

# Algorithmic Statistics Problem Set 1

## Problem 1 - Uniformity Testing

In this exercise you will complete a missing step from the proof of the uniformity testing lower bound from lecture 1.

Using facts from the lecture notes about Poissonization and  $\chi^2$  divergence, finish the proof of Paninski's sample complexity lower bound (poissonized variant) by proving the poissonized version of Lemma 5.2 from the lecture 1 notes. You may look up and use any other facts about Poisson random variables.

## Problem 2: Convergence of Linear Regression

In this problem we analyze one of the most fundamental problems in learning theory: linear regression. Before delving into the analysis it is worth noting a few of the differences between it and many of the other analyses of linear regression.

1. *Agnostic setting.* We do not assume that  $y = \langle x, \beta_{\text{gt}} \rangle + \text{noise}$  for some ground truth  $\beta_{\text{gt}}$ . Instead, we simply define the *best linear predictor*  $\beta^*$  by population least-squares:

$$\beta^* = \arg \min_{\beta} \mathbb{E} [(y - \langle x, \beta \rangle)^2].$$

This makes the results applicable to a wider variety of settings.

2. *Parameter convergence vs. risk convergence.* Here we study convergence of the parameters both in  $\ell_2$  norm (i.e.  $\|\hat{\beta} - \beta^*\|_2 \leq \varepsilon$ ), and in *prediction risk*:

$$\mathbb{E} [(y - \langle x, \hat{\beta} \rangle)^2] \leq \min_{\beta} \mathbb{E} [(y - \langle x, \beta \rangle)^2] + (\text{generalization error}).$$

3. *Fast vs. slow rates.* By exploiting matrix concentration and boundedness assumptions, we obtain sample complexity

$$n = \tilde{O}\left(\frac{d}{\varepsilon^2}\right),$$

sometimes called a *fast rate*. In contrast, “black-box” VC-dimension style theory typically yields only

$$n = \tilde{O}\left(\frac{d}{\varepsilon^4}\right),$$

a *slow rate*.

4. *Bounded Covariates and Residuals* One somewhat unorthodox assumption we make is that the covariates  $x_i$  and the labels  $y_i$  are bounded almost-surely (i.e., with probability 1), where many analyses would only assume something like subgaussian covariates and labels.

We make this assumption to simplify the use of matrix Bernstein (see below), but with a little more work, the same techniques can be used with just a tail bound on these quantities.

**Setting** Let  $(x_i, y_i)_{i=1}^n$  be i.i.d., with  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ , satisfying

- (*isotropy*)  $\mathbb{E}[x] = 0$  and  $\Sigma := \mathbb{E}[xx^\top] = I_d$ ,
- (*bounded norm*)  $\|x\|_2 \leq C\sqrt{d}$  almost surely, for an absolute constant  $C$ ,
- (*bounded labels*)  $|y| \leq 1$  almost surely.

**Definitions** Define the *population* quantities

$$\beta^* := \arg \min_{\beta \in \mathbb{R}^d} \left\{ \mathbb{E} [(y - \langle \beta, x \rangle)^2] \right\} = \mathbb{E} [\Sigma^{-1} xy] =: g, \quad \Sigma = I_d,$$

and the *empirical* quantities

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n x_i x_i^\top, \quad \hat{g} = \frac{1}{n} \sum_{i=1}^n x_i y_i, \quad \hat{\beta} = \hat{\Sigma}^{-1} \hat{g}.$$

**Matrix Bernstein Inequality** You may use the following without proof.

**Lemma 1** (Matrix Bernstein Inequality). *Let  $Z_1, \dots, Z_n$  be independent mean-zero, self-adjoint random matrices in  $\mathbb{R}^{d \times d}$  with  $\|Z_i\| \leq L$  almost surely, and set  $v = \left\| \sum_i \mathbb{E} [Z_i^2] \right\|$ . Then for all  $t \geq 0$ ,*

$$\Pr \left( \left\| \sum_{i=1}^n Z_i \right\| \geq t \right) \leq 2d \exp \left( -\frac{t^2}{2v + \frac{2}{3}Lt} \right).$$

**Goal** Our goal will be to prove a bound on the convergence of linear regression.

**Claim 2.** *Given  $n = \tilde{O}(\frac{d}{\varepsilon^2})$ , with probability  $1 - o(1)$ , the empirical regression solution  $\hat{\beta}$  is close to the population regression solution  $\beta^*$ , up to an  $L_2$  error of  $\|\hat{\beta} - \beta^*\| \leq \varepsilon$ .*

—

## 2A. Concentration of the covariance $\hat{\Sigma}$

Let  $Z_i := x_i x_i^\top - \Sigma = x_i x_i^\top - I_d$ . Note that  $\hat{\Sigma} - I_d = \frac{1}{n} \sum_{i=1}^n Z_i$ .

- (i) **Mean Zero.** Show that  $\mathbb{E} [Z_i] = 0$ .
- (ii) **Operator-norm Bound on  $Z_i$ .** Recall  $\|A\| := \sup_{\|u\|=1} \|Au\|_2$ .  
Prove that  $\|Z_i\| = O(d)$  almost surely.
- (iii) **Variance Property.** Show that

$$Z_i^2 = (x_i x_i^\top - I_d)^2 = (\|x_i\|_2^2 - 2) x_i x_i^\top + I_d,$$

and conclude that

$$\left\| \mathbb{E} [Z_i^2] \right\| = O(d), \quad \left\| \sum_{i=1}^n \mathbb{E} [Z_i^2] \right\| = O(nd).$$

- (iv) **Matrix Bernstein Inequality.** Apply Lemma 1 to  $\sum_{i=1}^n Z_i$  and show that with probability at least  $1 - \delta$ ,

$$\left\| \hat{\Sigma} - I_d \right\| \leq O \left( \sqrt{\frac{d \log \frac{2d}{\delta}}{n}} + \frac{d \log \frac{2d}{\delta}}{n} \right).$$

- (v) **Stability of the inverse.** Show that if  $\left\| \hat{\Sigma} - I_d \right\| \leq \eta < 1$ , then

$$\left\| \hat{\Sigma}^{-1} - I_d \right\| \leq \frac{\eta}{1 - \eta}.$$

—

## 2B. Concentration of the linear term $\hat{g}$

We need a high-probability bound on  $\|\hat{g} - g\|_2$ .

To this end, define  $v_i = x_i y_i - g$  and

$$V_i = \begin{pmatrix} 0 & v_i \\ v_i^\top & 0 \end{pmatrix}.$$

- (i) **Mean Zero.** Show that  $\mathbb{E}[V_i] = 0$ .
- (ii) **Operator-norm Bound.** Show that  $\|V_i\| = O(\sqrt{d})$  almost surely.
- (iii) **Variance Property.** Show that

$$\left\| \mathbb{E}[V_i^2] \right\| = O(d), \quad \left\| \sum_{i=1}^n \mathbb{E}[V_i^2] \right\| = O(nd).$$

- (iv) **Matrix Bernstein.** Conclude that with probability at least  $1 - \delta$

$$\|\hat{g} - g\|_2 = O\left(\sqrt{\frac{d \log(d/\delta)}{n}} + \frac{\sqrt{d} \log(d/\delta)}{n}\right).$$

—

## 2C. Concluding regression convergence

Use the previous parts of the question to prove Claim 2.

- (i) **Key identity.** Prove

$$\hat{\beta} - \beta^* = (\hat{\Sigma}^{-1} - I_d)g + \hat{\Sigma}^{-1}(\hat{g} - g).$$

- (ii) **Bounding each term.** Use the bounds from 2A and 2B.
- (iii) **Final conclusion.** Deduce that for  $n = \tilde{O}(d/\varepsilon^2)$ , with high probability  $\|\hat{\beta} - \beta^*\|_2 \leq \varepsilon$ .

—

## 2D. Extensions: Generalization Errors in General Settings

Explore how the conclusions change under weaker assumptions.

- (i) **Non-isotropic distributions.** Suppose instead  $\Sigma = \text{Cov}(x)$  is not the identity, but satisfies  $x^\top \Sigma^{-1} x = O(d)$  almost surely. Show a similar bound holds on

$$\left\| \Sigma^{1/2}(\hat{\beta} - \beta^*) \right\|.$$

- (ii) **Bounded residuals rather than bounded labels.** Suppose instead of the bounded label assumption, we are told that some “ground-truth”  $\beta_{\text{gt}}$  such that

$$y_i = \langle x_i, \beta_{\text{gt}} \rangle + e_i, \quad e_i \in [-1, 1].$$

Show that the same convergence analysis applies even though the labels are not necessarily bounded (*hint*: analyze the optimization problem  $\gamma^* = \operatorname{argmin}_{\gamma \in \mathbb{R}^d} \left\{ \mathbb{E} \left[ \left( y - \langle x, (\gamma - \beta_{\text{gt}})^2 \rangle \right) \right] \right\}$ ).

- (iii) **Bounding the generalization error.** Suppose  $|y - \langle \beta_{\text{gt}}, x \rangle| \leq \eta$  almost surely. Bound the excess risk

$$\mathbb{E} \left[ \left( y - \langle \hat{\beta}, x \rangle \right)^2 \right] - \mathbb{E} \left[ \left( y - \langle \beta_{\text{gt}}, x \rangle \right)^2 \right].$$

(*hint*: show that  $\mathbb{E}[x(y - \langle \beta_{\text{gt}}, x \rangle)]$ , and then express the excess risk in terms of  $\delta := \beta_{\text{gt}} - \beta$ )

## Problem 3: Learning Halfspaces

In this problem we will build towards the main theorem of Kalai–Klivans–Mansour–Servedio: *halfspaces can be weakly agnostically learned under the uniform distribution*.

### Background

- **Why agnostic learning?** Real-world data may be noisy or not generated exactly by our model class. Agnostic learning relaxes the assumptions: the data distribution is arbitrary, and we only require our learner to perform nearly as well as the best model in hindsight.

In other words, we still assume that our training samples  $(x, y)$  are iid, but we no longer assume that the distribution behaves like a model plus random noise.

- **Why halfspaces?** Halfspaces (linear threshold functions) are fundamental classifiers of the form

$$h(x) = \text{sign}(w \cdot x - \theta), \quad x \in \{\pm 1\}^n.$$

- **Main result (informal).** Kalai–Klivans–Mansour–Servedio (2005) show that halfspaces are agnostically learnable to error  $\text{opt} + \varepsilon$ . In this problem, you will prove a weaker guarantee: there exists an efficient *weak agnostic learner* for halfspaces under the uniform distribution, achieving error

$$O(\text{opt} + \varepsilon).$$

- **Setup.** We are given training data  $\{(x_i, y_i)\}_{i=1}^m$ , with  $x_i \in \{\pm 1\}^n$  features and  $y_i \in \{\pm 1\}$  labels. The key property we will exploit is that every halfspace can be well-approximated by a low-degree Fourier polynomial.

### 3A. Fourier characters on the hypercube

(i) **Definition and orthogonality.** For a subset  $S \subseteq [n]$ , define the Fourier character

$$\chi_S(x) = \prod_{i \in S} x_i.$$

Prove that for subsets  $S, T \subseteq [n]$ :

$$\mathbb{E}_{x \sim \{\pm 1\}^n} [\chi_S(x) \chi_T(x)] = \begin{cases} 1 & S = T, \\ 0 & S \neq T. \end{cases}$$

*Hint:* Show that  $\chi_S \cdot \chi_T = \chi_{S \Delta T}$ , where  $\Delta$  is symmetric difference.

(ii) **Low-degree feature map.** Define the feature map  $\phi : \{\pm 1\}^n \rightarrow \mathbb{R}^N$  by

$$\phi(x) = (\chi_S(x))_{S: 1 \leq |S| \leq d},$$

i.e. all non-empty characters of degree at most  $d$ . Show that under the uniform distribution:

- $\mathbb{E}_{x \sim \{\pm 1\}^n} [\phi(x)] = 0$ ,
- $\mathbb{E}_{x \sim \{\pm 1\}^n} [\phi(x) \phi(x)^\top] = I_N$ .

Thus,  $\phi$  is isotropic (see Problem 2).

### 3B. Learning low-degree Boolean functions with noise

We now show that if the true function is close to a low-degree polynomial, then regression against noisy samples still recovers a good hypothesis.

**Definition 1** (Low-degree Approximability). *Let  $f : \{\pm 1\}^n \rightarrow \{\pm 1\}$  be a Boolean function. We say that  $f$  can be  $\varepsilon$ -approximated by a degree  $d$  polynomial in  $\ell_2$  norm if there exists  $p^*$  in the span of  $\{\chi_S : |S| \leq d\}$  such that*

$$\mathbb{E}_{x \sim \{\pm 1\}^n} [(f(x) - p^*(x))^2] \leq \varepsilon.$$

**Claim 3.** *Let  $f, g : \{\pm 1\}^n \rightarrow \{\pm 1\}$  be Boolean functions, where  $f$  represents the optimal target and  $g$  represents noisy training labels. Suppose:*

1.  $f$  is  $\varepsilon$ -approximated by a degree  $d$  polynomial.
2. (Bounded noise)

$$\Pr_{x \sim \{\pm 1\}^n} [f(x) \neq g(x)] \leq \eta.$$

Let

$$p = \arg \min_{q \in \text{span}\{\chi_S : |S| \leq d\}} \mathbb{E}_{x \sim \{\pm 1\}^n} [(g(x) - q(x))^2].$$

Then

$$\Pr_{x \sim \{\pm 1\}^n} [|p(x) - g(x)| > 1/2] = O(\eta + \varepsilon).$$

(i) **Almost-triangle inequality.** Show that for any reals  $a, b, c$ ,

$$(a - c)^2 \leq 2((a - b)^2 + (b - c)^2).$$

(ii)  $p^*$  has small loss against  $g$ . Use the inequality from part (i) to prove:

$$\mathbb{E}_x [(p^*(x) - g(x))^2] \leq 2\mathbb{E}_x [(f(x) - g(x))^2] + 2\varepsilon.$$

Conclude that

$$\mathbb{E}_x [(p^*(x) - g(x))^2] \leq 8\eta + 2\varepsilon.$$

Why does this imply that the regression minimizer  $p$  must also have loss at most  $8\eta + 2\varepsilon$ ?

(iii) **Conclude the claim.** Explain why the bound on squared error implies that  $p$  misclassifies at most an  $O(\eta + \varepsilon)$  fraction of inputs. *Hint:* If  $|p(x) - g(x)| > 1/2$ , then the squared error is at least  $1/4$ .

### 3C. Agnostically learning halfspaces

We now combine everything.

It is a known fact (you may assume without proof) that every halfspace  $h(x) = \text{sign}(w \cdot x - \theta)$  can be  $\varepsilon$ -approximated in  $L_2$  under the uniform distribution by a polynomial of degree  $\text{poly}(1/\varepsilon)$ . Thus, halfspaces fall into the approximately low-degree setting of Problem 3B.

(i) **Finite-sample regression.** Using Problem 2 and part 3A, show that with  $m = \text{poly}(N, n)$  samples, regression recovers a hypothesis  $\hat{p}$  whose coefficient vector in  $\mathbb{R}^N$  satisfies

$$\|\hat{p} - p\|_2 < \frac{1}{2Nn}.$$

(i.e., bound the  $\ell_2$  norm of the  $\mathbb{R}^N$  vector whose  $S$ th entry is the difference in the coefficient of  $\chi_S$  in  $\hat{p}$  and in  $p$ )

**(ii) From coefficient closeness to classification error.** Show that the  $\ell_2$ -closeness of coefficient vectors implies  $\ell_\infty$  closeness of predictions

$$\max_{x \in \{\pm 1\}^n} \{|p(x) - \hat{p}(x)|\} = \|\hat{p} - p\|_\infty < \frac{1}{2}.$$

*Hint:* Use Cauchy–Schwarz together with the triangle inequality.

Finally, use this and the result of Problem 3B to conclude that  $\hat{p}$  misclassifies at most  $O(\eta + \varepsilon)$  fraction of the domain. Thus, there is a polynomial-time weak agnostic learner for halfspaces under the uniform distribution.