# PAC Learning Linear Thresholds from Label Proportions

Anand Brahmbhatt, Rishi Saket and Aravindan Raghuveer

Google Research India

# Problem Statement

❖ **PAC Learning**: For a function $f$: $\mathbb{R}^d \rightarrow$ **{0, 1}**, given **m** samples **(x, $f$(x))** where **x ~ $\mathcal{D}$**, find a hypothesis $h$ s.t. $\mathbf{Pr_{x \sim \mathcal{D}}[\mathit{h}(x) \neq \mathit{f}(x)] \leq \varepsilon}$ w.p. **1 - δ**. Efficient if **m ≤ O(poly(d,1/ε, log(1/δ)))**.

# Problem Statement

❖ **PAC Learning**: For a function $f$: $\mathbb{R}^d \rightarrow$ **{0, 1}**, given **m** samples **(x, $f$(x))** where **x ~ $\mathcal{D}$**, find a hypothesis $h$ s.t. $\text{Pr}_{x \sim \mathcal{D}}[h(\textbf{x}) \neq f(\textbf{x})] \leq \boldsymbol{\varepsilon}$ w.p. **1 - δ**. Efficient if **m ≤ O(poly(d,1/ε, log(1/δ)))**.

❖ **Learning from Label Proportions (LLP)**: Samples are "bags" of the form **({$x_1$, ..., $x_q$}, k)** where $\boldsymbol{\Sigma_i f(\textbf{x}_i) = \textbf{k}}$. Our study: $f \leftarrow$ Linear Threshold function (LTF) **[ $\mathbb{1}\{\textbf{r}^\textbf{T}\textbf{x} + \textbf{c} > \textbf{0}\}$ ]**

# Problem Statement

❖ **PAC Learning**: For a function $f$: $\mathbb{R}^d \rightarrow$ **{0, 1}**, given **m** samples **(x, $f$(x))** where **x ~ $\mathcal{D}$**, find a hypothesis $h$ s.t. **$\Pr_{x\sim\mathcal{D}}[h(x)\neq f(x)] \leq \varepsilon$** w.p. **1 - $\delta$**. Efficient if **m ≤ O(poly(d,1/ε, log(1/δ)))**.

❖ **Learning from Label Proportions (LLP)**: Samples are "bags" of the form **({$x_1$, ..., $x_q$}, k)** where **$\Sigma_i f(x_i)$ = k**.  Our study: $f \leftarrow$ Linear Threshold function (LTF) **[ $\mathbb{1}\{r^Tx + c > 0\}$ ]**

➢ LTFs are efficiently PAC learnable to arbitrary accuracy

# Problem Statement

❖ **PAC Learning**: For a function $f: \mathbb{R}^d \rightarrow \{0, 1\}$, given **m** samples **(x,** $f$**(x))** where **x ~ $\mathcal{D}$**, find a hypothesis $h$ s.t. $\mathbf{Pr_{x \sim \mathcal{D}}[h(x) \neq f(x)] \leq \varepsilon}$ w.p. **1 - δ**. Efficient if **m ≤ O(poly(d,1/ε, log(1/δ)))**.

❖ **Learning from Label Proportions (LLP)**: Samples are "bags" of the form **({$x_1$, ..., $x_q$}, k)** where $\mathbf{\Sigma_i f(x_i) = k}$. Our study: $f \leftarrow$ Linear Threshold function (LTF) **[ $\mathbb{1}\{r^Tx + c > 0\}$ ]**
  ➤ LTFs are efficiently PAC learnable to arbitrary accuracy

❖ **LLP Hardness [Saket 21, 22]**: Given a set of bags of size ≤ **q**, s.t. ∃ LTF consistent with all bags, NP-hard to find <u>**any**</u> LTF consistent with ≥ **(1/q + o(1))**-fraction of the bags.

# Problem Statement

❖ **PAC Learning**: For a function $f: \mathbb{R}^d \rightarrow$ **{0, 1}**, given **m** samples **(x, $f$(x))** where **x ~ $\mathcal{D}$**, find a hypothesis $h$ s.t. $\mathbf{Pr_{x \sim \mathcal{D}}[h(x) \neq f(x)] \leq \varepsilon}$ w.p. **1 - δ**. Efficient if **m ≤ O(poly(d,1/ε, log(1/δ)))**.

❖ **Learning from Label Proportions (LLP)**: Samples are "bags" of the form **({$x_1$, …, $x_q$}, k)** where $\mathbf{\Sigma_i f(x_i) = k}$. Our study: $f \leftarrow$ Linear Threshold function (LTF) **[ $\mathbb{1}\{r^T x + c > 0\}$ ]**
  ➤ LTFs are efficiently PAC learnable to arbitrary accuracy

❖ **LLP Hardness [Saket 21, 22]**: Given a set of bags of size ≤ **q**, s.t. ∃ LTF consistent with all bags, NP-hard to find <u>**any**</u> LTF consistent with ≥ **(1/q + o(1))**-fraction of the bags.

  **Question:** *What happens for natural/well-behaved distributions?*

# Problem Statement

❖ **PAC Learning**: For a function $f: \mathbb{R}^d \rightarrow \{0, 1\}$, given **m** samples $(x, f(x))$ where **x ~ 𝒟**, find a hypothesis $h$ s.t. $\Pr_{x\sim\mathcal{D}}[h(x)\neq f(x)] \leq \varepsilon$ w.p. **1 - δ**. Efficient if **m ≤ O(poly(d,1/ε, log(1/δ)))**.

❖ **Learning from Label Proportions (LLP)**: Samples are "bags" of the form $(\{x_1, ..., x_q\}, k)$ where $\Sigma_i f(x_i) = k$. Our study: $f \leftarrow$ Linear Threshold function (LTF) $[\ \mathbb{1}\{r^Tx + c > 0\}\ ]$
  ➢ LTFs are efficiently PAC learnable to arbitrary accuracy

❖ **LLP Hardness [Saket 21, 22]**: Given a set of bags of size ≤ **q**, s.t. ∃ LTF consistent with all bags, NP-hard to find **any** LTF consistent with ≥ **(1/q + o(1))**-fraction of the bags.

  **Question:** *What happens for natural/well-behaved distributions?*

❖ Bag Oracle for LTF $f$, **𝒟 = N(μ, Σ)** and fixed **q, k** :- **Ex($f$,𝒟, q, k)**
  ➢ Samples bag with **k** feature-vecs. from **𝒟|f(x)=1** and **q-k** from **𝒟|f(x)=0**.

# Our Results

| Homogeneous LTF Standard Gaussian | Homogeneous LTF Centered Gaussian | Non-Homogeneous LTF General Gaussian |
|---|---|---|
| $k \neq q/2$ | $\forall\, k \in \{1, \ldots, q-1\}$ | $\forall\, k \in \{1, \ldots, q-1\}$ |
| $\mathcal{D} := N(\mathbf{0}, \mathbf{I})$ | $\mathcal{D} := N(\mathbf{0}, \boldsymbol{\Sigma})$ | $\mathcal{D} := N(\mu, \boldsymbol{\Sigma})$ |
| $f(\mathbf{x}) := \mathbf{1}\{\mathbf{r}_*^\mathsf{T}\mathbf{x} > 0\}$ | $f(\mathbf{x}) := \mathbf{1}\{\mathbf{r}_*^\mathsf{T}\mathbf{x} > 0\}$ | $f(\mathbf{x}) := \mathbf{1}\{\mathbf{r}_*^\mathsf{T}\mathbf{x} + c_* > 0\}$ |
| $\hat{f}(\mathbf{x}) := \mathbf{1}\{\hat{\mathbf{r}}^\mathsf{T}\mathbf{x} > 0\}$ | $\hat{f}(\mathbf{x}) := \mathbf{1}\{\hat{\mathbf{r}}^\mathsf{T}\mathbf{x} > 0\}$ | $\hat{f}(\mathbf{x}) := \mathbf{1}\{\hat{\mathbf{r}}^\mathsf{T}\mathbf{x} + \hat{c} > 0\}$ |
| $O\left(\frac{d}{\varepsilon^2}\log\left(\frac{d}{\delta}\right)\right)$ | $O\left(\frac{d}{\varepsilon^4}\log\left(\frac{d}{\delta}\right)\left(\frac{\lambda_{\max}}{\lambda_{\min}}\right)^6 q^8\right)$ | $O\left(\frac{d}{\varepsilon^4}\log\left(\frac{d}{\delta}\right)\frac{O(\ell^2)}{\Phi(\ell)(1-\Phi(\ell))}\left(\frac{\lambda_{\max}}{\lambda_{\min}}\right)^4\left(\frac{\sqrt{\lambda_{\max}}+\|\mu\|_2}{\sqrt{\lambda_{\min}}}\right)^4 q^8\right)$ |

$$d \leftarrow \text{ dimension of the feature-vectors}, \quad \ell := -(c_* + \mathbf{r}_*^\mathsf{T}\mu)/\|\boldsymbol{\Sigma}^{1/2}\mathbf{r}_*\|_2$$

$$\lambda_{\max} \leftarrow \text{ maximum eigenvalue of } \boldsymbol{\Sigma}, \quad \lambda_{\min} \leftarrow \text{ minimum eigenvalue of } \boldsymbol{\Sigma}$$
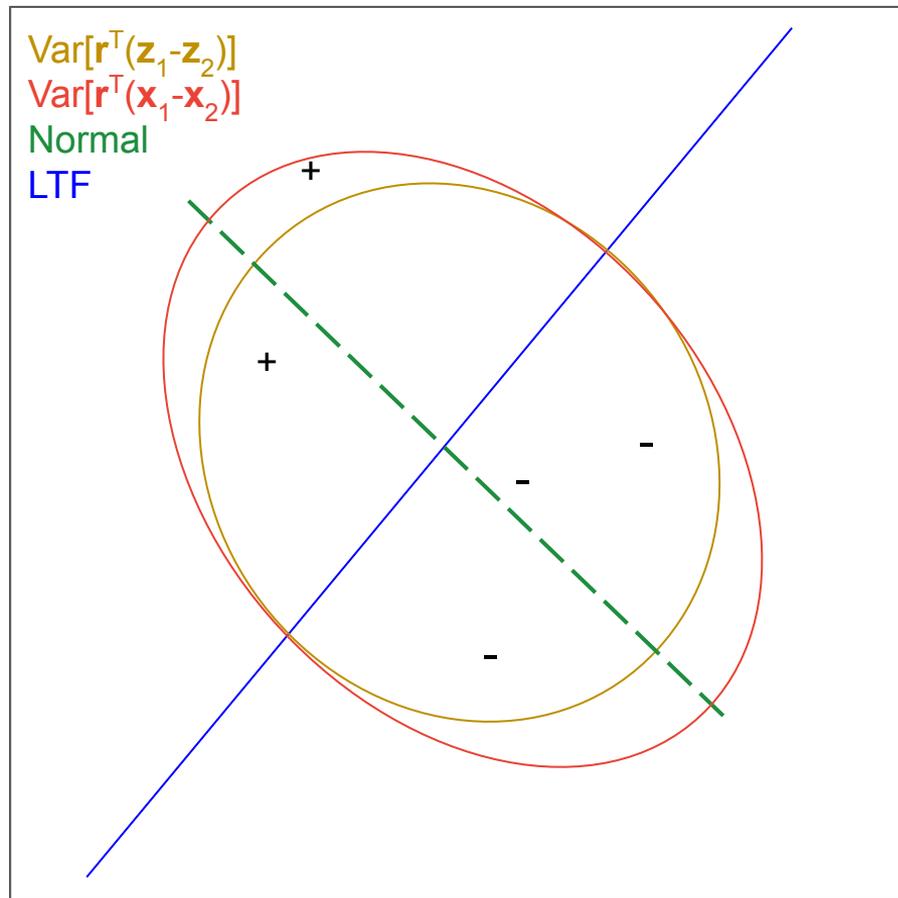
# Normal Estimation

❖ **Observation**: Sampling a pair of feature vectors from a bag

  ➤ $(\mathbf{x}_1, \mathbf{x}_2)$ independently u.a.r:
  $\Pr[f(\mathbf{x}_1) \neq f(\mathbf{x}_1)] = 2k(q\text{-}k)/q^2$

  ➤ $(\mathbf{z}_1, \mathbf{z}_2)$ pair u.a.r w/o replacement:
  $\Pr[f(\mathbf{z}_1) \neq f(\mathbf{z}_2)] = 2k(q\text{-}k)/q(q\text{-}1)$

**Theorem:** $\rho(\mathbf{r}) := \mathsf{Var}\left[\mathbf{r}^{\mathsf{T}}(\mathbf{z}_1 - \mathbf{z}_2)\right] / \mathsf{Var}\left[\mathbf{r}^{\mathsf{T}}(\mathbf{x}_1 - \mathbf{x}_2)\right]$
is maximized when $\mathbf{r} = \pm \mathbf{r}_*$.

**Proof Sketch:** Case $\mathcal{D} = N(\mathbf{0}, \mathbf{I})$ and $k = q/2$. For any $\mathbf{r} \in \mathbb{S}^{d-1}$

$$\rho(\mathbf{r}) = \begin{cases} 1 & \text{if } \mathbf{r}^{\mathsf{T}}\mathbf{r}_* = 0, \text{ i.e. } \mathbf{r} \text{ lies on the LTF,} \\ 1 + \frac{1}{q-1}\left(\frac{2}{\pi}\right) & \text{if } \mathbf{r} = \pm\mathbf{r}_*, \text{ i.e. } \mathbf{r} \text{ is aligned} \\ & \text{with the normal to the LTF,} \\ 1 + \frac{1}{q-1}\left(\frac{2}{\pi}\right)\cos^2\theta & \text{if } \mathbf{r}^{\mathsf{T}}\mathbf{r}_* = \cos\theta, \text{ i.e. the angle} \\ & \text{between } \mathbf{r} \text{ and } \mathbf{r}_* \text{ is } \theta. \end{cases}$$



$\mathsf{Var}[\mathbf{r}^{\mathsf{T}}(\mathbf{z}_1\text{-}\mathbf{z}_2)]$
$\mathsf{Var}[\mathbf{r}^{\mathsf{T}}(\mathbf{x}_1\text{-}\mathbf{x}_2)]$
Normal
LTF

# Homogeneous LTF with $N(0, \Sigma)$

❖ Let $\Sigma_B = E[(x_1 - x_2)(x_1 - x_2)^\top]$ and $\Sigma_D = E[(z_1 - z_2)(z_1 - z_2)^\top]$

❖ **Objective**: $\text{argmax}_{||r||=1} r^\top \Sigma_B r / r^\top \Sigma_D r = \Sigma_B^{-\frac{1}{2}} \textbf{PrincipalEigenVector}(\Sigma_B^{-\frac{1}{2}} \Sigma_D \Sigma_B^{-\frac{1}{2}})$

# Homogeneous LTF with $\mathbf{N(0, \Sigma)}$

❖ Let $\mathbf{\Sigma_B = E[(x_1 - x_2)(x_1 - x_2)^\top]}$ and $\mathbf{\Sigma_D = E[(z_1 - z_2)(z_1 - z_2)^\top]}$

❖ **Objective**: $\mathrm{argmax}_{||r||=1} \mathbf{r^\top \Sigma_B r / r^\top \Sigma_D r}$ = $\mathbf{\Sigma_B^{-\frac{1}{2}} PrincipalEigenVector(\Sigma_B^{-\frac{1}{2}} \Sigma_D \Sigma_B^{-\frac{1}{2}})}$

❖ **Stability Theorem**: The ratio maximization computed with high probability approximations of $\mathbf{\Sigma_B}$ and $\mathbf{\Sigma_D}$ is close to the normal with high probability.

# Homogeneous LTF with $N(0, \Sigma)$

❖ Let $\Sigma_B = E[(x_1-x_2)(x_1-x_2)^T]$ and $\Sigma_D = E[(z_1-z_2)(z_1-z_2)^T]$

❖ **Objective**: $\text{argmax}_{||r||=1} r^T \Sigma_B r / r^T \Sigma_D r \ = \ \Sigma_B^{-\frac{1}{2}}\textbf{PrincipalEigenVector}(\Sigma_B^{-\frac{1}{2}}\Sigma_D\Sigma_B^{-\frac{1}{2}})$

❖ **Stability Theorem**: The ratio maximization computed with high probability approximations of $\Sigma_B$ and $\Sigma_D$ is close to the normal with high probability.

❖ Geometric bound → Bound on sample error

➢ An algorithm to find a high probability estimator of $c_*$ given a high probability estimator of $r_*$

# Homogeneous LTF with $N(0, \Sigma)$

❖ Let $\Sigma_B = E[(x_1-x_2)(x_1-x_2)^\top]$ and $\Sigma_D = E[(z_1-z_2)(z_1-z_2)^\top]$

❖ **Objective**: $\text{argmax}_{||r||=1} r^\top\Sigma_B r/r^\top\Sigma_D r = \Sigma_B^{-\frac{1}{2}}\textbf{PrincipalEigenVector}(\Sigma_B^{-\frac{1}{2}}\Sigma_D\Sigma_B^{-\frac{1}{2}})$

❖ **Stability Theorem**: The ratio maximization computed with high probability approximations of $\Sigma_B$ and $\Sigma_D$ is close to the normal with high probability.

❖ Geometric bound → Bound on sample error

➢ An algorithm to find a high probability estimator of $c_*$ given a high probability estimator of $r_*$

❖ Bound on sample error → Bound on Generalization error

➢ **Generalization Error Bound**

# Homogeneous LTF with $\mathbf{N(0, \Sigma)}$

❖ Let $\mathbf{\Sigma_B = E[(x_1-x_2)(x_1-x_2)^T]}$ and $\mathbf{\Sigma_D = E[(z_1-z_2)(z_1-z_2)^T]}$

❖ **Objective**: $\text{argmax}_{||\mathbf{r}||=1}\mathbf{r^T\Sigma_B r/r^T\Sigma_D r}$ $=$ $\mathbf{\Sigma_B^{-\frac{1}{2}}PrincipalEigenVector(\Sigma_B^{-\frac{1}{2}}\Sigma_D\Sigma_B^{-\frac{1}{2}})}$

❖ **Stability Theorem**: The ratio maximization computed with high probability approximations of $\mathbf{\Sigma_B}$ and $\mathbf{\Sigma_D}$ is close to the normal with high probability.

❖ Geometric bound → Bound on sample error

   ➢ An algorithm to find a high probability estimator of $c_*$ given a high probability estimator of $\mathbf{r_*}$

❖ Bound on sample error → Bound on Generalization error

   ➢ **Generalization Error Bound**

❖ Sub-gaussian concentration bounds for thresholded Gaussians.

Thank You