

Test set reuse

Statistics prescribes the iron vault for test data. But the empirical reality of machine learning benchmarks couldn't be further from the prescription. Repeated adaptive testing brings theoretical risks and practical power.

5	Test set reuse	1
5.1	Test set reuse in machine learning benchmarks	2
	Adaptive data analysis	2
	Adaptive analysts	3
5.2	Guarantees of the holdout method under adaptivity	4
	Climbing the leaderboard without looking at the data	6
5.3	Alternatives to the holdout method	9
5.4	Freedman's paradox	11
	Variable selection	12
5.5	Notes	14

Source: The Emerging Science of Machine Learning Benchmarks. M. Hardt, 2025. URL: <https://mlbenchmarks.org>. Compiled on 2025-05-01.

In the previous chapter, we derived powerful theoretical guarantees for the holdout method. Unfortunately, these guarantees hold only under one-time use. If a researcher builds a model based on prior interactions with the holdout set, the holdout set loses its formal guarantees. In this chapter, we develop a theoretical model that accommodates how people actually use the holdout method in practice. We then work out generalization bounds in this setting, contrasting them with the results of the previous chapter.

5.1 *Test set reuse in machine learning benchmarks*

Machine learning benchmarks typically have a fixed test set known to researchers. Researchers freely use the test set for evaluation purposes. Scientific papers report new achievements on the test set. Subsequent work builds on prior evaluations. The central goal is to move ahead the *state of the art*, that is, to advance the performance of the best known model on the test set. In doing so, researchers work with the test set incrementally. They do what Duda and Hart called *training on the test set*. Researchers incrementally refine their model by repeated evaluations against the test set.

We call scientific analyses that depend on the test set *adaptive*. Results that in some way depend on the test set have been adapted to the test set. Likewise, *adaptivity* refers to the practice of using a test set incrementally and interactively to inform scientific analyses and processes, such as model building.

Adaptive data analysis

To reason about adaptivity, it is helpful to frame the problem as an interaction between two parties. One party holds the dataset S . Think of this party as implementing the holdout method. The other party—called *analyst*—can *query* the dataset by requesting an estimate of the risk $R(f)$ for a given predictor f on the dataset S . In reality, the two parties might be one and the same researcher. Nevertheless, conceptually it will be quite helpful to think of the problem as an interaction between two parties.

The standard holdout method returns the empirical risk $R_S(f)$ given a query function $f: \mathcal{X} \rightarrow \mathcal{Y}$. But we'll allow for other methods that don't necessarily output the empirical risk. It turns out that there are interesting alternatives to the standard holdout mechanism that enjoy stronger guarantees. Intuitively, these alternatives limit the amount of information about

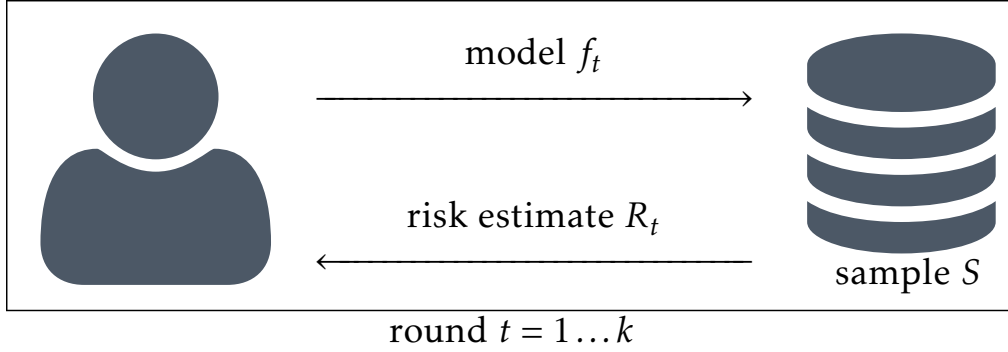


Figure 5.1: An adaptive analyst interacts with a dataset repeatedly.

the holdout set revealed by each estimate.

Throughout this chapter, we restrict our attention to the case of the zero-one loss and binary prediction, although the theory extends to other settings.

The two parties interact for some number k of rounds, thus creating a sequence of adaptively chosen predictors f_1, \dots, f_k . Keep in mind that this sequence depends on the dataset! In particular, when S is drawn at random, f_2, \dots, f_k become random variables, too, that are in general not independent of each other.

Adaptive analysts

Formally, an adaptive analyst is an algorithm \mathcal{A} that, given a sequence $f_1, R_1, \dots, f_t, R_t$ of queries and responses, returns a new query $f_{t+1}: \mathcal{X} \rightarrow \mathcal{Y}$, where we let $f_1 = \mathcal{A}(\emptyset)$ be the first query the analyst chooses. We assume that the analyst \mathcal{A} is a deterministic algorithm.

A useful idea is that we can represent an adaptive analyst \mathcal{A} as a tree. The root node is labeled by $f_1 = \mathcal{A}(\emptyset)$, i.e., the first function that the analyst queries without any context. The holdout mechanism then returns a response R_1 . This response is generally a sample quantity like the empirical risk. Note that the empirical risk $R_S(f_1)$ can only take on $n+1$ possible values. This is because we consider the zero-one loss, which can only take the values in the set $\mathcal{R} = \{0, 1/n, 2/n, \dots, 1\}$. So, it makes sense to assume that $R_1 \in \mathcal{R}$ takes values in the same set. We can always achieve this by rounding R_1 to the nearest multiple of $1/n$ without changing the response significantly.

Each possible response value $r \in \mathcal{R}$ creates a new child node in the tree corresponding to the function $f = \mathcal{A}(r)$ that the analyst queries when receiving

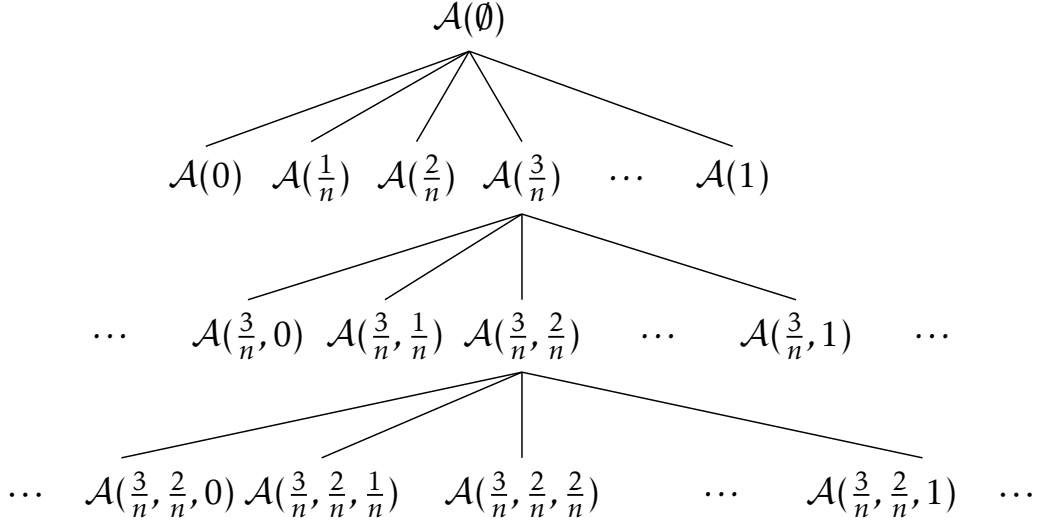


Figure 5.2: Constructing a tree of depth k and degree $n + 1$ representing an adaptive analyst. Each node corresponds to the predictor the analyst chooses based on the responses seen so far.

the response r to the first query f_1 . This gives us $n + 1$ children to the root node, one for each possible response $r \in \mathcal{R}$. Each child corresponds to one function. In this manner, we recursively continue the process until we have a tree of depth k and degree $n + 1$.

Note that this tree depends only on the analyst \mathcal{A} and how it responds to all the possible transcripts that can occur in the interaction with a holdout set. The tree does not depend on the random sample, however. The tree is therefore data-independent. It's an explicit representation of the algorithm \mathcal{A} . This property is a useful tool for proving generalization bounds in the adaptive setting.

5.2 Guarantees of the holdout method under adaptivity

We can now derive guarantees for the holdout method in the adaptive analyst model. The idea is to fix an analyst and apply the analysis from Chapter 4 to all the functions appearing in the tree corresponding to the analyst.

Proposition 1. *For any sequence of k adaptively chosen predictors f_1, \dots, f_k , the*

holdout method satisfies with probability $1 - \delta$,

$$\mathbb{P} \left\{ \max_{1 \leq t \leq k} |\Delta_S(f_t)| \leq \sqrt{\frac{k \log(4(n+1)/\delta)}{2n}} \right\} \geq 1 - \delta.$$

Proof. Fix an adaptive analyst \mathcal{A} and the corresponding tree. The tree is of size

$$1 + (n+1) + (n+1)^2 + \dots + (n+1)^k = \frac{(n+1)^{k+1} - 1}{n} \leq 2(n+1)^k.$$

Let F be the set of functions appearing at any of the nodes in the tree. Since each node has one function, we have

$$|F| \leq 2(n+1)^k.$$

Moreover, the set of functions is fully determined by the algorithm \mathcal{A} and therefore does not depend on any sample. Now, pick a random sample S of size n . The model selection guarantee for the holdout method from the previous chapter gives us for every $\delta > 0$,

$$\mathbb{P} \left\{ \max_{f \in F} |\Delta_S(f)| \leq \sqrt{\frac{\log(2|F|/\delta)}{2n}} \right\} \geq 1 - \delta.$$

Plugging in the upper bound on $|F|$,

$$\log(2|F|/\delta) \leq k \log(4(n+1)/\delta).$$

This completes the proof. □

In the typical regime where $k \geq \log(n)$ is at least logarithmic in n and $\delta \geq 1/n$ is not too small, the bound on the maximum error simplifies to $O(\sqrt{k/n})$. This contrasts with the bound $O(\sqrt{\log(k)/n})$ that we obtained in the non-adaptive setting.

Analyst	Error bound
non-adaptive	$O(\sqrt{\log(k)/n})$
adaptive	$O(\sqrt{k/n})$

The dependence on k in the adaptive setting is exponentially worse than in the non-adaptive setting. This raises the question if there's a way to do better.

Climbing the leaderboard without looking at the data

So far, the guarantee of the standard holdout mechanism that we have in the adaptive case is exponentially worse in k compared with the non-adaptive case. As it turns out, this is inevitable in the worst case. What we'll work out is a worst-case *lower bound* on the gap between risk and empirical risk in the adaptive setting. To prove such a lower bound it is enough to exhibit an adaptive analyst that forces a large gap.

Indeed, there is a fairly natural sequence of k adaptively chosen predictors, resembling the practice of ensembling, on which the empirical risk diverges from the risk by at least $\Omega(\sqrt{k/n})$. This matches the upper bound from the previous section up to a constant factor. In particular, with $k \approx n$ queries, we can force a constant generalization gap. The lower bound is *worst-case*: It only holds for this one analyst and doesn't say anything about the typical behavior of the holdout method.

Throughout, we focus on zero-one loss in a binary prediction problem. The core idea extends to many other settings.

Overfitting by ensembling:

1. Choose k random binary predictors $f_1, \dots, f_k: \mathcal{X} \rightarrow \{0, 1\}$.
2. Compute the set $I = \{i \in [k]: R_S[f_i] < 1/2\}$.
3. Output the predictor $f = \text{majority}\{f_i: i \in I\}$ that takes a majority vote over all the predictors computed in the second step.

The key idea of the algorithm is to select all the predictors that have accuracy strictly better than random guessing. This selection step creates a bias that gives each selected predictor an advantage over random guessing. The majority vote in the third step amplifies this initial advantage into a larger advantage that grows with k . If we want to be a bit more clever, we don't have to throw away any functions at all. If $R_S[f_i] > 1/2$, we know that flipping all the bits gives $R_S(1 - f_i) < 1/2$. So, we can include $1 - f_i$ in the ensemble. This saves us a factor two in the number of queries we make.

In practice, we can do a bit better still by weighting each function with its advantage over random guessing. The larger the advantage, the larger the weight of the function in the ensemble. This makes intuitive sense and leads

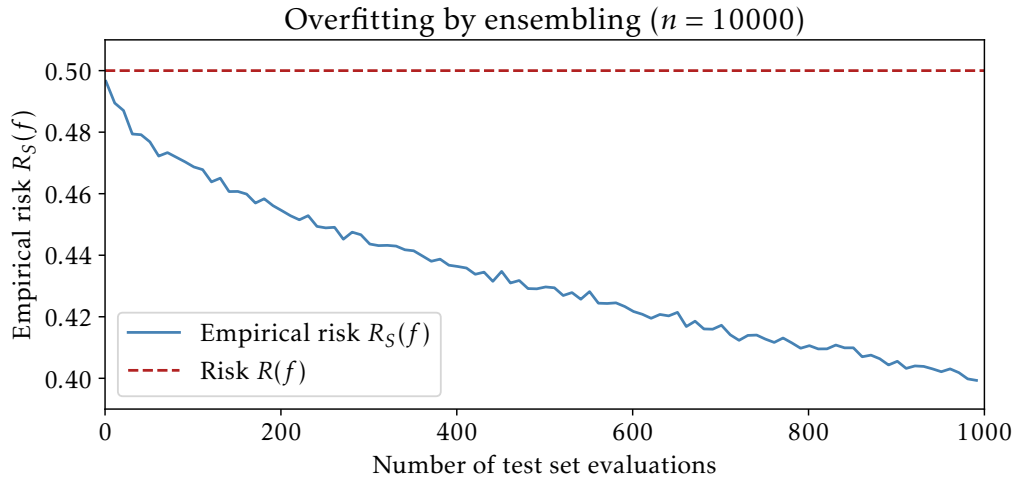


Figure 5.3: Overfitting by ensembling on a test set of size $n = 10000$. Thousand queries force a generalization gap of 10%.

to some improvements. What this algorithm essentially does is *boosting*. Boosting is a well-known technique for turning weak predictors into strong predictors.

The next proposition confirms that indeed this strategy finds a predictor whose empirical risk is bounded away from $1/2$ (random guessing) by a margin of $\Omega(\sqrt{k/n})$. Since the predictor does nothing but taking a majority vote over random functions, its risk is of course no better than $1/2$.

Proposition 2. *For sufficiently large $k \leq n$, overfitting by ensembling returns a predictor f whose classification error satisfies with probability $1/3$,*

$$R_S(f) \leq \frac{1}{2} - \Omega(\sqrt{k/n}).$$

In particular, $\Delta_S(f) \geq \Omega(\sqrt{k/n})$.

Proof. (Sketch)

The formal proof is a bit tedious, but it's good to develop the intuition.

On a randomly chosen function the number of errors on the test set S is distributed according to the binomial distribution $B(n, 1/2)$. We expect to see $n/2$ errors and one standard deviation is $\sqrt{n}/2$. This is because the variance of $B(n, p)$ is $np(1-p)$.

When we observe that $R_S(f_t) < 1/2$ for a randomly chosen f_t , we know that it makes fewer than $n/2$ errors on the test set. But apart from this condition, the errors are still randomly distributed. So, we get the distribution of $B(n, 1/2)$ conditional on falling below its mean. The mean of this conditional distribution is smaller than $n/2 - \sqrt{n}/2$. If it's below the mean, it'll be so by at least one standard deviation more than half the time.

This is the first part of the argument. Our selection step biases the functions to make about \sqrt{n} fewer errors on the test set than a random function.

How many functions do we select? In other words, how large is the index set I ? Since the binomial distribution is symmetric around its mean, the probability of any given function being selected is about $1/2$. So we expect the index set to be of size $m = k/2$. In fact, m concentrates around $k/2$. Therefore, we can be pretty sure that we select at least, say, $k/4$ functions. Up to constant factors, we can think of m and k as being the same.

To summarize, we know that we select many functions and each function we select is a bit better than random guessing.

The last step is to show that the majority vote amplifies this small advantage over random guessing. This is the same kind of argument we do when we reason about ensembling methods. Taking the majority vote over many experts—each a bit better than random guessing—results in a more accurate prediction.

So, consider the majority predictor $f = \text{majority}\{f_i : i \in I\}$. Fix a data point $(x, y) \in S$. What is the probability that f makes a mistake on x ? For convenience, assume $y = 1$. The case of $y = 0$ is the same. By definition, the majority classifier makes a mistake if more than half the functions in I vote 0. Formally,

$$\sum_{i \in I} f_i(x) < \frac{m}{2}.$$

What is the chance that this happens? The errors that each f_i makes on the test set are randomly located, since our selection step is invariant under permutation of the test set. Therefore, each function we select is correct on x with probability $1/2 + \epsilon$ for some $\epsilon > 1/\sqrt{n}$.

This means that the random variable $X = \sum_{i \in I} f_i(x)$ follows the binomial distribution $B(m, 1/2 + \epsilon)$. Its mean and variance are

$$\mathbb{E} X = \frac{m}{2} + \epsilon m, \quad \mathbb{V} X = m(1/2 + \epsilon)(1/2 - \epsilon) \leq m/4.$$

For the majority vote f to be wrong, the random variable X has to be below

its mean by an additive ϵm . Since a standard deviation is less than \sqrt{n} , this deviation is equivalent to $\epsilon\sqrt{m}$ in units of standard deviation. Note that

$$\epsilon\sqrt{m} = \sqrt{\frac{m}{n}} \approx \sqrt{\frac{k}{n}}.$$

Since we assume $k \leq n$, we're talking about less than one standard deviation. For large enough n and constant p , the binomial distribution $B(n, p)$ behaves very much like a normal distribution with mean np and variance $np(1-p)$. Within one standard deviation of its mean, the normal distribution acts a lot like a uniform distribution. In particular, we can calculate that a deviation of $\sqrt{k/n}$ below its mean has probability less than $1/2 - c\sqrt{k/n}$ up to some positive constant $c > 0$ in front of the $\sqrt{k/n}$ term.

So, we can conclude that the probability of a mistake by our majority vote is at most

$$\mathbb{P}\left\{\sum_{i \in I} f_i(x) < \frac{m}{2}\right\} \leq \frac{1}{2} - c\sqrt{\frac{k}{n}},$$

for some positive constant $c > 0$. Written differently, this means

$$R_S(f) \leq \frac{1}{2} - c\sqrt{\frac{k}{n}}.$$

That's what we wanted to show. This concludes the proof sketch.

□

Zooming out again, ensembling by majority voting shows that the problem of “training on the testing data” is real. In the worst case, holdout data can quickly lose its guarantees.

5.3 *Alternatives to the holdout method*

If this pessimistic bound manifested in practice, popular benchmark datasets would quickly become useless. You might wonder if there's anything we can do about the pessimistic lower bound we just encountered. So far, we've analyzed the standard holdout method that returns the exact empirical risk $R_S(f)$ given a query function f . There's hope that alternative mechanisms might have stronger guarantees. For example, we could release an approximate answer rather than an exact answer. We have to be a bit careful

though. The argument from the previous section extends to the case where the holdout method only gives approximate answers so long as these are within $o(1/\sqrt{n})$ from the exact answer. The reason is that $1/\sqrt{n}$ is roughly one standard deviation of the empirical risk. If all answers deviate less than one standard deviation, the argument is roughly the same. So, the solution isn't as easy as rounding the answer to four digits of precision. This is a common heuristic in machine learning competitions, but it has no formal guarantees.

Nevertheless, a related idea works. If we carefully add noise to each answer we can improve dependence on k from \sqrt{k} to $k^{1/4}$. This alternative holdout method permits a quadratic number $k = n^2$ of queries before it becomes useless in the worst-case.

The proof of this result is surprisingly tricky and requires several tools that are beyond the scope of this chapter. The key idea is based on the privacy guarantee *differential privacy*. If you can make the holdout mechanism differentially private, then no adaptive analyst can learn about the test set enough to overfit to it. This is a theorem that forms the basis for the alternative mechanism that uses noise addition to achieve the quadratic number of queries.

Unfortunately, researchers also showed that in the worst-case no holdout mechanism can give an answer to more than a quadratic number of queries with small constant error on every query. Keep in mind this is a worst-case impossibility result. It does not say what we should expect in practice under typical conditions.

The bounds in the previous chapter are optimistic. They are strong, but they require an assumption that is clearly violated in practice. The bounds in this chapter go about it from the other end. They allow for powerful adaptive analyses, like those you'll see in practice. The downside is that we end up with fairly pessimistic worst-case guarantees. If we typically encountered these bounds in practice, holdout sets would have a serious problem. Both perspectives are useful though. One tells us the best that we can hope for. The other shows the worst that could happen.

In the next chapters we'll look at the empirical phenomena around test set reuse in machine learning research. Then we'll return to theoretical arguments that better capture the empirical phenomena. Before we get there it's worth looking at the problem of adaptivity in the context of statistics.

5.4 Freedman's paradox

The problem with the holdout method we just saw has a close cousin in the field of statistics that goes by the name of Freedman's paradox. The statistician David Freedman pointed out this problem in 1983, although he thought that his observation was neither new, nor a paradox.

The problem routinely arises in the context of variable selection for statistical modeling. If we first select variables in a data-dependent way and then fit a model on the selected variables, the model will appear to be better than it is. At a high-level this is directly analogous to the ensembling procedure we discussed above. The details are different though and require a bit of statistical nomenclature that we'll develop next.

Freedman considers a linear regression problem

$$y = X\beta + \eta,$$

where the error term $\eta \sim N(0, \sigma^2)$ is sampled from a centered Gaussian distribution with variance σ^2 . The matrix X has shape $n \times d$ where n is the number of observations and d is the number of features. The vector $\beta \in \mathbb{R}^d$ corresponds to *unknown* the true solution of the equation system that we'd like to recover. We only observe the matrix X and the vector y .

Least squares regression solves the objective

$$\min_{\hat{\beta} \in \mathbb{R}^d} \|X\hat{\beta} - y\|^2.$$

Let $\hat{y} = X\hat{\beta}$ denote the optimal solution to least squares. Due to the noise in the system of equations, we can't hope to recover β exactly.

The situation statisticians really want to avoid is that $\beta = 0$ but somehow we come away thinking that $\beta \neq 0$. The case $\beta = 0$ corresponds to the situation where there is no signal in the data. The vector y is just a random normal vector with no dependence on X . Statisticians call this a *null model*. A null model is simply a data-generating distribution corresponding to the case where there is nothing to be discovered. In contrast, the case $\beta \neq 0$ means that there is some true relationship between X and y that we would like to discover and report.

To test if we're dealing with the null model, statisticians use hypothesis tests. A common one is the F -test that is based on the *observed F-value*

$$F^{\text{obs}} = \frac{\|\hat{y} - \bar{y}\|^2/d}{\|y - \hat{y}\|^2/(n-d-1)}.$$

Here, \bar{y} is the mean $\bar{y} = \sum_{i=1}^n y_i$ of our observations. We call F^{obs} the *observed* F -value, since it's the one we compute from our sample.

The observed F -value is the *test statistic*. The way the hypothesis test works is that we're going to reject the null model if the observed F -value is above a certain threshold, call it τ . Under the null model, the observed F -value follows an F -distribution. It doesn't matter what exactly this distribution is. It's just some continuous distribution we can compute.

The probability that the F -value we observe exceeds some threshold τ is given by $1 - \text{CDF}(\tau)$, where CDF is the cumulative density function of the F -distribution. The probability that the F -value exceeds the value we observed is therefore

$$p = 1 - \text{CDF}(F^{\text{obs}}).$$

Here, the tail probability p is called the *p-value* of the test. The p -value is the probability under the null model that the F -distribution exceeds the value F^{obs} we observe.

That means p -values are tail probabilities under the null model.

A fact called *integral probability theorem* implies that $\text{CDF}(F^{\text{obs}})$ follows the uniform distribution $\text{Uniform}([0, 1])$ over the interval $[0, 1]$. Therefore, the distribution of a p -value under the null model is the uniform distribution. This is the only fact about p -values we'll need.

So far, everything is working as intended. The p -values we see under the null model are uniformly distributed as they should. The chance of seeing a p -value as small as 0.05 is, by construction, at most 5%. This is the now infamous *significance level* at which statisticians will *reject the null hypothesis*. The idea is that seeing $p \leq 0.05$ renders the null model implausible.

Variable selection

Here's the problem. Our goal is not only to avoid rejecting the null model when $\beta = 0$. This alone would be easy: Never reject anything. But our goal is also to actually reject the null model when $\beta \neq 0$ and so we should. The higher the dimension d of our regression problem, the harder it generally is to find signal in the data. Therefore, it is common practice to first select some promising variables from the set of all variables, before fitting a model on the selected variables.

To select variables, we can again perform hypothesis tests. Specifically, we can test individual coefficients in the regression model to see if they are

significantly large using the T -statistic:

$$T_j^{\text{obs}} = \hat{\beta}_j / s_j, \quad \text{with} \quad s_j^2 = \hat{\sigma}^2 (X^T X)^{-1}_{jj}$$

Call the p -value associated with the j -th test p_j . That is, p_j is the probability of seeing a value as extreme as T_j under the null model, i.e., $\beta = 0$ and Gaussian errors. Again, under the null model all these p -values are uniformly distributed.

Now consider the following two step procedure:

Regression after variable selection:

1. Select all variables with $p_j < 0.25$. Call that index set I .
2. Solve a regression problem on X_I , the $n \times |I|$ matrix where we retain only the columns in X corresponding to selected variables.

Since the p -values in the first step are all uniformly distributed, the selection step simply picks a random subset of roughly $d/4$ variables. Any variable is equally likely to be chosen. Starting from $\beta = 0$, of course, there's still no signal in the linear system given by X_I and y . So, we still shouldn't reject the null hypothesis.

What Freedman observed, however, is that the p -value corresponding to an F -test on the final model is sharply biased towards smaller p -values. We're therefore much more likely to reject the null hypothesis ($\beta = 0$) than we should be.

The problem Freedman discovered is now often called *inference after selection*. We first select variables and then we'd like to do *inference*, i.e., compute p -values or confidence intervals, on the selected variables. There's nothing specific about F -tests or T -tests in this example. It's a fundamental problem with two-stage data-dependent analyses.

A safe way to do this is to perform the two steps on independent data sets. So, in a sense, sample splitting solves this two step problem. Statisticians have also come up with a number of clever methods to do the two steps correctly on the same sample.

More broadly, Freedman's paradox relates to a practice that later became known as *p-hacking*. The term describes the practice deliberately choosing an analysis in a data-dependent way so as to find a significant p -value. We'll return to this problem in the next chapter. Freedman's observation foreshadowed problems to come.

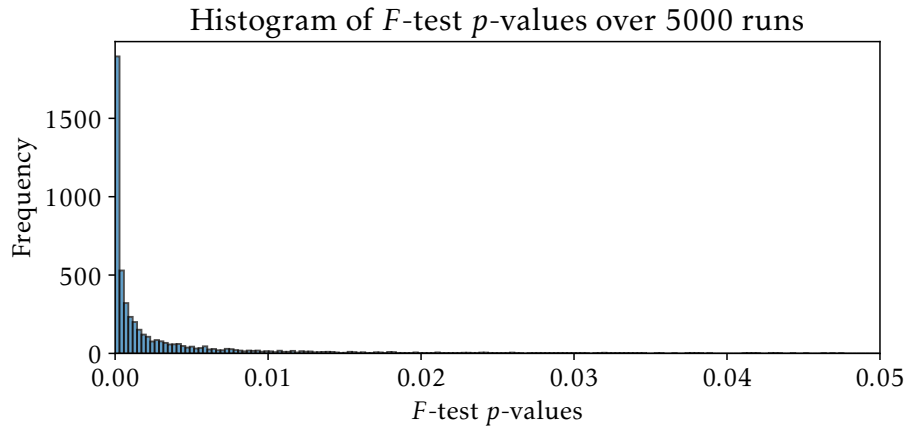


Figure 5.4: Biased p -values in Freedman’s two-stage regression with $n = 1000$ samples and $d = 50$ variables

5.5 Notes

Dwork et al.^{1–3} initiated the study of *adaptive data analysis*. Tools from differential privacy give variants of the holdout method that have stronger guarantees under adaptive use. The idea is to add a sufficient amount of noise to each empirical risk computation on the holdout data. Whereas the plain holdout method supports linear number of queries in the worst case, a holdout method based on noise addition can support a quadratic number of queries. Bassily et al.⁴ improved the error bounds compared with the suboptimal bounds by Dwork et al. Under cryptographic hardness assumption, however, no holdout method can support more than a quadratic number of queries.⁵

Blum and Hardt⁶ study holdout reuse in the context of machine learning benchmarks. The key observation they make is that ranking is different from evaluation: If the goal is to identify the best model from a sequence of model evaluations, a variant of the holdout method supports an exponential number of queries. We will return to this result in a later chapter. In the same paper, they also point out how the basic holdout method fails in the adaptive setting due to overfitting via ensembling. Hardt provided a proof of this proposition in a subsequent paper.⁷ Dwork et al.² gave a similar argument for how adaptively fitting a linear model can lead to the same deviation between risk and empirical risk. Feldman, Frostig, and Hardt

further develop the connection between boosting and overfitting.⁸

Freedman described the problem with p -values after selection in a short 1983 paper.⁹ Freedman neither called it a paradox, nor did he believe that the observation was new. Nevertheless, his note has been quite influential. There's now much work on this problem in statistics under the name *post-selection inference*¹⁰ or *inference after selection*¹¹. Hastie, Tibshirani, and Friedman discuss a variant of Freedman's paradox in the context of cross validation in Chapter 7.10.2 of their book.¹²

While adaptivity brings theoretical challenges, many recognize the need for permitting adaptivity in statistical analysis. In particular, John Tukey and George Box both advocated for allowing incremental progress in data analysis.^{13,14} See also the argument by Gelman and Loken¹⁵ that we'll return to in the next chapter.

Bibliography

- [1] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Roth. The reusable holdout: Preserving validity in adaptive data analysis. *Science*, 349(6248):636–638, 2015.
- [2] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Leon Roth. Preserving statistical validity in adaptive data analysis. In *ACM Symposium on Theory of Computing (STOC)*, pages 117–126, 2015.
- [3] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toni Pitassi, Omer Reingold, and Aaron Roth. Generalization in adaptive data analysis and holdout reuse. *Neural Information Processing Systems (NeurIPS)*, 28, 2015.
- [4] Raef Bassily, Kobbi Nissim, Adam Smith, Thomas Steinke, Uri Stemmer, and Jonathan Ullman. Algorithmic stability for adaptive data analysis. In *ACM Symposium on Theory of Computing (STOC)*, pages 1046–1059, 2016.
- [5] Moritz Hardt and Jonathan Ullman. Preventing false discovery in interactive data analysis is hard. In *Annual Symposium on Foundations of Computer Science (FOCS)*, pages 454–463. IEEE, 2014.
- [6] Avrim Blum and Moritz Hardt. The ladder: A reliable leaderboard for machine learning competitions. In *International Conference on Machine Learning (ICML)*, pages 1006–1014. PMLR, 2015.
- [7] Moritz Hardt. Climbing a shaky ladder: Better adaptive risk estimation. *arXiv preprint arXiv:1706.02733*, 2017.
- [8] Vitaly Feldman, Roy Frostig, and Moritz Hardt. The advantages of multiple classes for reducing overfitting from test set reuse. In *International Conference on Machine Learning (ICML)*, pages 1892–1900. PMLR, 2019.
- [9] David A Freedman and David A Freedman. A note on screening regression equations. *the american statistician*, 37(2):152–155, 1983.
- [10] Jason D Lee, Dennis L Sun, Yuekai Sun, and Jonathan E Taylor. Exact post-

- selection inference, with application to the lasso. *The Annals of Statistics*, 44(3):907, 2016.
- [11] William Fithian, Dennis Sun, and Jonathan Taylor. Optimal inference after model selection. *arXiv preprint arXiv:1410.2597*, 2014.
 - [12] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction (Corrected 12th printing)*. Springer, 2017.
 - [13] John W Tukey. We need both exploratory and confirmatory. *The american statistician*, 34(1):23–25, 1980.
 - [14] George Box. Scientific method: The generation of knowledge and quality. *Quality Progress*, 30(1):47, 1997.
 - [15] Andrew Gelman and Eric Loken. The garden of forking paths: Why multiple comparisons can be a problem, even when there is no “fishing expedition” or “p-hacking” and the research hypothesis was posited ahead of time. *Department of Statistics, Columbia University*, 348(1-17):3, 2013.