $\theta$  and  $\phi$ , and the covariance matrix  $K(\mathbf{Z} \mid \gamma)$  of the nonparametric effect of fracking and location on oil production.

## B Expected Present Discounted Value of Oil Production

I compute *ex post* expectations for all wells, and I compute *ex ante* expectations for wells fracked by firms with sufficiently large information sets. I require that a firm's information set has at least 50 wells and at least 300 well-months of production. This limits the set of wells I can analyze, and the earliest wells with information sets this large do not appear until the fourth quarter of 2007.



I compute  $\mathbb{E}[DOP_{ij}]$  using both expectation operators for a 10 by 10 grid of possible frack designs  $j$ , with sand use between 0 and 650 lbs per foot and water use between 0 and 750 gals per foot. These grid points cover 95% of observed sand choices and 99% of observed water choices. By the normality assumptions in the production function model, the joint distribution of log-production for well  $i$  under fracking design  $j$  over  $T$  months of existence (call this  $\log \tilde{Y}_{ij}$ ) is multivariate normal, with mean  $\mu_{ij}$  and covariance  $\Sigma_{ij}$  given by:

$$\begin{aligned}\mu_{ij} &= \tilde{\mathbf{X}}_i \theta + \tilde{g}(Z_{ij}) \\ \Sigma_{ij} &= \tilde{\mathbf{X}}_i \Sigma^\theta \tilde{\mathbf{X}}_i^\top + (\sigma_\epsilon^2 + \sigma_{g,ij}^2) \mathbf{1}_T + \sigma_\nu^2 \mathbf{I}_T\end{aligned}$$

where  $\tilde{\mathbf{X}}_i$  is a matrix of well  $i$ 's static characteristics and a vector of log-age values from 1 month to  $T$  months,  $\tilde{g}(\cdot)$  is the estimated GPR,  $Z_{ij}$  is the a vector of design ( $S_j, W_j$ ) and latitude and longitude for well  $i$ ,  $\Sigma^\theta$  is the covariance matrix for the estimates of  $\theta$ ,  $\sigma_{g,ij}^2$  is the estimated variance of the GPR at  $Z_{ij}$ ,  $\mathbf{1}_T$  is a  $T$  by  $T$  matrix of ones, and  $\mathbf{I}_T$  is a  $T$  by  $T$  identity matrix. With this construction, I am assuming that the variances for  $\epsilon$  and  $\nu$  are estimated perfectly (i.e., there is no term in  $\Sigma_{ij}$  that accounts for variance in those estimates).

Because  $\log \tilde{Y}_{ij}$  is multivariate normal, the distribution of the *level* of production over time,  $\tilde{Y}_{ij}$ , is multivariate log-normal with the same parameters. The mean vector and covariance matrix of this distribution are:

$$\begin{aligned}\tilde{\mu}_{ij} &= \exp\left(\mu_{ij} + \frac{1}{2} \mathcal{D}(\Sigma_{ij})\right) \\ \left[\tilde{\Sigma}_{ij}\right]_{kl} &= \exp\left([\mu_{ij}]_k + [\mu_{ij}]_l + \frac{1}{2}([\Sigma_{ij}]_{kk} + [\Sigma_{ij}]_{ll})\right) (\exp([\Sigma_{ij}]_{kl}) - 1)\end{aligned}$$

where  $\mathcal{D}(\cdot)$  represents the diagonal vector of a square matrix and  $[M]_{xy}$  is the  $(x, y)$ -th entry of a matrix  $M$ .<sup>36</sup> Finally,  $\mathbb{E}[DOP_{ij}]$  is:

$$\mathbb{E}[DOP_{ij}] = \sum_{t=1}^T \rho^t \tilde{\mu}_{ijt}$$

A similar calculation is available for the variance of present discounted oil production:

$$\begin{aligned}\mathbb{V}[DOP_{ij}] &= \mathbb{V}\left[\sum_{t=1}^T \rho^t \tilde{Y}_{ijt}\right] \\ &= \sum_{t_1=1}^T \sum_{t_2=1}^T \rho^{t_1+t_2} \left[\tilde{\Sigma}_{ij}\right]_{t_1, t_2}\end{aligned}$$

## C Weighted Gaussian Process Estimates

Recall that the mean and variance of the Gaussian process estimates of  $f$  at the point  $\tilde{Z}$  are given by:

$$\begin{aligned}\mathbb{E}\left[f(\tilde{Z}) \mid g, \mathbf{Z}, \gamma\right] &= k(\tilde{Z} \mid \gamma)^\top K(\gamma)^{-1} g \\ \mathbb{V}\left[f(\tilde{Z}) \mid g, \mathbf{Z}, \gamma\right] &= k(\tilde{Z} \mid \gamma)^\top K(\gamma)^{-1} k(\tilde{Z} \mid \gamma)\end{aligned}$$

<sup>36</sup>By the properties of the log-normal distribution, the mean and standard deviation of production are closely related, with the standard deviation equal to the mean times the exponent of the variance minus 1. This means that the ‘‘correlation’’ between the mean and standard deviation of production, computed across designs  $j$  will be positive by construction.

where  $k(\tilde{Z} | \gamma) = (k(Z_1, \tilde{Z} | \gamma) \dots k(Z_N, \tilde{Z} | \gamma))^\top$ ,  $K(\gamma)$  is the matrix of pairwise kernel distances for each point in  $\mathbf{Z}$  and  $g = (g_1 \dots g_N)^\top$ . To compute a *weighted* mean and variance, I introduce a weighting matrix function,  $L(\lambda)$ , and compute a weighted estimate of the mean and variance:

$$\begin{aligned}\mathbb{E} [f(\tilde{Z}) | g, \mathbf{Z}, \gamma, \lambda] &= k(\tilde{Z} | \gamma)^\top L(\lambda)^\top K(\gamma)^{-1} g \\ \mathbb{V} [f(\tilde{Z}) | g, \mathbf{Z}, \gamma, \lambda] &= k(\tilde{Z} | \gamma)^\top L(\lambda)^\top K(\gamma)^{-1} L(\lambda) k(\tilde{Z} | \gamma)\end{aligned}$$

The weighting matrix function  $L(\lambda)$  biases these estimates towards a firm's own experiences when  $\lambda$  is closer to 0 and towards other firms' experiences when  $\lambda$  is closer to 1. In particular, if  $(k_0(\gamma), K_0(\gamma), g_0)$  are the subsets of  $k(\gamma), K(\gamma), g$  computed using only the firm's own wells, and  $(k_1(\gamma), K_1(\gamma), g_1)$  are the subsets computed using only other firms' wells, then the weighted estimates satisfy 3 relationships:

1. At  $\lambda = 0$ , the weighted estimates are equal to the estimates computed using the subset of wells the firm operated:

$$\begin{aligned}k(\tilde{Z} | \gamma)^\top L(0)^\top K(\gamma)^{-1} g &= k_0(\tilde{Z} | \gamma)^\top K_0(\gamma)^{-1} g_0 \\ k(\tilde{Z} | \gamma)^\top L(0)^\top K(\gamma)^{-1} L(0) k(\tilde{Z} | \gamma) &= k_0(\tilde{Z} | \gamma)^\top K_0(\gamma)^{-1} k_0(\tilde{Z} | \gamma)\end{aligned}$$

2. At  $\lambda = \frac{1}{2}$ , the weighted estimates are equal to the unweighted estimates:

$$\begin{aligned}k(\tilde{Z} | \gamma)^\top L\left(\frac{1}{2}\right)^\top K(\gamma)^{-1} g &= k(\tilde{Z} | \gamma)^\top K(\gamma)^{-1} g \\ k(\tilde{Z} | \gamma)^\top L\left(\frac{1}{2}\right)^\top K(\gamma)^{-1} L\left(\frac{1}{2}\right) k(\tilde{Z} | \gamma) &= k(\tilde{Z} | \gamma)^\top K(\gamma)^{-1} k(\tilde{Z} | \gamma)\end{aligned}$$

3. At  $\lambda = 1$ , the weighted estimates are equal to the estimates computed using the subset of wells the firm did not operate:

$$\begin{aligned}k(\tilde{Z} | \gamma)^\top L(1)^\top K(\gamma)^{-1} g &= k_1(\tilde{Z} | \gamma)^\top K_1(\gamma)^{-1} g_0 \\ k(\tilde{Z} | \gamma)^\top L(1)^\top K(\gamma)^{-1} L(1) k(\tilde{Z} | \gamma) &= k_1(\tilde{Z} | \gamma)^\top K_1(\gamma)^{-1} k_1(\tilde{Z} | \gamma)\end{aligned}$$

At intermediate values of  $\lambda$ ,  $L(\lambda)$  interpolates between these extremes. To accomplish this,  $L(\lambda)$  takes this form:

$$L(\lambda) = \begin{bmatrix} L_1(\lambda) \mathbf{I}_{n_0} & L_2(\lambda) K_{01}(\gamma) K_{11}(\gamma)^{-1} \\ L_3(\lambda) K_{10}(\gamma) K_{00}(\gamma)^{-1} & L_4(\lambda) \mathbf{I}_{n_1} \end{bmatrix}$$

where  $n_0$  is the number of wells in the firm's information set that it operated,  $n_1$  is the number of wells that other firms operated, the matrices  $K_{00}(\gamma), K_{01}(\gamma), K_{10}(\gamma), K_{11}(\gamma)$  are submatrices of  $K(\gamma)$ :

$$K(\gamma) = \begin{bmatrix} K_{00}(\gamma) & K_{01}(\gamma) \\ K_{10}(\gamma) & K_{11}(\gamma) \end{bmatrix}$$

and the functions  $L_1, L_2, L_3, L_4$  are

$$\begin{aligned}L_1(\lambda) &= 1 + \lambda - 2\lambda^2 \\ L_2(\lambda) &= -\lambda + 2\lambda^2 \\ L_3(\lambda) &= 1 - 3\lambda + 2\lambda^2 \\ L_4(\lambda) &= 3\lambda - 2\lambda^2\end{aligned}$$

Thus,  $L(\lambda)$  is a quadratic interpolation between  $L(0)$ , which selects out the firm's own wells, and  $L(1)$ , which selects out all other firms' wells.