

Goal

Establishing a strong notion of impossibility of learning

- ▶ Binary supervised classification : $\mathcal{T}_n = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$.
- ▶ Class conditional densities are Gaussian : $\mathcal{N}(\mu_{\pm}, \Sigma)$.
- ▶ High Dimensional Setting : $d/n \rightarrow \infty$ as $n \rightarrow \infty$.
- ▶ The optimal Bayes error does not scale with n and d .
- ▶ Can we lower bound worst case probability of error for any classifier asymptotically?

Preliminaries

- ▶ Training set \mathcal{T}_n , with n iid samples from $p(\mathbf{x}, y|\theta)$.
- ▶ Learning rule $\hat{y}_{\mathcal{T}_n} : \mathbb{R}^d \rightarrow \{-1, 1\}$.
- ▶ Conditional Probability of error : $P_{e|\mathcal{T}_n, \theta} = \Pr_{(\mathbf{x}, y)}(\hat{y}_{\mathcal{T}_n}(\mathbf{x}) \neq y | \mathcal{T}_n, \theta)$.
- ▶ Probability of error : $P_{e|\theta} = \mathbb{E}_{\mathcal{T}_n}(P_{e|\mathcal{T}_n, \theta})$.
- ▶ Bayes optimal error : $P_{e|\theta}^* = \min_{\hat{y}} P_{e|\theta}$.

- ▶ Set of feasible parameters to keep $P_{e|\theta}^*$ fixed :

$$\Theta(\alpha) = \{\theta : P_{e|\theta}^* = Q(\alpha/2)\}$$

- ▶ Worst case probability of error :

$$P_e(n, d, \Theta(\alpha), \hat{y}) = \sup_{\theta \in \Theta(\alpha)} P_{e|\theta}$$

- ▶ Goal : Evaluation of $\liminf_{(d, n/d) \rightarrow (\infty, 0)} P_e(n, d, \Theta(\alpha), \hat{y}_{\mathcal{T}_n})$.

- ▶ It is generally known that if the number of samples n is much less than the VC-dimension V , the learning is impossible in the *distribution free* setting.

Theorem (Anthony and Bartlett, 1999)

If $P_e^* = c$ and $n < V/(320\epsilon^2)$, then for all supervised learning rules $\hat{y}_{\mathcal{T}_n}$, and all $\epsilon, \delta \in (0, 1/64)$,

$$\sup_{p(\mathbf{x}, y) : P_e^* = c} \Pr(P_{e|\mathcal{T}_n, \theta} - c \geq \epsilon) \geq \delta$$

Major Differences in our Results

- ▶ The set of distributions is not as large as all distributions and is not as small as only a single distribution.
- ▶ The worst case is evaluated only on a subset of Gaussian distributions in our results.

Feasible Sets for Gaussian Distribution

- ▶ $\theta = (\mu_+, \mu_-, \Sigma)$ and $P_{e|\theta}^* = Q\left(\frac{1}{2}\|\Sigma^{-\frac{1}{2}}(\mu_+ - \mu_-)\|_2\right)$
 - ▶ $\Theta_0(\alpha) = \{(\mu_+, \mu_-, \Sigma) : \|\Sigma^{-\frac{1}{2}}(\mu_+ - \mu_-)\|_2 = \alpha\}$
 - ▶ A canonical subset of $\Theta_0(\alpha)$, which is easier to analyze
- $$\Theta_{\text{Sphere}}(\alpha) := \{(\mathbf{h}, -\mathbf{h}, \beta^2 \mathbf{I}) : \|\mathbf{h}\|_2 = 1, \beta = 2/\alpha\}.$$

Theorem (Main Result)

For any sequence of classifiers $\hat{y}_{\mathcal{T}_n}$, and $\alpha \geq 0$, we have

$$\liminf_{(d, n/d) \rightarrow (\infty, 0)} P_e(n, d, \Theta_{\text{Sphere}}(\alpha), \hat{y}_{\mathcal{T}_n}) \geq \frac{1}{2}$$

Discussion

- ▶ Even in the easier scenario of Θ_{Sphere} , there is no hope to get meaningful classification.
- ▶ Trivial to show that the same result holds for $\Theta_0(\alpha) \supseteq \Theta_{\text{Sphere}}(\alpha)$.
- ▶ It is consistent with the previous results :
 - ▶ If ML estimation of μ_{\pm} and Σ are plugged into the optimal Bayes classifier, error is asymptotically not less than half [Bickel and Levina, 2004].
 - ▶ Suppose that $\Theta_1(\alpha) = \{(\mathbf{h}, -\mathbf{h}, \mathbf{h}\mathbf{h}^\top + \sigma^2 \mathbf{I}) : P_{e|\mathbf{h}}^* = Q(\alpha/2)\}$. If ML estimation of \mathbf{h} is plugged into the optimal Bayes classifier, the performance is no better than a random coin toss asymptotically [Orten et al, 2011].

Necessary Condition on Θ for Learnability

- ▶ $\Theta_{\text{subset}} := \{(\mathbf{h}, -\mathbf{h}, \beta^2 \mathbf{I}) \in \Theta_{\text{Sphere}}, \mathbf{h} \in \mathcal{H} \subseteq \mathcal{S}^{d-1}\}$.
- ▶ Let $\text{vol}(\mathcal{H}) \triangleq \Pr_{H \sim U(\mathcal{S}^{d-1})}(H \in \mathcal{H})$.

Corollary

Suppose that $\lim_{d \rightarrow \infty} \text{vol}(\mathcal{H})$ exists. If for a sequence of classifiers $\hat{y}_{\mathcal{T}_n}$,

$$\limsup_{(d, n/d) \rightarrow (\infty, 0)} P_e(n, d, \Theta_{\text{Sphere}}, \hat{y}_{\mathcal{T}_n}) = \frac{1}{2}$$

and

$$\limsup_{(d, n/d) \rightarrow (\infty, 0)} P_e(n, d, \Theta_{\text{subset}}, \hat{y}_{\mathcal{T}_n}) < \frac{1}{2}$$

then

$$\lim_{d \rightarrow \infty} \text{vol}(\mathcal{H}) = 0.$$

Discussion

- ▶ Consistent classifiers exist if \mathbf{h} belongs to either of these sparsity classes [Orten et al, 2011]:
 - ▶ Assume that sorted absolute values of components of \mathbf{h} ($h_{(1)}, \dots, h_{(d)}$) decay exponentially or polynomially fast :

$$\mathcal{H}_{\text{exp}} = \{\mathbf{h} : |h_{(k)}| = M_1(d)\alpha^k, 0 < \alpha < 1\}$$

$$\mathcal{H}_{\text{poly}} = \{\mathbf{h} : |h_{(k)}| = M_2(d)k^{-\beta}, \beta > 0.5\}$$

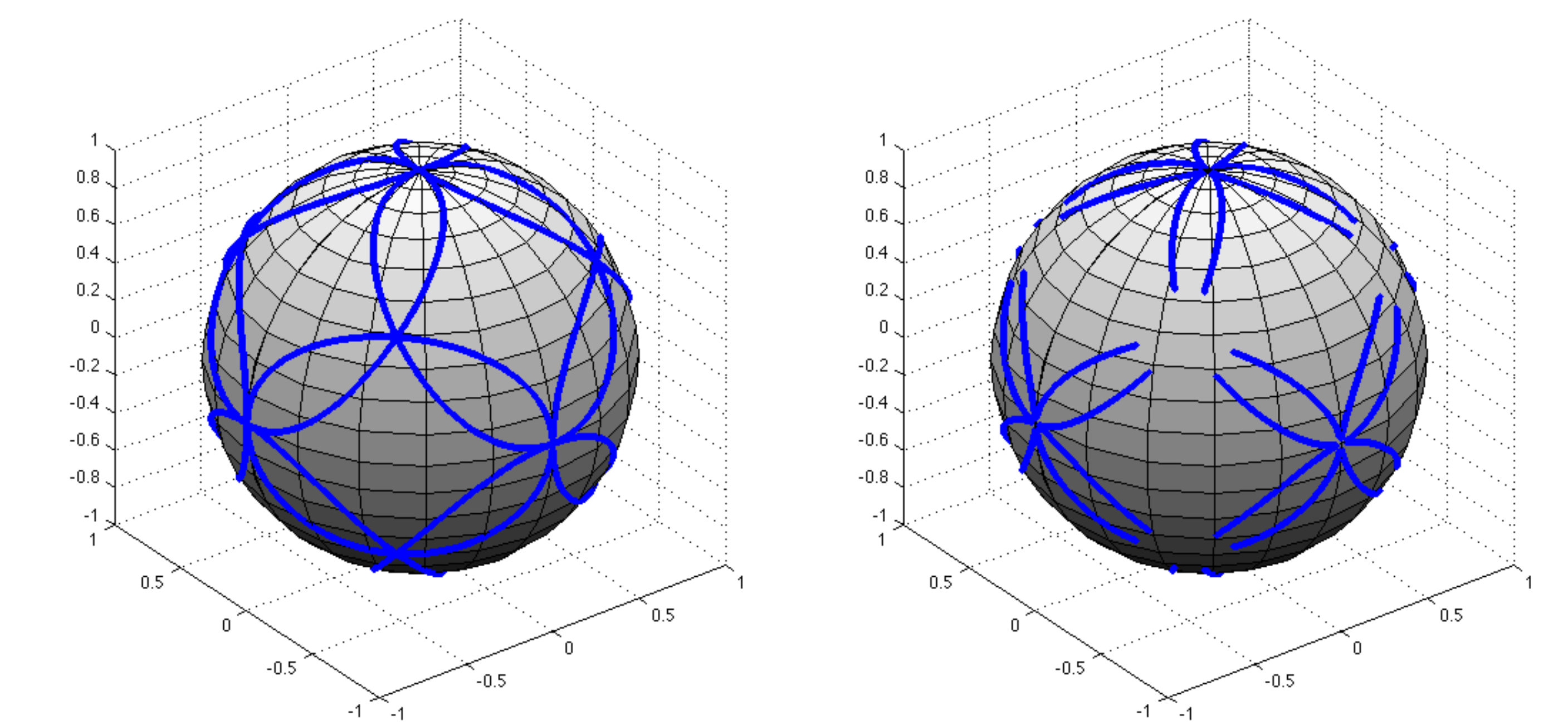


Figure: Left : \mathcal{H}_{exp} - Right : $\mathcal{H}_{\text{poly}}$ for $d = 3$.