

STAT260: Robust Statistics Course Notes

Feynman Liang*

Department of Statistics, UC Berkeley

Last updated: September 14, 2019

Contents

1	9/3/2019	2
1.1	Minimum distance functional	2
1.2	Midpoint lemma and resilience	4
1.3	Orlicz norms	7
2	9/5/2019	8
2.1	Recap	8
2.2	Concentration Inequalities and Composition	8
2.3	Failure of composition of higher moments and Rosenthal's inequality	10
2.4	Exponential tails and Chernoff bounds	10
2.5	Bounded random variables	12
2.6	Aside: Cumulants are additive	13
2.7	Max of n sub-Gaussians	13
3	9/10/2019	13
3.1	Bounding suprema via concentration	13
3.2	Warmup: max of sub-Gaussian	14
3.3	Maximum eigenvalue of random matrix	14
3.4	VC inequality and Symmetrization	16
4	9/12/2019	19
4.1	Recap	19
4.2	VC dimension of half spaces	20
4.3	Finite sample analysis of MDF via Generalized KS distance	21
	Bibliography	24

*feynman@berkeley.edu

1 9/3/2019

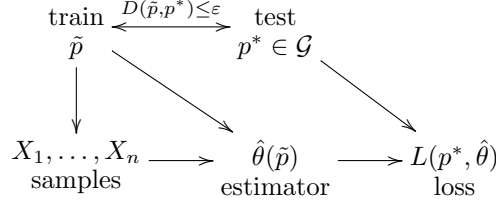


Figure 1: Overview of the framework. Training distribution \tilde{p} differs from test distribution p^* by some discrepancy $D(\tilde{p}, p^*) \leq \epsilon$. We constrain $p^* \in \mathcal{G}$ to encode distributional assumptions. Given an estimator $\hat{\theta}$ trained using samples $X_1, \dots, X_n \sim \tilde{p}$, we want to control the loss $L(p^*, \hat{\theta})$ incurred at test time.

1.1 Minimum distance functional

Introduced in Donoho et al. (1988), the minimum distance functional is one way to produce robust estimators which easily generalizes and also leverages distributional assumptions in \mathcal{G} .

Definition 1 (*Minimum distance functional*)

The *minimum distance functional* (MDF) is

$$\hat{\theta}(\tilde{p}) = \theta^*(q) = \operatorname{argmin}_{\theta} L(q, \theta) \text{ where } q = \operatorname{argmin}_{q \in \mathcal{G}} D(\tilde{p}, q) \quad (1)$$

In other words, q is the projection (under D) of \tilde{p} onto \mathcal{G} and $\hat{\theta}$ is the estimator obtained by using q as the training distribution.

One nice property of the MDF is that we can bound it using a supremum over nearby pairs $p, q \in \mathcal{G}$ satisfying $D(p, q) \leq 2\epsilon$. This is useful because we eliminate \tilde{p} and focus the theory around \mathcal{G} .

Proposition 2 (*Modulus of continuity bound*)

If D is a pseudometric (metric without requirement $d(x, y) = 0 \implies x = y$), then the cost $L(p^*, \hat{\theta}(\tilde{p}))$ of the MDF (Definition 1) is bounded by:

$$\mathfrak{m}(\mathcal{G}, 2\epsilon, D, L) = \sup_{\substack{p, q \in \mathcal{G} \\ D(p, q) \leq 2\epsilon}} L(p, \theta^*(q)) \quad (2)$$

\mathfrak{m} is called the **modulus of continuity**.

Proof. First fix $p = p^* \in \mathcal{G}$

$$\mathfrak{m} \geq \sup_{g \in \mathcal{G}: D(p^*, g) \leq 2\epsilon} L(p^*, \theta^*(g)) \quad (3)$$

Next, let $q = \operatorname{argmin}_{g \in \mathcal{G}} D(g, \tilde{p})$ be the projection of \tilde{p} onto \mathcal{G} as in Definition 1. Then since $D(p^*, \tilde{p}) \leq \epsilon$ by assumption and $p^* \in \mathcal{G}$, we have

$$D(q, \tilde{p}) = \min_{g \in \mathcal{G}} D(g, \tilde{p}) \leq D(p^*, \tilde{p}) \leq \epsilon \quad (4)$$

The following drawing visualizes the argument.

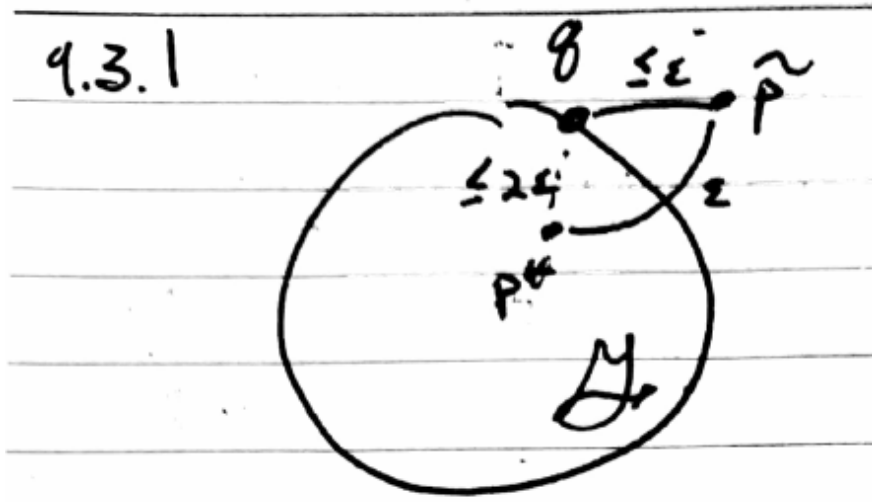


Figure 2: Given $D(p^*, \tilde{p}) \leq \varepsilon$, $p^* \in \mathcal{G}$, and q is the projection of \tilde{p} onto \mathcal{G} under D , we must have $D(\tilde{p}, q) \leq \varepsilon$ and by triangle inequality $D(p^*, q) \leq 2\varepsilon$

So $D(p^*, q) \leq 2\varepsilon$ and we can conclude

$$\mathfrak{m} \geq L(p^*, \theta^*(q)) \quad (5)$$

□

For now, we will specialize to the case $D = \text{TV}$ and $L(p, \theta) = \|\theta - \mu(p^*)\|_2$. Consider a Gaussian distributional assumption $\mathcal{G}_{\text{gauss}} = \{\mathcal{N}(\mu, I) : \mu \in \mathbb{R}^d\}$.

Lemma 3

$$\text{TV}(\mathcal{N}(\mu, I), \mathcal{N}(\mu', I)) \asymp \Theta(\min(\|\mu - \mu'\|_2, 1))$$

Therefore

$$\mathfrak{m}(\mathcal{G}_{\text{gauss}}, \varepsilon) = \sup_{\substack{p, q \in \mathcal{G} \\ \text{TV}(p, q) \leq 2\varepsilon}} \|\mu(p) - \mu(q)\|_2 = \Theta(\varepsilon) \quad (6)$$

for sufficiently small ε .

Proof. We first prove the 1D case. By translational symmetry, we can translate both distributions while preserving $\|\mu - \mu'\|_2 =: u$ so that wlog we may assume the two distributions are $p = \mathcal{N}(\frac{u}{2}, 1)$ and $q = \mathcal{N}(-\frac{u}{2}, 1)$. Then

$$\text{TV}(p, q) = \frac{1}{2} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} |e^{(t+u/2)^2/2} - e^{(t-u/2)^2/2}| dt \quad (7)$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-u/2}^{u/2} e^{-t^2/2} dt \quad (8)$$

where the last equality follows by cancelling the probability mass in the following picture:

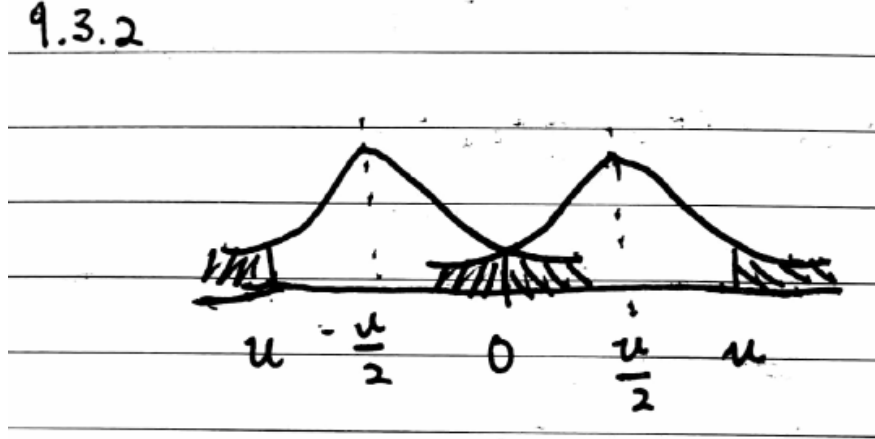


Figure 3: Both Gaussians exhibit identical $\pm \frac{u}{2}$ tails with opposite signs in the expression for TV, so the TV is equivalent to the area in $[-u/2, u/2]$ drawn out by the pointwise max between the two PDFs. By symmetry, this is just twice the area inside $[-u/2, u/2]$ which after cutting and pasting integration areas (and cancelling the $1/2$ in definition of TV) is equal to the probability mass between $[-u/2, u/2]$ for a Gaussian.

Note that $e^{-t^2/2} \geq \frac{1}{2}$ if $|t| < 1$, so $\text{TV} = \Omega(\min(u, 1))$ which can be seen by splitting the integral and examining the two cases where $\frac{u}{2} > 1$ (which yields the 1) and $\frac{u}{2} < 1$ (which yields the u).

Similarly, $e^{-t^2/2} \leq 1$ for all $t > 0$ so $\text{TV} = O(\min(u, 1))$.

To generalize to higher dimensions, note identity covariance implies rotational invariance so we can rotate and translate such that the two means are on the first coordinate axis and separated by $\|\mu - \mu'\| = |\mu_1 - \mu'_1|$. In particular, $\mu_i = 0$ for $i \neq 1$ hence in the TV expression they can be factored out and integrated to 1 to reduce to the 1D case. \square

1.2 Midpoint lemma and resilience

As a less restrictive family, consider distributions with bounded covariance:

$$\mathcal{G}_{\text{cov}}(\sigma) = \{p : \mathbb{E}_p[(X - u)(X - u)'] \preceq \sigma^2 I\} \quad (9)$$

We begin with an important lemma which will be used to prove the modulus of continuity for \mathcal{G}_{cov} and generalized in the following section.

Lemma 4 (Midpoint lemma)

If $\text{TV}(p, q) \leq \varepsilon$ then exists a **midpoint** distribution r such that $r \leq \min\{\frac{p}{1-\varepsilon}, \frac{q}{1-\varepsilon}\}$ and

1. $r(x) \leq \frac{p(x)}{1-\varepsilon}$ for all x
2. r is an ε -**deletion** of p (obtained by deleting ε mass from p)
3. $r = p|_E$ for $p(E) \geq 1 - \varepsilon$ where $E \mid X$ has probability 1 if $p(x) \leq q(x)$ and $\frac{q(x)}{p(x)}$ if $p(x) > q(x)$

Proof. The midpoint distribution is given by $r = \frac{\min(p, q)}{1 - \text{TV}(p, q)}$ and is obtained from p by deleting probability mass from q and renormalizing.

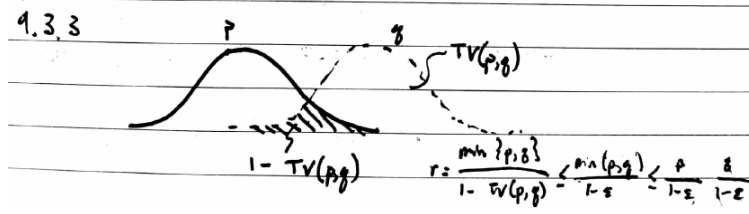


Figure 4: The midpoint distribution $r = \frac{\min(p, q)}{1 - \text{TV}(p, q)}$ can be reached from both p and q by deleting ϵ -mass and renormalizing.

Specifically, we delete $q(x) - p(x)$ mass from all points in $\{x : q(x) > p(x)\}$, the integral of which is precisely equal to the total variation distance. This means that we must renormalize by $1 - \epsilon$ to ensure r is a proper distribution. \square

Corollary 5

$$\mathbf{m}(\mathcal{G}_{cov}(\sigma), \epsilon) = O(\sigma\sqrt{\epsilon})$$

Proof. Take $p, q \in \mathcal{G}_{cov}$ such that $\text{TV}(p, q) \leq \epsilon$. By Lemma 4, there exists a midpoint distribution $r = p|_E$ for which

$$\mathbb{E}_r[X - \mu(p)] = \mathbb{E}_p[X - \mu(p) \mid \underbrace{E}_{1-\epsilon}] = \frac{-\epsilon}{1-\epsilon} \mathbb{E}_p[X - \mu(p) \mid \underbrace{E^c}_{\epsilon}] \quad (10)$$

where the last equality follows from

$$0 = \mathbb{E}_p[X - \mu(p)] = \underbrace{p(E)}_{1-\epsilon} \mathbb{E}_p[X - \mu \mid E] + \underbrace{p(E^c)}_{\epsilon} \mathbb{E}_p[X - \mu \mid E^c] \quad (11)$$

(This is a common trick for moving from conditioning on an event to conditioning on its complement in zero mean functionals).

(Chebyshev in \mathbb{R}^d) By linearity of expectation and Jensen's inequality

$$\|\mathbb{E}_p[X - \mu(p) \mid E^c]\|_2 = \sup_{\|v\|_2 \leq 1} \langle \mathbb{E}_p[X - \mu(p) \mid E^c], v \rangle \quad (12)$$

$$= \sup_{\|v\|_2 \leq 1} \mathbb{E}_p[\langle X - \mu(p), v \rangle \mid E^c] \quad (13)$$

$$\leq \sup_{\|v\|_2 \leq 1} \sqrt{\mathbb{E}_p[\langle X - \mu(p), v \rangle^2 \mid E^c]} \quad (14)$$

Note $\mathbb{E}_p[\langle X - \mu(p), v \rangle^2] = \text{Var}_p[\langle X - \mu(p), v \rangle] = v^\top \text{Cov}_p(X) v \leq \sigma^2$ so

$$\|\mathbb{E}_p[X - \mu(p) \mid E^c]\|_2 \leq \sqrt{\frac{\sigma^2}{\Pr[E^c]}} = \frac{\sigma}{\sqrt{\epsilon}} \quad (15)$$

As a result, we have

$$\|\mu(r) - \mu(p)\|_2 = \|\mathbb{E}_r[X - \mu(p)]\|_2 \leq \frac{\epsilon}{1-\epsilon} \frac{\sigma}{\sqrt{\epsilon}} \leq 2\sigma\sqrt{\epsilon} \quad (16)$$

for $\epsilon < 1/2$. A similar argument involving q gives $\|\mu(r) - \mu(q)\|_2 \leq 2\sigma\sqrt{\epsilon}$ so by triangle inequality $\|\mu(p) - \mu(q)\|_2 \leq 4\sigma\sqrt{\epsilon}$. \square

Remark 6. Unlike the trimmed mean, there is no dependence on d here. This means that the MDF remains a good robust estimator even in high dimensions!

The above proof utilizes two key ingredients:

- The midpoint property of TV; both p and q are close to some ε -deletion r
- The bounded tails (second moment) of \mathcal{G}_{cov} , which is used to control how close $\mu(r)$ and $\mu(p)$ are in Eq. (15)

The previous proof can be suitably generalized to yield a modulus of continuity bound for other families:

Definition 7 (Resilient distribution)

A distribution is (ρ, ε) -resilient if

$$r \leq \frac{p}{1 - \varepsilon} \implies \|\mathbb{E}_r[X] - \mathbb{E}_p[X]\|_2 \leq \rho \quad (17)$$

In other words, for any (not just midpoint) ε -deletion r the mean does not change in norm by more than ρ . Equivalently (e.g. when p does not have a density) we can view $r = p|_E$ for an event E and use the definition

$$p(E) \geq 1 - \varepsilon \implies \|\mathbb{E}_p[X|E] - \mathbb{E}_p[X]\| \leq \rho \quad (18)$$

We let $\mathcal{G}_{TV}(\rho, \varepsilon)$ be the set of all (ρ, ε) -resilient distributions.

Remark 8. This definition is only applicable for mean estimation under squared error loss.

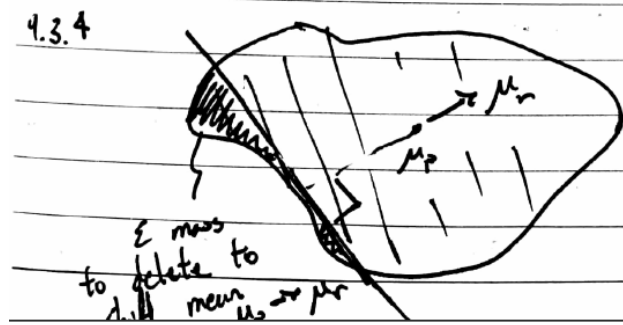


Figure 5: Deleting ε mass from a resilient distribution p shifts the mean by a controlled amount $\|\mu_p - \mu_r\|_2 \leq \rho$.

Example 9

Corollary 5 shows $\mathcal{G}_{cov}(\sigma) \subset \mathcal{G}_{TV}(2\sigma\sqrt{\varepsilon}, \varepsilon)$

Example 10

Lemma 3 shows $\mathcal{G}_{gauss}(\sigma) \subset \mathcal{G}_{TV}(\varepsilon\sqrt{\log \frac{1}{\varepsilon}}, \varepsilon)$

Combining with Proposition 2, for squared error loss we can say

Corollary 11 (Modulus of continuity bound for resilient distributions)

$$m(\mathcal{G}_{TV}(\rho, \varepsilon), \varepsilon) \leq 2\rho \quad (19)$$

Proof. For any $p, q \in \mathcal{G}_{TV}$, use Lemma 4 to get a midpoint distribution and then Eq. (17) with triangle inequality to control the squared error loss. \square

So we can always project onto the family of resilient distributions to get a MDF estimator which has loss independent of d .

1.3 Orlicz norms

Definition 12 (*Orlitz function / norm*)

An **Orlicz function** $\psi : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ is

1. Convex
2. Non-decreasing
3. $\psi(0) = 0$, $\psi(x) \rightarrow \infty$ as $x \rightarrow \infty$

Given an Orlicz function ψ , the **Orlicz norm** or ψ -norm of a random variable X is

$$\|X\|_\psi = \inf \left\{ t : \mathbb{E} \psi \left(\frac{|X|}{t} \right) \leq 1 \right\} \quad (20)$$

For multivariate $X \in \mathbb{R}^d$, define

$$\|X\|_\psi = \inf \left\{ t > 0 : \sup_{v \in \mathcal{S}^{d-1}} \|\langle X, v \rangle\|_\psi \leq t \right\} \quad (21)$$

In other words, X has bounded ψ -norm if all of its one dimensional projections do.

Let $\mathcal{G}_\psi(\sigma) = \{X : \|X - \mathbb{E}[X]\|_\psi \leq \sigma\}$.

Example 13

$\psi(x) = x^k$ gives $\|X\|_\psi = (\mathbb{E}[|X|^k])^{1/k}$, which looks like L_p norms. In fact, these are precisely distributions with bounded k th moments.

For $\psi(x) = x^2$, we have $\mathcal{G}_\psi(\sigma) = \mathcal{G}_{cov}(\sigma)$.

Definition 14 (*Sub-Gaussian/Sub-Exponential*)

For $\psi_2(x) = e^{x^2} - 1$, $\mathcal{G}_{\psi_2}(\sigma)$ are called the σ -sub-Gaussian random variables.

For $\psi_1(x) = e^x - 1$, $\mathcal{G}_{\psi_1}(\sigma)$ are called the σ -sub-exponential random variables.

The next proposition shows that any distribution with bounded Orlicz norm is resilient.

Proposition 15 (*Bounded Orlicz norm implies resilience*)

$$\begin{aligned} \mathcal{G}_\psi(\sigma) &\subset \mathcal{G}_{TV}(2\sigma\epsilon\psi^{-1}(\frac{1}{\epsilon}), \epsilon) \text{ if } \epsilon < 1/2. \\ \psi(x) \rightarrow \psi^{-1}(x) &= \sqrt{x} \rightarrow \epsilon\psi^{-1}(1/\epsilon) = \sqrt{\epsilon} \end{aligned}$$

Proof.

$$\|\mathbb{E}_r[X] - \mathbb{E}_p[X]\|_2 = \|\mathbb{E}_p[X - \mu \mid \underbrace{E}_{p(E)=1-\epsilon}]\|_2 = \frac{\epsilon}{1-\epsilon} \|\mathbb{E}_p[X - \mu \mid E^c]\| \quad (22)$$

Focusing in on the expectation term

$$\|\mathbb{E}_p[X - \mu \mid E^c]\|_2 = \sup_{\|v\|_2=1} \mathbb{E}_p[\langle X - \mu, v \rangle \mid E^c] \quad (23)$$

By Jensen's inequality, convexity of ψ (equivalently concavity of ψ^{-1}), definition of multivariate Orlicz norm

(Eq. (21)), and monotonicity of ψ , we have

$$\|\mathbb{E}_p[X - \mu \mid E^c]\|_2 = \sup_{\|v\|_2=1} \sigma \left(\mathbb{E}_p \left[(\sigma\psi^{-1} \circ \psi) \left(\frac{|\langle X - \mu, v \rangle|}{\sigma} \right) \mid E^c \right] \right) \quad (24)$$

$$\leq \sup_{\|v\|_2=1} \sigma\psi^{-1} \left(\mathbb{E}_p \left[\psi \left(\frac{|\langle X - \mu, v \rangle|}{\sigma} \right) \mid E^c \right] \right) \quad (25)$$

$$\leq \sup_{\|v\|_2=1} \sigma\psi^{-1} \left(\underbrace{\mathbb{E}_p \left[\psi \left(\frac{|\langle X - \mu, v \rangle|}{\sigma} \right) \right]}_{\leq 1} \underbrace{\frac{1}{\Pr[E^c]}}_{\frac{1}{\varepsilon}} \right) \quad (26)$$

$$\leq \sigma\psi^{-1} \left(\frac{1}{\varepsilon} \right) \quad (27)$$

□

2 9/5/2019

2.1 Recap

- Minimum distance functionals: good error, bounded by modulus of continuity \mathfrak{m}
- Resilience \implies bounded \mathfrak{m}
- Bounded Orlicz ψ -norm \implies resilience

This lecture:

- Want to analyze X_1, \dots, X_n
- “The empirical average converges to the mean if n is large”
- Two steps:
 1. Show **concentration inequality**: bound variation of p in terms of σ
 2. Show **composition property**: σ gets smaller as we take more independent samples

2.2 Concentration Inequalities and Composition

Example 16

A slot machine has expected payout of \$5 and always pays out positive.

Question: What is the maximum probability of $\geq \$100$?

Answer: 5%, by letting $P(X = \$0) = 0.95$ and letting $P(X = \$100) = 5\%$.

The preceding example is an instance of Markov's Inequality:

Theorem 17 (*Markov's Inequality*)

If $X \geq 0$ has bounded first moment, then

$$\Pr[X \geq t\mathbb{E}[X]] \leq \frac{1}{t} \quad (28)$$

Proof.

$$t\mathbb{E}[X] \mathbb{1}\{X \geq t\mathbb{E}[X]\} \leq X \quad (29)$$

Take expectation of both sides and rearrange. □

Markov's Inequality has a nice composition property:

Theorem 18 (Composition of Markov for sums)

Let $X_1, X_2 \sim p$ with mean μ .

$$\Pr \left[\frac{X_1 + X_2}{2} \geq t\mu \right] \leq \frac{1}{t} \quad (30)$$

This is because $\mathbb{E}[(X_1 + X_2)/2] = \mu = \mathbb{E}[X_1] = \mathbb{E}[X_2]$.

We can apply Markov's Inequality to $Z = f(X)$ for $f : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$ and get a family of inequalities (provided $\mathbb{E}[f(X)] < \infty$). For example, taking $Z = (X - \mu)^2$ and assuming $\mathbb{E}[Z] = \text{Var}[X] = \sigma^2 < \infty$ yields

Theorem 19 (Chebyshev's inequality)

$$\Pr[|X - \mu| \geq t\sigma] \leq \frac{1}{t^2} \quad (31)$$

Analogous to Theorem 18 (Composition of Markov for sums), a composition property for Chebyshev's inequality would require a composition property involving variances:

Theorem 20 (Variances add for independent RVs)

If $\{X_i\}_{i=1}^n$ are independent, then

$$\text{Var} \left[\sum_{i=1}^n X_i \right] = \sum_{i=1}^n \text{Var}[X_i] \quad (32)$$

Example 21 (Concentration of empirical average)

Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p$ with mean μ and variance σ^2 . Let $S = \sum_i^n X_i$ and $\frac{S}{n}$ the empirical average. Then

$$\text{Var}[S/n] = n \text{Var}[X/n] = n \frac{\sigma^2}{n^2} = \frac{\sigma^2}{n} \quad (33)$$

Hence, the standard deviation of the empirical average $\frac{S}{n}$ is $\sigma_{avg} = \frac{\sigma}{\sqrt{n}}$. Chebyshev's inequality then yields

Corollary 22

$$\Pr \left[\left| \frac{S}{n} - \mu \right| \geq t \frac{\sigma}{\sqrt{n}} \right] \leq \frac{1}{t^2} \quad (34)$$

The t^{-2} quadratic decay in Corollary 22 is tight, as the following proposition shows:

Proposition 23 (Lower bounds for Chebyshev)

There exists X_1, \dots, X_n pairwise independent, bounded in $[-1, 1]$, mean zero, variance one, such that

$$\Pr \left[\sum_{i=1}^n X_i = n \right] = \frac{1}{n} \quad (35)$$

Consequently, Corollary 22 (with $\mu = 0$, $\sigma = 1$, and $t = \sqrt{n}$) is tight.

Proof. Flip k independent coins and let $n = 2^k$. For any subset $\emptyset \subsetneq S \subset [k]$, define the random variable

$$X_S = \begin{cases} 1 & \# \text{ heads in } S \text{ is even} \\ -1 & \# \text{ heads in } S \text{ is odd} \end{cases} \quad (36)$$

X_S is mean zero, variance one, bounded $[-1, 1]$, and pairwise independent (since the coin flips are). The event $\{\sum_{i=1}^n X_i = n\}$ occurs iff all of the coins land tails, which occurs with probability $2^{-k} = \frac{1}{n}$. \square

2.3 Failure of composition of higher moments and Rosenthal's inequality

To try to extend Chebyshev's inequality, we can consider applying Markov's Inequality to $Z = f(X) = (X - \mu)^4$ to get:

Theorem 24

$$\Pr[|X - \mu| \geq t\mathbb{E}[Z]^{1/4}] \leq \frac{1}{t^4} \quad (37)$$

However, the composition property fails here since for $\mathbb{E}[X_1] = \mathbb{E}[X_2] = 0$ we find

$$\mathbb{E}[(X_1 + X_2)^4] = \mathbb{E}[X_1^4] + \mathbb{E}[X_2^4] + \cancel{4\mathbb{E}[X_1^3]\mathbb{E}[X_2]}^0 + \cancel{4\mathbb{E}[X_2^3]\mathbb{E}[X_1]}^0 + \underbrace{6\mathbb{E}[X_1^2 X_2^2]}_{\geq 0} \quad (38)$$

Thus, the fourth moment of a sum can be larger than the sum of the fourth moments.

In general, higher moments don't add. One method to work around this is to work with cumulants (see Section 2.6). An alternative method is through Rosenthal's inequality:

Lemma 25 (*Rosenthal's inequality*)

If X_1, \dots, X_n are independent mean zero random variables with finite p th moments, then

$$\mathbb{E}\left[\left|\sum_{i=1}^n X_i\right|^p\right] \leq O(p)^p \sum_{i=1}^n \mathbb{E}[|X_i|^p] + O(\sqrt{p})^p \left(\sum_{i=1}^n \mathbb{E}[X_i^2]\right)^{p/2} \quad (39)$$

How can we use Rosenthal's inequality? Suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \pi$ with $\mathbb{E}[|X|^p] = k^p$ and $\mathbb{E}[X^2] = \sigma^2$. Let $S = \sum_{i=1}^n X_i$. Then

$$\mathbb{E}[|S|^p] \leq O(p)^p n k^p + O(\sqrt{p})^p (n \sigma^2)^{p/2} \quad (40)$$

$$\mathbb{E}[|S|^p]^{1/p} \leq O(p k n^{1/p} + \sqrt{p} \sigma n^{1/2}) \quad (41)$$

$$\mathbb{E}\left[\left|\frac{S}{n}\right|^p\right]^{1/p} \leq O(p k n^{-(1-\frac{1}{p})} + \sqrt{p} \sigma n^{-1/2}) \quad (42)$$

Hence, all of the p th moments of the empirical mean $\frac{S}{n}$ decay in n , so the empirical mean concentrates about the population mean as the number of samples $n \rightarrow \infty$.

2.4 Exponential tails and Chernoff bounds

Another approach which can yield exponential tail bounds is through the Moment Generating Function.

Definition 26 (*Moment Generating Function*)

Let X be a random variable with bounded moments. The *moment generating function* (MGF)

of X is

$$m_X(\lambda) = \mathbb{E} \exp(\lambda(X - \mu)) = 1 + \lambda \mathbb{E}[X] + \frac{\lambda^2}{2} \mathbb{E}[X^2] + \frac{\lambda^3}{6} \mathbb{E}[X^3] + \dots \quad (43)$$

MGFs satisfy a desirable composition property enabling us to easily compute the MGF of a sum in terms of the MGFs of the summands:

Lemma 27 (Composition property for MGFs)

If X_1, \dots, X_n are independent, then

$$m_{\sum_{i=1}^n X_i}(\lambda) = \prod_{i=1}^n m_{X_i}(\lambda) \quad (44)$$

Proof. Exponential of sum is product of exponentials, independence of X_i allows splitting of \mathbb{E} . □

Another strong advantage of working with moment generating functions is that we can use them to get exponentially decaying tail bounds:

Theorem 28 (Chernoff's bound)

For $\lambda \geq 0$,

$$\Pr[X - \mu \geq t] \leq \inf_{\lambda \geq 0} m_X(\lambda) e^{-\lambda t} \quad (45)$$

Proof. $X - \mu \geq t$ implies $\exp(\lambda(X - \mu)) \geq e^{\lambda t}$. The same technique used to prove Chebyshev's inequality (with $f(x) = e^{\lambda x}$) gives

$$\Pr[\exp(\lambda(X - \mu)) \geq e^{\lambda t}] \leq e^{-\lambda t} m_X(\lambda) \quad (46)$$

□

Example 29 (sub-exponential Chernoff bound)

Recall from Definition 14 (Sub-Gaussian/Sub-Exponential) that σ -sub-exponential means bounded Orlicz norm $\|X - \mu\|_\psi = \mathbb{E} \left[\psi \left(\frac{|X - \mu|}{\sigma} \right) \right] \leq 1$ for $\psi(x) = e^x - 1$. Chernoff's bound then implies

$$\mathbb{E}[\exp(|X - \mu|/\sigma) - 1] \leq 1 \quad (47)$$

$$\mathbb{E}[\exp(|X - \mu|/\sigma)] \leq 2 \quad (48)$$

$$m_X(1/\sigma) = \mathbb{E} \exp \left(\frac{x - \mu}{\sigma} \right) \leq \mathbb{E} \exp \left(\frac{|x - \mu|}{\sigma} \right) \leq 2 \quad (49)$$

$$\Pr[X - \mu \geq t] \leq 2 \exp(-t/\sigma) \quad (50)$$

This explains the name “sub-exponential”: the tail probabilities decay faster than an exponential.

Example 30 (sub-Gaussian Chernoff bound)

Recall from Definition 14 (Sub-Gaussian/Sub-Exponential) that σ -sub-Gaussian means bounded Orlicz norm $\|X - \mu\|_\psi$ with $\psi(x) = e^{x^2} - 1$. Hence, $\mathbb{E}[\exp((X - \mu)^2/\sigma^2)] \leq 2$ and

$$m_X(\lambda) = \mathbb{E} \exp(\lambda(X - \mu)) \leq \exp(\lambda^2 \sigma^2 / 4) \mathbb{E} \exp((X - \mu)^2 / \sigma^2) \leq 2 \exp(\lambda^2 \sigma^2 / 4) \quad (51)$$

where we have used inequality $ab \leq \frac{a^2}{4} + b^2$ to conclude

$$\lambda(X - \mu) \leq \frac{\lambda^2 \sigma^2}{4} + \frac{(X - \mu)^2}{\sigma^2} \quad (52)$$

Remark 31. We can also show

$$m_X(\lambda) \leq \exp\left(\frac{1}{2}\lambda^2(\sigma')^2\right) \quad (53)$$

where $\sigma' \leq \sqrt{3}\sigma$. This is sometimes taken as definition of σ' -sub-Gaussian.

Applying Chernoff's bound shows

$$\Pr[X - \mu \geq t] \leq \inf_{\lambda \geq 0} m_X(\lambda)e^{-\lambda t} \quad (54)$$

$$\leq \inf_{\lambda \geq 0} \exp\left(\frac{1}{2}\lambda^2(\sigma')^2 - \lambda t\right) \quad (55)$$

$$= \exp\left(-\frac{t^2}{2(\sigma')^2}\right) \quad (56)$$

This explains the name “ σ' -sub-Gaussian”: the tail probabilities are decaying faster than those of a Gaussian distribution with variance σ' .

By Lemma 27 (Composition property for MGFs), we have that the sum $S = \sum_i^n X_i$ of σ' -sub-Gaussian RVs is itself $\frac{\sigma'}{\sqrt{n}}$ -sub-Gaussian and satisfies the tail bound

$$\Pr\left[\frac{S}{n} - \mu \geq t\right] \leq \exp\left(-\frac{nt^2}{2(\sigma')^2}\right) = \exp\left(-\frac{nt^2}{6\sigma^2}\right) \quad (57)$$

This yields our desired exponential rate of concentration.

2.5 Bounded random variables

Bounded RVs are sub-Gaussian, but we can get tighter bounds than the previous example. Let $X - \mu \in [-M, M]$. Then

$$\mathbb{E} \exp \frac{|X - \mu|}{M^2 / \log 2} \leq \mathbb{E} \exp \log 2 = 2 \quad (58)$$

Hence X is sub-Gaussian with parameter $\sigma = \sqrt{\frac{M^2}{\log 2}}$ and we can use Eq. (66) to get tail bounds. More generally:

Corollary 32 (*Hoeffding's inequality*)

Let $X_1, \dots, X_n \in [a, b]$ be bounded independent mean zero random variables. Then

$$\Pr\left[\frac{S}{n} - \mu \geq t\right] \leq \exp\left(-\frac{2nt^2}{(a - b)^2}\right) \quad (59)$$

Proof. Bound MGF (tighter than what we are doing here) and apply Chernoff's bound. \square

Hoeffding's inequality shows that an empirical average of independent bounded random variables converges to its mean at a rate of $\frac{1}{\sqrt{n}}$ with tails that decay at least as fast as Gaussians. Compare this against the $\frac{1}{n}$ rate for sub-exponentials we found in Example 29 and the quadratic $\frac{1}{t^2}$ tails from Chebyshev's inequality (which only required finite second moments).

2.6 Aside: Cumulants are additive

We saw in Section 2.3 that fourth moments are additive. While Lemma 27 (Composition property for MGFs) provides a convenient composition property for moment generating functions, the existence of MGFs requires all moments of the random variable to be bounded. In particular, this excludes random variables with fat tails.

To construct additive quantities, we can start with MGF (multiplicative) and take log (which is additive)

$$K_X(\lambda) = \log \mathbb{E} \exp(\lambda(X - \mu)) \quad (60)$$

$$= \log \left(1 + \mathbb{E}[(X - \mu)^2] \frac{\lambda^2}{2} + \mathbb{E}[(X - \mu)^3] \frac{\lambda^3}{6} + \dots \right) \quad (61)$$

$$= 1 + \sum_{n=1}^{\infty} \frac{\kappa_n(X)}{n!} \lambda^n \quad (62)$$

This leads to the cumulant generating function:

Definition 33 (*Cumulants*)

The *cumulant generating function* for a random variable X is

$$K_X(\lambda) = \log \mathbb{E} \exp(\lambda X) = 1 + \sum_{n=1}^{\infty} \frac{\kappa_n(X)}{n!} \lambda^n \quad (63)$$

$\kappa_n(X)$ is called the n th cumulant of X .

Notice $K_{X+Y}(\lambda) = K_X(\lambda) + K_Y(\lambda)$ so we have additivity of the CGF and consequentially

$$\kappa_4(X + Y) = \kappa_4(X) + \kappa_4(Y) \quad (64)$$

Contrast this to Eq. (38).

However, computing the cumulants require Taylor expanding log using the infinite series in Eq. (61) as the argument and are laborious to work with. To handle heavy tails, it may be easier to use Rosenthal's inequality instead.

2.7 Max of n sub-Gaussians

Let $X_1, \dots, X_n \sim p$, p is σ -sub-Gaussian. A simple union bound shows:

Theorem 34 (*Max of sub-gaussian bound*)

$$\Pr[X_1 \vee \dots \vee X_n \geq t] \leq \sum_{i=1}^n \Pr[X_i \geq t] \leq ne^{-\frac{t^2}{2\sigma^2}} \quad (65)$$

So in particular if $t \gg \sigma\sqrt{\log n}$, then its not likely the max will exceed t .

3 9/10/2019

3.1 Bounding suprema via concentration

The typical type of quantity we will focus on here is

$$\underbrace{\sup_{v \in V}}_{\text{bound by discretization}} \underbrace{\frac{1}{n} \sum_{i=1}^n (X_i(v) - \mathbb{E}[X(v)])}_{\text{bound for fixed } v \text{ via concentration}} \quad (66)$$

When V is finite, a simple union bound can be applied. To deal with infinitely large $|V|$, we will need to first discretize V into a finite set.

3.2 Warmup: max of sub-Gaussian

Suppose $\{X_i\}_{i=1}^n \stackrel{\text{iid}}{\sim} p$ where p is mean zero and σ -sub-Gaussian. How big is $\max_{i=1}^n X_i$?

Lemma 35

With probability $\geq 1 - \delta$

$$\max_{i=1}^n X_i \in O\left(\sigma\sqrt{\log n + \log \frac{1}{\delta}}\right) \quad (67)$$

Proof. By union bound, iid, and sub-Gaussian Chernoff bound

$$\Pr\left[\max_{i=1}^n X_i \geq t\right] \leq n \Pr[X_1 \geq t] \leq n \exp\left(-\frac{t^2}{2\sigma^2}\right) \quad (68)$$

To ensure this failure event occurs with probability $\leq \delta$, we need

$$n \exp\left(-\frac{t^2}{2\sigma^2}\right) \leq \delta \quad (69)$$

$$\frac{t^2}{2\sigma^2} = \log n + \log \frac{1}{\delta} \quad (70)$$

$$t \leq \sigma\sqrt{2\left(\log n + \log \frac{1}{\delta}\right)} \quad (71)$$

□

If instead we were interested in $\max_{i=1}^n |X_i|$, then a union bound on the two tail events $\{X_i \geq t\}$ and $\{-X_i \geq t\}$ (note $-X_i$ is still sub-Gaussian) gives

$$\Pr\left[\max_{i=1}^n |X_i| \geq t\right] \leq 2n \exp\left(-\frac{t^2}{2\sigma^2}\right) \quad (72)$$

$$\max_{i=1}^n |X_i| \in O\left(\sigma\sqrt{2\left(\log 2 + \log n + \log \frac{1}{\delta}\right)}\right) \quad (73)$$

In later the next section, we will see how we can “reduce” an infinitely large V into an exponentially large N after which we will use the same technique to bound the event $\{\max_{i \in N} X_i \geq t\}$. To get concentration, we will need the exponential tail bound to dominate the now exponentially large $n = |N|$ arising from the union bound over N .

3.3 Maximum eigenvalue of random matrix

Suppose $\{X_i \in \mathbb{R}^d\}_{i=1}^n \stackrel{\text{iid}}{\sim} p$ with p zero mean and σ -sub-Gaussian.

Recall from Eq. (21) (Orlitz function / norm) and Definition 14 (Sub-Gaussian/Sub-Exponential) that $X \in \mathbb{R}^d$ is σ -sub-Gaussian if all its one dimensional projections are, that is:

$$\|X\|_\psi + 1 = \sup_{v \in \mathcal{S}^{d-1}} \|\langle v, X \rangle\|_\psi + 1 = \sup_{v \in \mathcal{S}^{d-1}} \mathbb{E} \exp\left(\frac{\langle v, X \rangle^2}{\sigma^2}\right) \leq 2 \quad (74)$$

We are interested in the (random) empirical covariance matrix

$$M = \frac{1}{n} \sum_{i=1}^n X_i X_i^\top \quad (75)$$

Specifically, we would like to understand how big $\|M\| = \lambda_{\max}(M)$ is.

Proposition 36

With probability $\geq 1 - \delta$

$$\|M\| = O\left(\sigma^2 \left(1 + \frac{d}{n} + \frac{\log \frac{1}{\delta}}{n}\right)\right) \quad (76)$$

Remark 37. Proposition 36 shows that:

- As $n \rightarrow \infty$, $\|M\| = O(\sigma^2)$ and does not depend on d .
- The population covariance operator norm $\|\mathbb{E}X_iX_i^\top\| = O\left(\frac{\sigma^2}{n} \log \frac{1}{\delta}\right)$ is attained if $d = \Theta(n)$ (i.e. if the dimension grows at the same rate as n)

To relate back to the two-step strategy outlined in Eq. (66), note

$$\|M\| = \sup_{v \in \mathcal{S}^{d-1}} v^\top M v = \sup_{v \in \mathcal{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n \langle X_i, v \rangle^2 \quad (77)$$

This quantity looks promising as it is the sum of independent sub-Gaussian RVs.

Since $\langle X_i, v \rangle$ is σ -sub-Gaussian, $\langle X_i, v \rangle^2$ is σ^2 -sub-exponential (Definition 14, or equivalently Eq. (74)) and for any fixed $v \in \mathcal{S}^{n-1}$

$$\mathbb{E} \exp\left(\frac{\langle X_i, v \rangle^2}{\sigma^2}\right) \leq 2 \quad (78)$$

$$\mathbb{E} \exp\left(\frac{n}{\sigma^2} \frac{1}{n} \sum_{i=1}^n \langle X_i, v \rangle^2\right) = \prod_{i=1}^n \mathbb{E} \exp\left(\frac{\langle X_i, v \rangle^2}{\sigma^2}\right) \leq 2^n \quad (79)$$

where we used Composition property for MGFs for the equality in the second line.

By Theorem 28, for fixed $v \in \mathcal{S}^{d-1}$

$$\Pr[v^\top M v \geq t] \leq 2^n \exp\left(\frac{-nt}{2\sigma^2}\right) \quad (80)$$

So we have accomplished the first step (showing the individual terms inside the sup concentrate for fixed v).

For the second step, we will take a sufficiently small discretization of the unit ball $\{\|v\| \leq 1\}$:

Lemma 38

There exists a finite set $N \subset \mathbb{R}^d$ with $|N| \leq 9^d$ and

$$\sup_{v \in \mathcal{S}^{d-1}} v^\top M v \leq 2 \sup_{v \in N} v^\top M v \quad (81)$$

Applying Lemma 38, a union bound, Eq. (80), and bounding the failure probability by δ shows that

$$\Pr[\|M\| \geq t] = \Pr\left[\sup_{v \in \mathcal{S}^{d-1}} v^\top M v \geq t\right] \leq 9^d 2^n \exp\left(\frac{-nt}{2\sigma^2}\right) = \delta \quad (82)$$

$$\frac{nt}{2\sigma^2} = d \log 9 + n \log 2 + \log \frac{1}{\delta} \quad (83)$$

$$t = O\left(\sigma^2 \left(\frac{d}{n} + 1 + \frac{\log 1/\delta}{n}\right)\right) \quad (84)$$

Proof of Lemma 38. Let N be a maximal packing of $\text{Ball}_1(0)$ in \mathbb{R}^d such that $\|u - v\|_2 \geq \frac{1}{4}$ for all $u \neq v \in N$.

As shown in Fig. 6, if we place a $1/8$ -radius ball at all the points in N then (1) all the balls are disjoint and (2) the union of all the balls is contained in $\text{Ball}_{9/8}(0)$. Therefore, by the (converse of the) pigeonhole principle, $|N| \leq \frac{\text{Vol}(\text{Ball}_{9/8}(0))}{\text{Vol}(\text{Ball}_{1/8}(0))} = 9^d$.

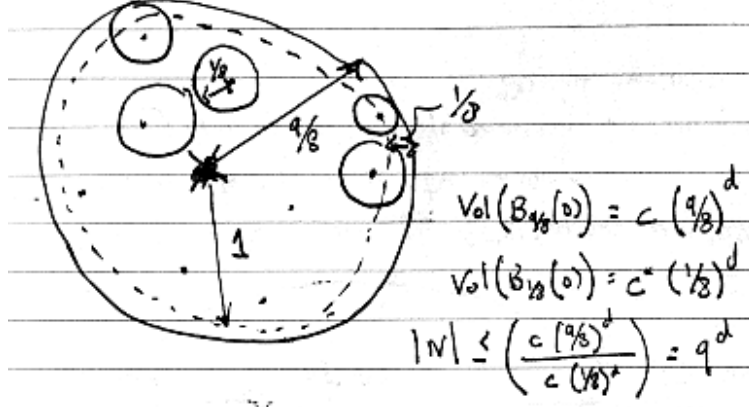


Figure 6: $1/8$ -radius balls centered at all packing points are disjoint, the union of all these balls is contained in $B_{9/8}(0)$, so the cardinality $|N| \leq \left(\frac{9/8}{1/8}\right)^d = 9^d$.

Let $v \in \mathcal{S}^{d-1}$ maximize $v^\top M v$ and $u \in N$ such that $\|u - v\|_2 \leq \frac{1}{4}$. Such a u must exist, otherwise $N \cup \{v\}$ is a larger $1/4$ -packing which contradicts maximality of N .

$$\|M\| - |u^\top M u| = |v^\top M v| - |u^\top M u| \quad (85)$$

$$\leq |v^\top M v - u^\top M u| \quad (86)$$

$$= |(u + v)^\top M (u - v)| \quad (87)$$

$$\leq \underbrace{\|u + v\|_2}_{\leq 2} \|M\| \underbrace{\|u - v\|_2}_{\leq 1/4} \quad (88)$$

$$\leq \frac{1}{2} \|M\| \quad (89)$$

Hence $\|M\| \leq 2u^\top M u \leq 2 \sup_{u \in N} u^\top M u$ as desired. \square

3.4 VC inequality and Symmetrization

In this section, we will see how a family of events with certain geometric structure (which we will quantify using VC-dimension) converges to its expectation at a rate dependent on the geometry. In the process, we will encounter the technique of **symmetrization** (Prof. Steinhardt calls it “bring your own randomness”) used to add additional randomness which will be required to get concentration.

Let \mathcal{H} be a collection of functions $f : \mathcal{X} \rightarrow \{0, 1\}$ and $\{X_i \in \mathcal{X}\}_{i=1}^n \stackrel{\text{iid}}{\sim} p$. For $f \in \mathcal{H}$, let

$$\nu(f) = \mathbb{E}_{x \sim p}[f(x)] = \Pr_{x \sim p}[f(X) = 1] \quad (90)$$

$$\nu_n(f) = \frac{1}{n} \sum_{i=1}^n f(X_i) = \frac{1}{n} \#\{i : f(X_i) = 1\} \quad (91)$$

be the population and empirical averages respectively.

Question: How big is the discrepancy

$$D_n = \sup_{f \in \mathcal{H}} |\nu_n(f) - \nu(f)| \quad (92)$$

Easy case: $|\mathcal{H}| < \infty$. Since $f(X_i)$ is bounded, apply Hoeffding's inequality to the sum of independent bounded random variables to get:

$$D_n = \max_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E}f(X)) \quad (93)$$

$$\Pr \left[\frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E}f(X)) \geq t \right] \leq \exp(-2nt^2) \quad (94)$$

A subsequent union bound over $|\mathcal{H}|$ reveals $t = O\left(\sqrt{\frac{1}{2n} (\log|\mathcal{H}| + \log \frac{1}{\delta})}\right)$

More common case: $|\mathcal{H}| = \infty$. In this situation, we will bound D_n using the geometry of \mathcal{H} . To do so, we will quantify the geometry using the following definitions:

Definition 39 (Shattering number / VC dimension)

The *shattering number* of \mathcal{H} is

$$V_{\mathcal{H}}(\{x_i\}_{i=1}^n) = \# \text{ distinct}\{(f(x_1), \dots, f(x_n)) : f \in \mathcal{H}\} \quad (95)$$

$$V_{\mathcal{H}}(n) = \max_{|S|=n} V_{\mathcal{H}}(S) \quad (96)$$

It measures the number of possible ways to assign $\{0, 1\}$ labels to x_i which can be perfectly fit by $f \in \mathcal{H}$.

The **VC dimension**

$$vc(\mathcal{H}) = \max\{n : V_{\mathcal{H}}(n) = 2^n\} \quad (97)$$

It measures the largest cardinality n such that for any set of points S with cardinality $|S| = n$ and any $\{0, 1\}$ labelling of those points, some $f \in \mathcal{H}$ can perfectly fit it.

The shattering number is useful because instead of taking $\sup_{f \in \mathcal{H}}$ of a term involving f only through $\{f(X_i)\}_{i=1}^n$, we can instead take the supremum over $\{f(X_i)\}_{i=1}^n$ directly and only deal with $V_{\mathcal{H}}(n)$ terms.

Example 40 (VC dimension of half spaces)

Let $\mathcal{X} = \mathbb{R}^d$, $\mathcal{H} = \text{half spaces} = \{f(x) = \mathbb{1}[\langle v, x \rangle \geq \tau] : v \in \mathbb{R}^d, \tau \in \mathbb{R}\}$. Then $vc(\mathcal{H}) = d + 1$.

We will see a full proof later in Proposition 43, but for now consider an example where $d = 2$. We can always separate 3 points by drawing a line, so $vc(\mathcal{H}) \geq 3$. However, with 4 points there can be crossings (see Example 40) which cannot be shattered.

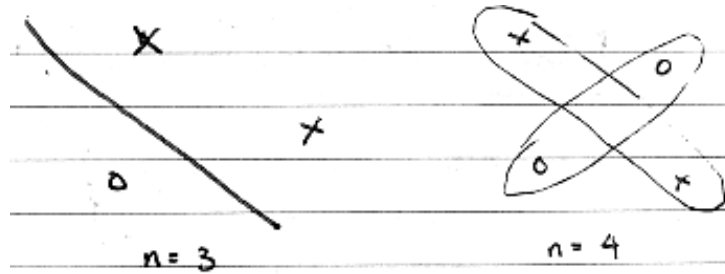


Figure 7: $n = 3$ can always be shattered by a line, but the crossings possible when $n = 4$ prevent this.

Clearly by definition $V_{\mathcal{H}}(n) = 2^n$ for all $n \leq vc(\mathcal{H})$. When $n > vc(\mathcal{H})$, by Eq. (97) (Shattering number / VC dimension) we have $V_{\mathcal{H}}(n) < 2^n$. The following lemma quantifies this and shows that the shattering number is actually significantly smaller (growing polynomially in n rather than exponentially):

Lemma 41 (Sauer-Shelah)

If $vc(\mathcal{H}) = d$, then $V_{\mathcal{H}}(n) \leq 2n^d$.

While we will use this without proof, Sauer-Shelah is the main reason why VC dimension is useful for us: it allows us to convert the infinite supremum over $f \in \mathcal{H}$ into a finite supremum over $O(n^{c(\mathcal{H})})$ many terms of the form $\{f(X_i)\}_{i=1}^d$.

Theorem 42 (VC inequality)

With probability $\geq 1 - \delta$

$$D_n = O\left(\sqrt{\frac{vc(\mathcal{H}) + \log \frac{1}{\delta}}{n}}\right) \quad (98)$$

Proof. We will show something weaker, namely:

$$\mathbb{E}D_n \leq O\left(\frac{vc(\mathcal{H}) \log n}{n}\right) \quad (99)$$

The $\log \frac{1}{\delta}$ tail bound follows from McDiarmid's inequality, and removing the extra $\log n$ refines the argument we will give using a tool called chaining.

Incorrect proof path: Notice that

$$D_n = \sup_{f \in \mathcal{H}} \left| \underbrace{\frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E}f(X))}_{\Pr[\cdot \geq t] \leq \exp(-2nt^2)} \right| \quad (100)$$

So Hoeffding's inequality can be used to control the term inside the supremum. Let $vc(\mathcal{H}) = d$. By Lemma 41, there are only $O(n^d)$ distinct $(f(X_1), \dots, f(X_n))$ so a union bound implies $t = O\left(\sqrt{\frac{d \log n + \log \frac{1}{\delta}}{2n}}\right)$

This is incorrect because applying Sauer-Shelah requires us to condition on a specific realization of $\{X_i\}_{i=1}^n$ (after which we know there are at most $V_{\mathcal{H}}(n)$ distinct values of $(f(X_1), \dots, f(X_n))$). After conditioning, there's no randomness left for applying Hoeffding's inequality to get concentration.

Solution: Introduce additional randomness using **symmetrization**. Introduce independent copies X'_i and note

$$\mathbb{E}[D_n] = \mathbb{E}_{X_1, \dots, X_n} \left[\sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f(X) \right| \right] \quad (101)$$

$$= \mathbb{E}_{X_1, \dots, X_n} \left[\sup_{f \in \mathcal{H}} \left| \mathbb{E}_{X'_1, \dots, X'_n} \left[\frac{1}{n} \sum_{i=1}^n (f(X_i) - f(X'_i)) \right] \right| \right] \quad (102)$$

$|\cdot|$ is convex, so by Jensen's inequality

$$\mathbb{E}[D_n] \leq \mathbb{E}_X \left[\sup_{f \in \mathcal{H}} \mathbb{E}_{X'} \left[\left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - f(X'_i)) \right| \right] \right] \quad (103)$$

Also, $\sup_y \mathbb{E}f(X, y) \leq \mathbb{E} \sup_y f(X, y)$ for any function f (since $\mathbb{E}f(X, y) \leq \mathbb{E} \sup_y f(X, y)$ then take supremum on left-hand side, or see Fatou-Lebesgue theorem) hence we can move $\mathbb{E}_{X'}$ out of $\sup_{f \in \mathcal{H}}$ to get

$$\mathbb{E}[D_n] \leq \mathbb{E}_{X, X'} \left[\sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - f(X'_i)) \right| \right] \quad (104)$$

Here is where the randomness from symmetrization is added: since $f(X_i) - f(X'_i) \stackrel{d}{=} \varepsilon_i(f(X_i) - f(X'_i))$ for $\varepsilon_i \sim \text{Rad}$

$$\mathbb{E}[D_n] \leq \mathbb{E}_{X, X', \varepsilon} \left[\sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (f(X_i) - f(X'_i)) \right| \right] \quad (105)$$

Condition on X, X' and let $f(X_i) = a \in V_{\mathcal{H}}(\{x_1, \dots, x_n\})$ and $f(X'_i) = b \in V_{\mathcal{H}}(\{x'_1, \dots, x'_n\})$. Then

$$\sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (f(X_i) - f(X'_i)) \right| = \sup_{a, b} \underbrace{\left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (a_i - b_i) \right|}_{\Pr[|\cdot| \geq t] \leq 2 \exp(-\frac{nt^2}{2})} \quad (106)$$

Now we can apply Hoeffding's inequality (picking up an extra factor of 2 because of the absolute value, see Eq. (72)) to the independent, zero-mean (since $\mathbb{E}\varepsilon_i = 0$), bounded (since a_i, b_i , and ε_i are all bounded) random (since ε_i is still random) variables and union bound over the $O(n^{2d})$ (by Sauer-Shelah, squared since there is both $f(X)$ and $f(X')$) distinct $f(X)$ and $f(X')$

$$\Pr \left[\sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (f(X_i) - f(X'_i)) \right| \geq t \mid X, X' \right] \leq (2n^{2d}) 2 \exp \left(-\frac{nt^2}{2} \right) \quad (107)$$

$$(108)$$

This tail probability is small if $t \gg \sqrt{\frac{d \log n}{n}}$, so the expectation over ε in Eq. (105) is of the same order and we have

$$\mathbb{E}[D_n] \leq \mathbb{E}_{X, X'} \left[\mathbb{E}_{\varepsilon} \left[\sup_{f \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (f(X_i) - f(X'_i)) \right| \mid X, X' \right] \right] = O \left(\sqrt{\frac{d \log n}{n}} \right) \quad (109)$$

□

Discretization to a representative set ("fingerprinting") is how previous sections worked. The complication here is that to apply Lemma 41 we had to condition on X_i and remove the randomness. The secret sauce was to add randomness back using the ε_i in symmetrization.

why?? Try
 $\mathbb{E}X = \int P(X \geq t) dt$ for
 $X \geq 0$

4 9/12/2019

4.1 Recap

- Bounded $\mathbb{E} \sup_{v \in V} X(v)$ where $X(v)$ concentrates and V is finite or could be well approximated by a finite set
 - Top eigenvalue of random covariance matrix
 - VC inequality and symmetrization
- Debt: VC-dim of halfspaces is $d + 1$ (Example 40)

Today, we will:

- Pay off debt: prove the VC dimension of half spaces is $d + 1$
- Give a finite-sample analysis of Definition 1 (Minimum distance functional)
 - Weaken TV to $\widetilde{\text{TV}}$
 - Bound Modulus of continuity bound via "mean crossing lemma"
 - $\widetilde{\text{TV}}(\tilde{p}, \tilde{p}_n) \rightarrow 0$ as $n \rightarrow \infty$

4.2 VC dimension of half spaces

In Example 40 (VC dimension of half spaces) we claimed that $vc(\mathcal{H}) = d + 1$ for the family of half spaces (i.e. linear decision boundaries)

$$\mathcal{H} = \{\mathbb{1}\{\langle v, x \rangle \geq \tau\} : v \in \mathbb{R}^d, \tau \in \mathbb{R}\} \quad (110)$$

We previously showed it geometrically for the case when $d = 2$. Here, we will generalize this to higher dimensions.

Proposition 43 (VC dimension of half spaces)

No $d + 2$ set of points in \mathbb{R}^d can be shattered by any $f \in \mathcal{H}$.

Proof. Fix $\{x_i\}_{i=1}^{d+2} \in \mathbb{R}^d$ distinct. We will find two sets $S_+, S_- \subset \{x_1, \dots, x_{d+2}\}$ such that $S_+ \cap S_- = \emptyset$ but $\text{conv}(S_+) \cap \text{conv}(S_-) \neq \emptyset$. This is sufficient because every $f = \mathbb{1}\{\langle v, x \rangle \geq \tau\} \in \mathcal{H}$ can be identified with a half-space (of the points classified +1 by f)

$$H = f^{-1}(\{1\}) = \{x \in \mathbb{R}^d : \langle v, x \rangle \geq \tau\} \quad (111)$$

and by convexity of H

$$S_+ \subset H \implies \text{conv}(S_+) \subset H \quad (112)$$

Hence, if f correctly classifies all of S_+ then it must also misclassify some $x \in S_+ \cap S_- \subset S_-$.

Consider the linear system

$$\sum_{i=1}^{d+2} a_i x_i = 0, \quad \sum_{i=1}^{d+2} a_i = 0 \quad (113)$$

or equivalently in matrix form

$$\underbrace{\begin{bmatrix} \vdots & \vdots & \dots & \vdots \\ x_1 & x_2 & \dots & x_{d+2} \\ \vdots & \vdots & \dots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix}}_{(d+1) \times (d+2)} \begin{bmatrix} a_1 \\ \vdots \\ a_{d+2} \end{bmatrix} = \mathbf{0} \quad (114)$$

By the rank-nullity theorem, the null-space must have dimension ≥ 1 hence there exists at least one solution \mathbf{a} . Let

$$S_+ = \{i : a_i > 0\}, \quad S_- = \{i : a_i < 0\} \quad (115)$$

Then by Eq. (113)

$$\underbrace{\sum_{i \in S_+} \underbrace{\frac{a_i}{A}}_{\in [0,1]} x_i}_{\in \text{conv}(S_+)} = \sum_{i \in S_-} \underbrace{\frac{a_i}{A}}_{\in [0,1]} x_i \quad \text{where} \quad A = \sum_{i \in S_+} a_i = \sum_{i \in S_-} (-a_i) \quad (116)$$

This gives us a point in $\text{conv}(S_+) \cap \text{conv}(S_-)$. □

Remark 44. The geometric result that “any set of $d + 2$ points in \mathbb{R}^d can be partitioned into two disjoint sets whose convex hulls intersect” is known as **Radon’s theorem** on convex sets.

4.3 Finite sample analysis of MDF via Generalized KS distance

Recall Definition 1 (Minimum distance functional) projects \tilde{p} on to \mathcal{G} under some discrepancy D . Previously we worked with $D = \text{TV}$, which works fine if \tilde{p} is a continuous distribution (e.g. $\tilde{p} = \mathcal{N}(\mu, I)$ in Lemma 3). However, when we only have a finite number of samples we can only form the empirical distribution

$$\tilde{p}_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}, \quad X_i \sim \tilde{p} \quad (117)$$

Here, TV is inadequate because $\text{TV}(\tilde{p}_n, p) = 1$ almost surely for any continuous distribution p (this is because $\Pr_{X \sim p}[X = X_i] = 0$) so it's not clear how to project onto a continuous family such as $\mathcal{G}_{\text{gauss}}$. Moreover, in many cases $\text{TV}(\tilde{p}_n, \tilde{p}) = 1$ even as $n \rightarrow \infty$.

To address this issue, we can consider relaxing TV to a weakening $\widetilde{\text{TV}}$ which is more forgiving. We have two desiderata for $\widetilde{\text{TV}}$:

1. The modulus $\mathbf{m}(\mathcal{G}, \varepsilon, \widetilde{\text{TV}})$ remains small, so that Proposition 2 (Modulus of continuity bound) still gives a good result
2. $\widetilde{\text{TV}}(\tilde{p}, \tilde{p}_n) \rightarrow 0$ as $n \rightarrow \infty$, so that $\widetilde{\text{TV}}$ detects convergence of (discrete) empirical distributions to a (possibly continuous) population distribution

Remark 45. The two desiderata are competing. We want $\widetilde{\text{TV}}$ to be large in (1) so that $A = \{(p, q) \in \mathcal{G} : \widetilde{\text{TV}}(p, q) \leq \varepsilon\}$ is small and hence $\mathbf{m} = \sup_{(p, q) \in A} L(p, \theta^*(q))$ is small. At the same time, in (2) we would like $\widetilde{\text{TV}}$ to be small to avoid the failure of TV in detecting $\tilde{p}_n \rightarrow \tilde{p}$ (e.g. Glivenko-Cantelli ensures that the cumulative distribution functions converge uniformly).

Proposition 46

Suppose $\widetilde{\text{TV}}$ is a pseudometric such that $\widetilde{\text{TV}} \leq \text{TV}$. Let $\hat{\theta}_{\widetilde{\text{TV}}}(p) = \theta^*(q)$ where $q \in \arg\min_{q \in \mathcal{G}} \widetilde{\text{TV}}(p, q)$ (the Minimum distance functional under $\widetilde{\text{TV}}$). Then

$$L(p^*, \hat{\theta}_{\widetilde{\text{TV}}}(\tilde{p}_n)) \leq \mathbf{m}(\mathcal{G}, 2\varepsilon', \widetilde{\text{TV}}) \quad (118)$$

where $\varepsilon' = \varepsilon + \widetilde{\text{TV}}(\tilde{p}, \tilde{p}_n)$ (and $\widetilde{\text{TV}}(p^*, \tilde{p}) \leq \varepsilon$ as per the conventions outlined in Fig. 1)

Proof. By Proposition 2 (Modulus of continuity bound)

$$L(p^*, \hat{\theta}_{\widetilde{\text{TV}}}(\tilde{p}_n)) \leq \mathbf{m}(\mathcal{G}, 2\widetilde{\text{TV}}(p^*, \tilde{p}_n), \widetilde{\text{TV}}, L) \quad (119)$$

Since $\widetilde{\text{TV}}$ is a pseudometric, by the triangle inequality

$$\widetilde{\text{TV}}(p^*, \tilde{p}_n) \leq \underbrace{\widetilde{\text{TV}}(p^*, \tilde{p})}_{\leq \varepsilon} + \widetilde{\text{TV}}(\tilde{p}, \tilde{p}_n) \quad (120)$$

□

How do we construct $\widetilde{\text{TV}}$?

Definition 47 (Generalized Kolmogorov-Smirnov distance)

For a family of functions $\mathcal{H} = \{f : \mathcal{X} \rightarrow \mathbb{R}\}$, the *generalized Kolmogorov-Smirnov distance* induced by \mathcal{H} is

$$\widetilde{\text{TV}}_{\mathcal{H}}(p, q) = \sup_{f \in \mathcal{H}, \tau \in \mathbb{R}} \left| \Pr_p[f(X) \geq \tau] - \Pr_q[f(X) \geq \tau] \right| \quad (121)$$

Remark 48. For $f \in \mathcal{H}$ and $\tau \in \mathbb{R}$, if we define the event $E_{f,\tau} = \{f(X) \geq \tau\}$ then notice

$$\widetilde{\text{TV}}_{\mathcal{H}}(p, q) = \sup_{E_{f,\tau}} \left| \Pr_p[E_{f,\tau}] - \Pr_q[E_{f,\tau}] \right| \leq \sup_{E \text{ meas}} \left| \Pr_p[E] - \Pr_q[E] \right| = \text{TV}(p, q) \quad (122)$$

So $\widetilde{\text{TV}}$ is indeed dominated by TV as required by Proposition 46.

What \mathcal{H} should we pick? The answer depends on what we are trying to estimate (i.e. choice of $L(p, \theta)$). For now, consider mean estimation (i.e. $L(p, \theta) = \|\theta - \mu(p)\|_2$). One intuition is that knowledge of the one dimensional projections ($\mathbb{E} \langle v, x \rangle$ for all $v \in \mathbb{R}^d$) allows us to know $\mathbb{E}[X]$, so it's reasonable to consider

$$\mathcal{H} = \mathcal{H}_{\text{lin}} = \{x \mapsto \langle v, x \rangle : v \in \mathbb{R}^d\} \quad (123)$$

To bound the modulus, recall that previously if $p, q \in \mathcal{G}_{\text{TV}}$ are Resilient distributions then Corollary 11 gave us

$$\text{TV}(p, q) \leq \varepsilon \implies \|\mu(p) - \mu(q)\|_2 \leq 2\rho \quad (124)$$

Similarly, here we will also restrict our distributional assumptions to be within resilient distributions: $\mathcal{G} \subset \mathcal{G}_{\text{TV}}$.

We need to show our two desiderata:

1. The modulus is bounded:

$$p, q \in \mathcal{G} \subset \mathcal{G}_{\text{TV}}(\rho, \varepsilon) \text{ and } \widetilde{\text{TV}}(p, q) \leq \varepsilon \implies \|\mu(p) - \mu(q)\|_2 \leq \sigma = 2\rho \quad (125)$$

2. $\widetilde{\text{TV}}(\tilde{p}, \tilde{p}_n)$ is small, specifically:

$$\widetilde{\text{TV}}(\tilde{p}, \tilde{p}_n) = O\left(\sqrt{d/n}\right) \quad (126)$$

Proof of Eq. (125). Previously we used Midpoint lemma to find an ε -deletion $r \leq \min\left\{\frac{p}{1-\varepsilon}, \frac{q}{1-\varepsilon}\right\}$ close to both p and q in the sense that

$$\|\mu(p) - \mu(r)\|_2 \leq \rho \text{ and } \|\mu(q) - \mu(r)\|_2 \leq \rho \quad (127)$$

After which a triangle inequality completed the proof.

Unfortunately, we don't know of a way to find a single midpoint distribution under $\widetilde{\text{TV}}$. Instead, we will use the following key property:

Lemma 49 (Mean crossing property)

Suppose $\widetilde{\text{TV}}(p, q) \leq \varepsilon$. For any $v \in \mathbb{R}^d$, there exists ε -deletions $r_p \leq \frac{p}{1-\varepsilon}$ and $r_q \leq \frac{q}{1-\varepsilon}$ such that

$$\mathbb{E}_{r_q} \langle v, x \rangle \leq \mathbb{E}_{r_p} \langle v, x \rangle \quad (128)$$

In other words, after deleting ε mass to create r_q and r_p , the means are shifted such that they cross.

If we have the ε deletions $r_p \leq \frac{p}{1-\varepsilon}$ and $r_q \leq \frac{q}{1-\varepsilon}$ from Lemma 49 (Mean crossing property), then

$$\underbrace{\mathbb{E}_p \langle v, x \rangle}_{=\langle v, \mu_p \rangle} \leq \mathbb{E}_{r_p}[\langle v, x \rangle] + \rho \quad \text{resilience of } p \quad (129)$$

$$\leq \mathbb{E}_{r_q}[\langle v, x \rangle] + \rho \quad \text{mean crossing} \quad (130)$$

$$\leq \underbrace{\mathbb{E}_q[\langle v, x \rangle]}_{=\langle v, \mu_q \rangle} + 2\rho \quad \text{resilience of } q \quad (131)$$

Hence

$$\langle v, \mu_p - \mu_q \rangle \leq 2\rho \quad (132)$$

for all $\|v\|_2 = 1$. Therefore $\|\mu_p - \mu_q\|_2 \leq 2\rho$. \square

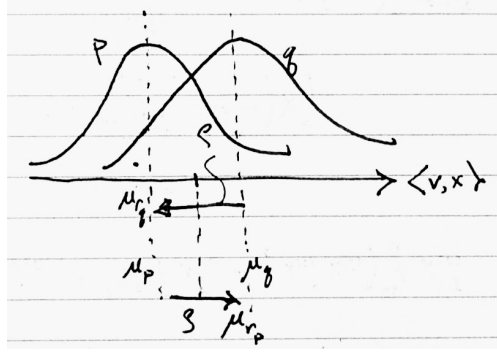


Figure 8: Resilience allows us to perform an ε -deletion to move from $\mu_p \rightarrow \mu_{r_p}$ and $\mu_q \rightarrow \mu_{r_q}$ and pick up a factor of $+2\rho$. Mean crossing allows us to relate μ_{r_p} and μ_{r_q} .

Proof of Mean crossing property. Consider Fig. 8, which visualizes the 1D projections of p and q in the v direction. To make $\langle v, \mu_{r_q} \rangle$ cross over $\langle v, \mu_{r_p} \rangle$, we would like to shift the mean of q to the left and the mean of p to the right as much as possible. Thus, delete ε mass from the right tail of q (and delete the left tail of p). Then

$$\Pr_{r_p}[\langle v, x \rangle \geq \tau] \geq \frac{\Pr_p[\langle v, x \rangle \geq \tau]}{1 - \varepsilon} \geq \frac{\Pr_q[\langle v, x \rangle \geq \tau] - \varepsilon}{1 - \varepsilon} = \Pr_{r_q}[\langle v, x \rangle \geq \tau] \quad (133)$$

where the first inequality is because r_p is p with the left tail deleted and renormalized by $1 - \varepsilon$, the second from $\Pr_q[\langle v, x \rangle \geq \tau] - \Pr_p[\langle v, x \rangle \geq \tau] \leq \widehat{\text{TV}}(p, q) \leq \varepsilon$, and the third from r_q being formed by deleting ε from the right tail of q and renormalizing by $1 - \varepsilon$.

We have shown that the right tail probabilities of r_p are always larger than those of r_q , i.e. r_p **stochastically dominates** r_q . As a consequence, $\mathbb{E}_{r_p}[\langle v, x \rangle] \geq \mathbb{E}_{r_q}[\langle v, x \rangle]$. \square

Proof of Eq. (126). Notice

$$\widehat{\text{TV}}_{\mathcal{H}}(p, q) = \sup_{v \in \mathbb{R}^d, \tau \in \mathbb{R}} \underbrace{\left| \Pr_p[\langle v, x \rangle \geq \tau] - \Pr_q[\langle v, x \rangle \geq \tau] \right|}_{\text{max discrepancy on halfspaces}} \quad (134)$$

By VC inequality and Proposition 43 (VC dimension of half spaces)

$$\widehat{\text{TV}}_{\mathcal{H}}(\tilde{p}, \tilde{p}_n) \leq O\left(\sqrt{\frac{vc(\text{half spaces})}{n}}\right) = O\left(\sqrt{\frac{d + \log \frac{1}{\delta}}{n}}\right) \quad (135)$$

with probability $\geq 1 - \delta$. \square

Consequences:

- For $(\rho, \varepsilon + O(\sqrt{d/n}))$ -resilient distributions, we can estimate mean with error 2ρ
- For bounded covariance, Corollary 5 gave us $\rho(\varepsilon) = O(\sqrt{\varepsilon})$ hence

$$L(p^*, \tilde{\theta}_{\widehat{\text{TV}}}(\tilde{p}_n)) \leq O\left(\sqrt{\varepsilon + \sqrt{d/n}}\right) \quad (136)$$

The lower bound $\sqrt{\varepsilon}$ is what we get in the infinite sample $n \rightarrow \infty$ limit, and $\sqrt{d/n}$ when $\varepsilon \rightarrow 0$, so we would like $\sqrt{\varepsilon} + \sqrt{d/n}$. The slack in the bound comes from $n \gg d/\varepsilon^2$, whereas we would need $n \gg \frac{d}{\varepsilon}$ but this analysis doesn't give it to us.

- For sub-Gaussians, $\rho(\varepsilon) = O(\varepsilon\sqrt{\log(1/\varepsilon)})$. When $n \gg \frac{d}{\varepsilon^2}$ we get $O(\varepsilon\sqrt{\log(1/\varepsilon)})$.

In general, this analysis holds for $n \gg d/\varepsilon^2$: whenever this holds, we can do as well as if we had infinite data. The analysis is tight in d but loose in ε .

Bibliography

Donoho, D. L., R. C. Liu, et al.

1988. The” automatic” robustness of minimum distance functionals. *The Annals of Statistics*, 16(2):552–586.

\widetilde{TV} similar to Tukey median, may be useful for challenge problem