

Robust Mean Estimation

Why Should We Care About Robustness

- model misspecification
- crowdsourcing data
 - * Quality vs Quantity tradeoff.
- Data poisoning attacks

Contamination Models

① No Contamination

Samples x_1, \dots, x_n are drawn from $N(\mu, I_d)$
independently & identically.

② Strong Contamination

Samples x_1, \dots, x_n drawn as before, except
adversary may remove ϵn points &
add ϵn points as she wants!

Contamination Models

① No Contamination

Samples x_1, \dots, x_n are drawn from $N(\mu, I_d)$
independently & identically.

Goal: output $\hat{\mu}$
close to true μ

② Strong Contamination

Samples x_1, \dots, x_n drawn as before, except
adversary may remove ϵn points &
add ϵn points as she wants!

No Contamination

Thm (Chernoff - Hoeffding)

Let X_1, \dots, X_n be (sub) Gaussian r.v.'s & independent
w/ mean $\mathbb{E} X_i = \mu$.

$$\Rightarrow \Pr \left[\frac{1}{n} \sum X_i - \mu \geq t \right] \leq O(\exp(-nt^2))$$

ex Bernoulli

ex Gaussian

ex
Bounded

No Contamination

Thm (Chernoff - Hoeffding)

Let X_1, \dots, X_n be (sub) Gaussian r.v.'s & independent
w/ mean $\mathbb{E} X_i = \mu$.

$$\Rightarrow \Pr \left[\frac{1}{n} \sum X_i - \mu \geq t \right] \leq O(\exp(-nt^2))$$

ex Bernoulli

ex Gaussian

ex Bounded

- After $n \geq \Omega(1/\varepsilon^2)$ samples, estimate up to $\pm \varepsilon$ in 1-D

- After $n \geq \Omega(d/\varepsilon^2)$ samples, " " " " in d-D

Is the Mean Robust?

No! A single corrupted point causes drift
arbitrarily far from μ !

(take some point to ∞ !)

but ...

Is the Mean Robust?

No! A single corrupted point causes drift
arbitrarily far from μ !

(take some point to ∞ !)

but...

the median naturally ignores points far from μ

↖ or any order statistic

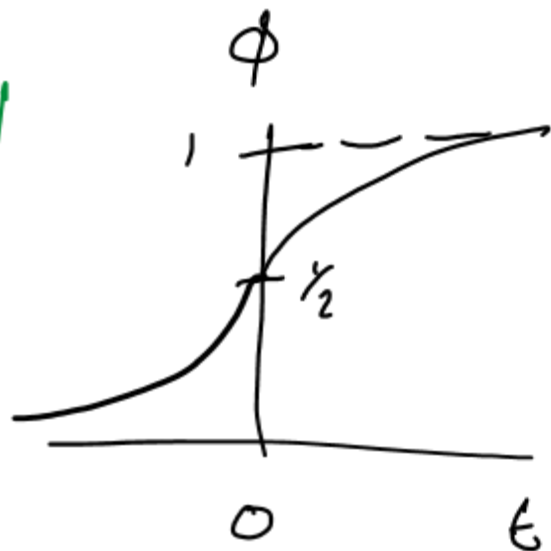
⇒ up next: the median achieves $O(\epsilon)$ error in 1-D.

Analysis of the Median

Thm Let ϕ be the Gaussian CDF,

$$\mathbb{P}[|\text{med}(S) - \mu| > t\sigma] \leq 2 \exp(-2n(\phi(t) - \frac{1}{2} - \varepsilon)^2)$$

$$\Rightarrow \mathbb{P}[|\text{med}(S) - \mu| > (1+b)\sqrt{2\pi}\varepsilon] \leq 2 \exp(-\Omega(b^2 n \varepsilon^2))$$

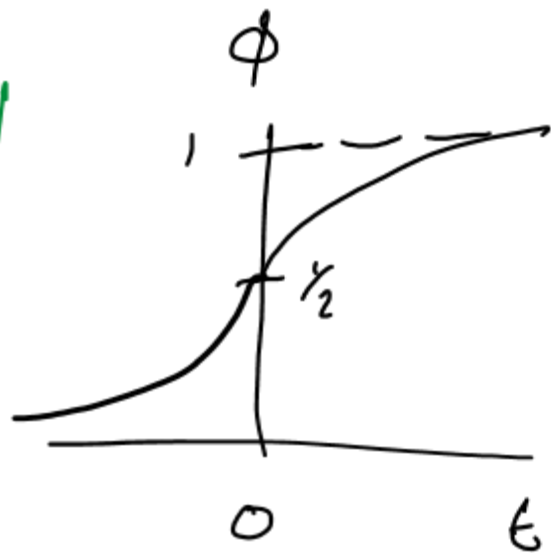


Analysis of the Median

Thm Let ϕ be the Gaussian CDF,

$$\mathbb{P}[|\text{med}(S) - \mu| > t\sigma] \leq 2 \exp(-2n(\phi(t) - \frac{1}{2} - \varepsilon)^2)$$

$$\Rightarrow \mathbb{P}[|\text{med}(S) - \mu| > (1+b)\sqrt{2\pi}\varepsilon] \leq 2 \exp(-\Omega(\ln \varepsilon^2))$$

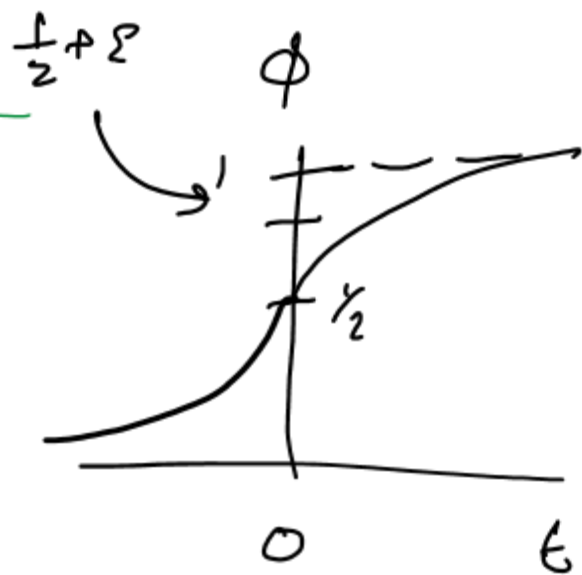


Thus, $|\text{med}(S) - \mu| = \underbrace{O(\varepsilon)}_{\text{optimal rate!!}}$ as long as $n = \Omega(1/\varepsilon^2)$.

Analysis of the Median

Thm Let ϕ be the Gaussian CDF,

$$\mathbb{P}[|\text{med}(S) - \mu| > t\sigma] \leq 2 \exp(-2n(\phi(t) - \frac{1}{2} - \varepsilon)^2)$$



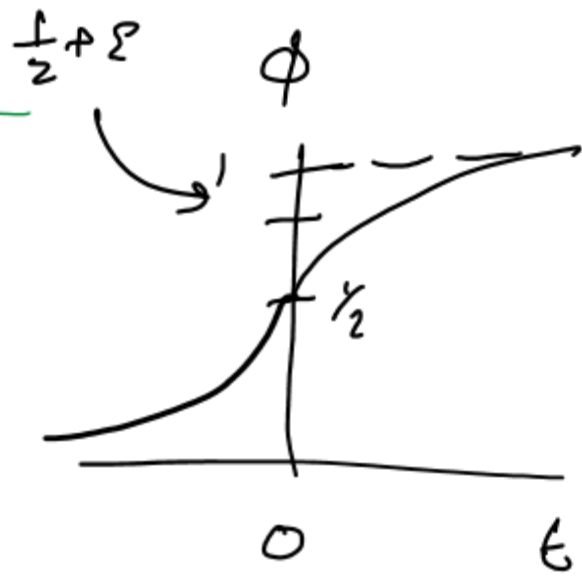
pf After ε -corruption, the median can move only up to $\frac{1}{2} + \varepsilon$ of the original CDF

$\Rightarrow |\text{med}(S) - \mu| > t\sigma$ if & only if $\frac{1}{2} - \varepsilon$ fraction of pts exceed $\mu + t\sigma$

Analysis of the Median

Thm Let ϕ be the Gaussian CDF,

$$\mathbb{P}[|\text{med}(S) - \mu| > t\sigma] \leq 2 \exp(-2n(\phi(t) - \frac{1}{2} - \varepsilon)^2)$$



pf After ε -corruption, the median can move only up to $\frac{1}{2} + \varepsilon$ of the original CDF

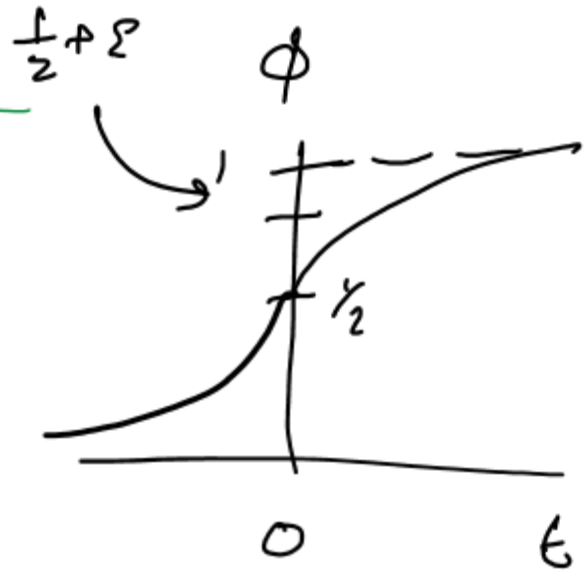
$\Rightarrow |\text{med}(S) - \mu| > t\sigma$ if & only if $\frac{1}{2} - \varepsilon$ fraction of pts exceed $\mu + t\sigma$
the probability of each event is $1 - \phi(t)$:

$$\mathbb{P}\left[\frac{1}{n} \sum Y_i \geq 1 - \phi(t) + \varepsilon\right] \leq \exp(-2n\varepsilon^2)$$

Analysis of the Median

Thm Let ϕ be the Gaussian CDF,

$$\mathbb{P}[|\text{med}(S) - \mu| > t\sigma] \leq 2 \exp(-2n(\phi(t) - \frac{1}{2} - \varepsilon)^2)$$



pf After ε -corruption, the median can move only up to $\frac{1}{2} + \varepsilon$ of the original CDF

$\Rightarrow |\text{med}(S) - \mu| > t\sigma$ if & only if $\frac{1}{2} - \varepsilon$ fraction of pts exceed $\mu + t\sigma$
the probability of each event is $1 - \phi(t)$:

$$\mathbb{P}\left[\frac{1}{n} \sum Y_i \geq 1 - \phi(t) + \varepsilon\right] \leq \exp(-2n\varepsilon^2)$$

\Rightarrow taking $\varepsilon = \phi(t) - \frac{1}{2} - \varepsilon$ gives the desired result!

How can we generalize to d -dimensions?

Note: the goal is to learn μ to min L_2 error

How can we generalize to d -dimensions?

Note: the goal is to learn μ to $\min L_2$ error

① compute median in all directions...

... this gives $O(\epsilon)$ in each, so $O(\epsilon\sqrt{d})$ total.

How can we generalize to d -dimensions?

Note: the goal is to learn μ to $\min L_2$ error

① Compute median in all directions...

... this gives $O(\epsilon)$ in each, so $O(\epsilon\sqrt{d})$ total.

② "Tukey Median" gives $O(\epsilon)$ total, but requires
computing ∞ -many medians*

How can we generalize to d -dimensions?

Note: the goal is to learn μ to $\min L_2$ error

① Compute median in all directions...

... this gives $O(\epsilon)$ in each, so $O(\epsilon \sqrt{d})$ total.

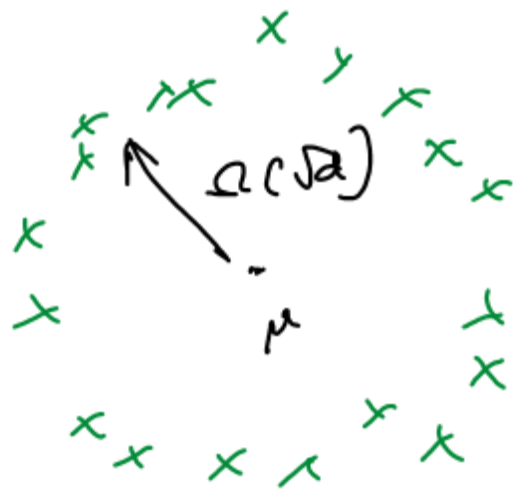
② "Tukey Median" gives $O(\epsilon)$ total, but requires
computing ∞ -many medians*

③ Polytine alternatives that give $O(\epsilon \sqrt{\log(1/\epsilon)})$ total !!

High-Dimensional Intuitions

- What is the "adversarial" strategy in \mathbb{R}^d ?

Fact: $X \sim N(\mu, I_d)$ concentrates around $\Theta(d)$
from μ , w/ error $O(1)$



High-Dimensional Intuitions

- What is the "adversarial" strategy in \mathbb{R}^d ?

Fact: $X \sim N(\mu, I_d)$ concentrates around $\Theta(d)$
from μ , w/ error $O(1)$

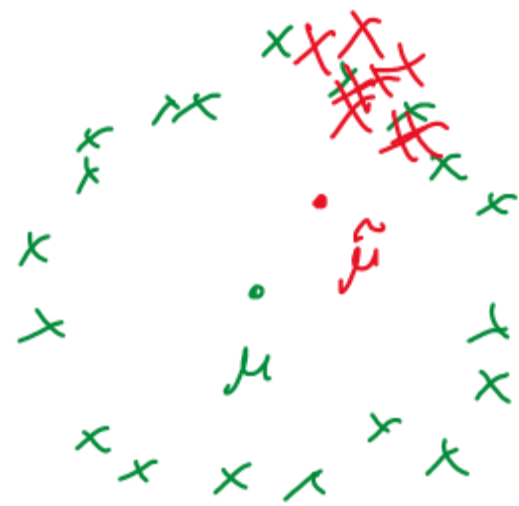


Rather than putting points at ∞ , the adversary
has to "clump" points to shift the mean!

High-Dimensional Intuitions

- What is the "adversarial" strategy in \mathbb{R}^d ?

Fact: $X \sim N(\mu, I_d)$ concentrates around $\Theta(d)$
from μ , w/ error $O(1)$



Rather than putting points at ∞ , the adversary has to "clump" points to shift the mean!

* If $\mu \rightarrow \tilde{\mu}$ by corruption, the "variance" in direction $\tilde{\mu} - \mu$ must be large

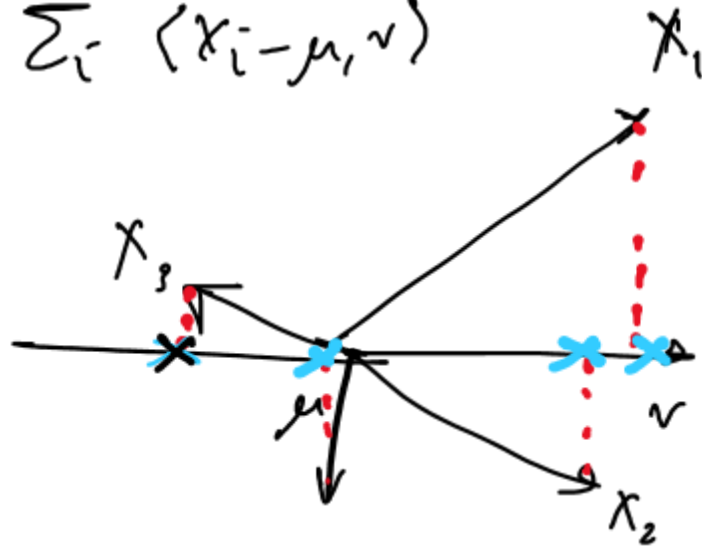
Covariance Function

Definition The covariance matrix of $x_1, \dots, x_n \in \mathbb{R}^d$ is

$$\Sigma = \frac{1}{n} \sum_i (x_i - \mu)(x_i - \mu)^T; \quad \mu = \frac{1}{n} \sum_{i=1} x_i$$

The covariance naturally encodes variance along all directions

$$v^T \Sigma v = \frac{1}{n} \sum_i v^T (x_i - \mu)(x_i - \mu)^T v = \frac{1}{n} \sum_i \langle x_i - \mu, v \rangle^2$$



Covariance Function

Definition The covariance matrix of $x_1, \dots, x_n \in \mathbb{R}^d$ is

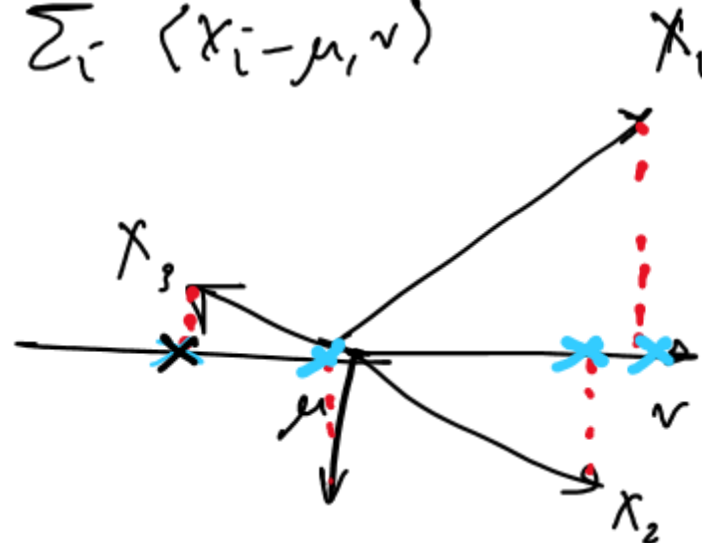
$$\Sigma = \frac{1}{n} \sum_i (x_i - \mu)(x_i - \mu)^T; \quad \mu = \frac{1}{n} \sum_{i=1} x_i$$

The covariance naturally encodes variance along all directions

$$v^T \Sigma v = \frac{1}{n} \sum_i v^T (x_i - \mu)(x_i - \mu)^T v = \frac{1}{n} \sum_i \langle x_i - \mu, v \rangle^2$$

So, the "top principal component" is

$\arg \max_{\|v\|=1} v^T \Sigma v$, the direction of
greatest variance.



High-Dimensional Intuitions

If the corrupted points move μ by distance δ

\Rightarrow corrupted points have projection onto v
at least δ/ε



High-Dimensional Intuitions

If the corrupted points move μ by distance δ

\Rightarrow corrupted points have projection onto v
at least δ/ϵ .

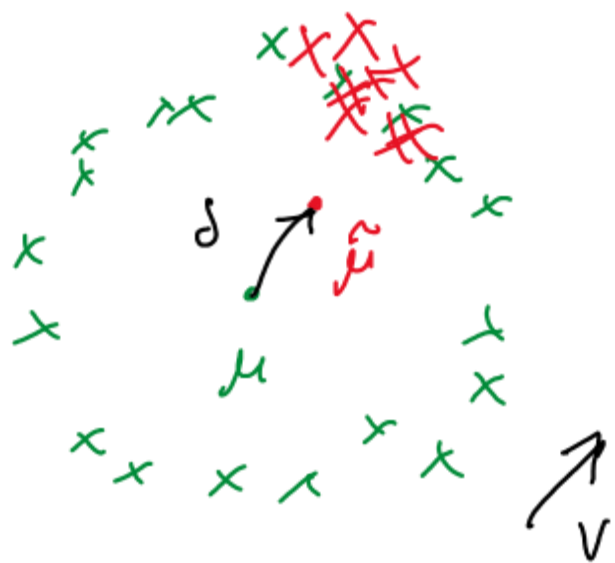
\Rightarrow corrupted points have $\epsilon \cdot (\delta/\epsilon)^2 = \delta^2/\epsilon$ variance onto v .



High-Dimensional Intuitions

If the corrupted points move μ by distance δ

\Rightarrow corrupted points have projection onto v
at least δ/ϵ .



\Rightarrow corrupted points have $\epsilon \cdot (\delta/\epsilon)^2 = \delta^2/\epsilon$ variance onto v .

Conjecture: A "random" set satisfies that if the

corrupted version has "small" operator norm

$\Rightarrow \|\mu_T - \mu_X\|_2$ is "small".

Stability -

We formalize the notion of "random" as a stable set

Defⁿ (ϵ, δ) -stable wrt X for $\epsilon \in (0, \frac{1}{2})$, $\delta \geq \epsilon$ if

$\forall v \in S^{d-1}$, $S' \subseteq S$ of size at least $(1-\epsilon)|S|$,

1 $\left| \frac{1}{|S'|} \sum_{x \in S'} \langle v, x - \mu_x \rangle \right| \leq \delta \leftarrow$

2 $\left| \frac{1}{|S'|} \sum_{x \in S'} \langle v, x - \mu_x \rangle^2 - 1 \right| \leq \delta^2 / \epsilon \leftarrow$

variance
bound

mean bound

Intuition: Ability to "match" moments
on a sample up to distortions

Stability -

We formalize the notion of "random" as a stable set

Defⁿ (ϵ, δ) -stable w.r.t X for $\epsilon \in (0, \frac{1}{2})$, $\delta \geq \epsilon$ if

$\forall v \in S^{d-1}$, $S' \subseteq S$ of size at least $(1-\epsilon)|S|$,

$$1 \quad \left| \frac{1}{|S'|} \sum_{x \in S'} \langle v, x - \mu_X \rangle \right| \leq \delta \quad \rightsquigarrow \epsilon \sqrt{\log(1/\epsilon)}$$

$$2 \quad \left| \frac{1}{|S'|} \sum_{x \in S'} \langle v, x - \mu_X \rangle^2 - 1 \right| \leq \delta^2 / \epsilon \quad \rightsquigarrow \epsilon \cdot \log(1/\epsilon)$$

Prop: A subgaussian distribution is $(\epsilon, O(\epsilon \sqrt{\log 1/\epsilon}))$ -stable



Thm (Conjecture)

Let S be (ϵ, δ) set wrt X , let T be ϵ -corrupted of S

if μ_T, Σ_T are mean & variance, $\|\Sigma_T\| \leq 1 + \lambda$

$$\Rightarrow \|\mu_T - \mu_X\|_2 \leq O(\delta + \sqrt{\epsilon \lambda})$$

Cor If W is a distribution differing from S by TV of ϵ .

$$\text{if } \|\Sigma_W\|_2 \leq 1 + \lambda \Rightarrow \|\mu_W - \mu_X\|_2 \leq O(\delta + \sqrt{\epsilon \lambda})$$

Intuition: W is a "weight" corresponding to each datapoint if "confidence in S "

Defⁿ $C = \{ W : W \text{ supp on } T, W(u) \leq \frac{1}{|T|(1-\epsilon)} \}$



Defⁿ $C = \left\{ W : W \text{ supp on } T, W(n) \leq \frac{1}{|T|(1-\epsilon)} \right\}$



Prop If, for some $W \in C$, Σ_W has no large eigenvalue

$\Rightarrow \mu_W$ is a good approx to μ_X
(by previous lemma)

Defⁿ $\mathcal{C} = \left\{ W : W \text{ supp on } T, W(u) \leq \frac{1}{|T|(1-\epsilon)} \right\}$



* Take $W^* = \text{unif}(T \cap S)$

By property: $\|\Sigma_{W^*}\| \leq 1 + \delta^2/\epsilon$

gives l_2 error $O(\delta)$ for
S being $(3\epsilon, \delta)$ -stable!

Prop If, for some $W \in \mathcal{C}$, Σ_W has no large eigenvalue

$\Rightarrow \mu_W$ is a good approx to μ_X
(by previous lemma)

Summary

1 Consistency condition w/ removing samples

* (ϵ, δ) -stable set

Summary

1 Consistency condition w/ removing samples

* (ϵ, δ) -stable set

2 Consistency condition happens whp on subGaussian distr

* for n large - $(\epsilon, O(\epsilon \sqrt{\log 1/\epsilon}))$ -stable

3 Corruption of consistency spectrally certifies closeness in l_2

* Any TV-close corruption shows a covariance $\rightarrow l_2$ reduction

Summary

1 Consistency condition w/ removing samples

* (ϵ, δ) -stable set

2 Consistency condition happens whp on subGaussian distr

* for n large - $(\epsilon, O(\epsilon \sqrt{\log 1/\epsilon}))$ -stable

3 Corruption of consistency spectrally certifies closeness in l_2

* Any TV-close corruption shows a covariance $\rightarrow l_2$ reduction

4 Produce a convex set \leadsto the optimum being the **set** of uncorrupted points.

Towards Robust Algorithms

Reduced mean estimation to a "convex" problem on the
set \mathcal{C} of weight distributions

strictly speaking,
not true, but
technique applies.

Towards Robust Algorithms

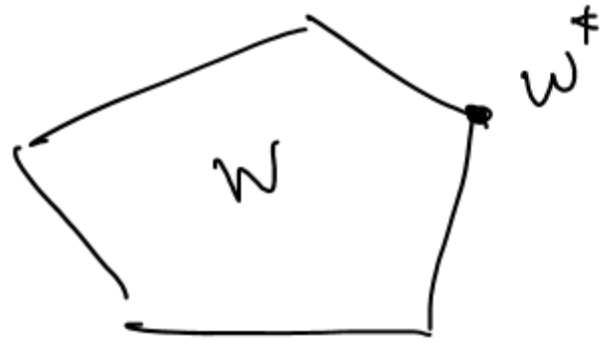
Reduced mean estimation to a "convex" problem on the set \mathcal{C} of weight distributions

Use the "ellipsoid method" to approach W^* .

Towards Robust Algorithms

Reduced mean estimation to a "convex" problem on the set \mathcal{C} of weight distributions

Use the "ellipsoid method" to approach w^* .

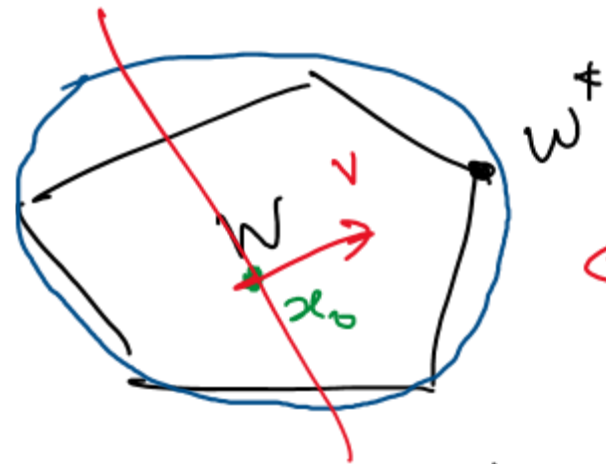


feasible set Δ
optima

Towards Robust Algorithms

Reduced mean estimation to a "convex" problem on the set \mathcal{C} of weight distributions

Use the "ellipsoid method" to approach w^* .



feasible set &
optima

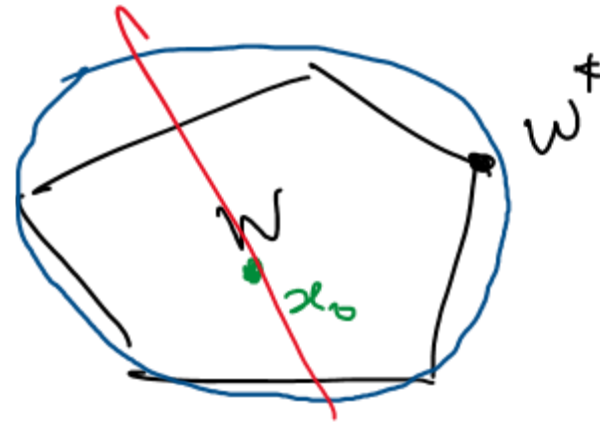
& center of ellipse

separation oracle
gives a direction st
 $\langle v, x_0 \rangle < \langle v, w^* \rangle$

Towards Robust Algorithms

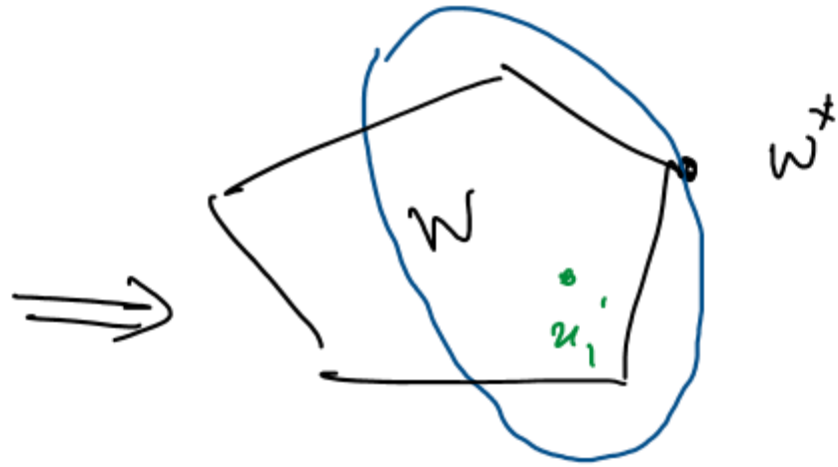
Reduced mean estimation to a "convex" problem on the set \mathcal{C} of weight distributions

Use the "ellipsoid method" to approach w^* .



feasible set &
optima

& center of ellipse

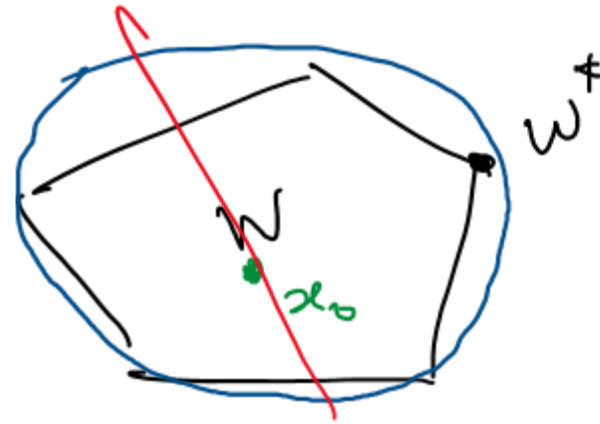


new feasible set
& ellipse.

Towards Robust Algorithms

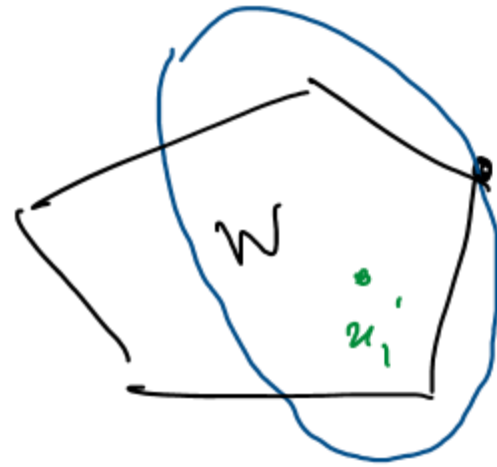
Reduced mean estimation to a "convex" problem on the set \mathcal{C} of weight distributions

Use the "ellipsoid method" to approach W^*



feasible set &
optima

& center of ellipse



new feasible set
& ellipse.

Without detail

\exists a separation oracle

for this nonconvex problem

\Rightarrow poly(d) iters suffice
to approach W^* .

Towards Robust Algorithms

An even simpler (but harder to analyze) approach is u/a filtering

Towards Robust Algorithms

An even simpler (but harder to analyze) approach is u/a filtering

Alg

1 If no directions have large variance, output empirical mean

By ~~prop~~ on stable
sees



Towards Robust Algorithms

An even simpler (but harder to analyze) approach is via filtering

Alg

- 1 If no directions have large variance, output empirical mean
 - 2 Throw out data along the high variance direction that is "too far" from the mean
- ← "trimmed mean" analogue

Towards Robust Algorithms

An even simpler (but harder to analyze) approach is via filtering

Alg

- 1 If no directions have large variance, output empirical mean
- 2 Throw out data along the high variance direction that is "too far" from the mean

Analysis: Each time we throw out data, the density of good points increases!

LEMMA 3.3 (INFORMAL). Let $S = G \cup E \setminus S_r$ be an ε -corrupted set of points from $\mathcal{N}(\mu, I)$ of size at least $\tilde{\Omega}(d / \varepsilon^2)$. Then, with probability at least 0.99 and after a simple preprocessing step, the filter satisfies the following property: Given any $S' \subseteq S$ satisfying $\Gamma(G, S') \leq 2\varepsilon$, the filter either

- (a) outputs $\hat{\mu}$ so that $\|\mu - \hat{\mu}\|_2 \leq O(\varepsilon \sqrt{\log 1/\varepsilon})$, or
- (b) outputs T so that $\Gamma(G, T) \leq \Gamma(G, S') - \varepsilon/\alpha$, where $\alpha = d \log(d/\varepsilon) \log(d \log(d/\varepsilon))$.

iteration complexity

value of δ

Thanks For Listening!