

CSE 6740 A/ISyE 6740: Computational Data Analysis: Introductory lecture

Nisha Chandramoorthy

September 5, 2023

Last time

- ▶ Recap of shrinkage by ridge regression

Last time

- ▶ Recap of shrinkage by ridge regression
- ▶ Geometric view of compression by LASSO

Last time

- ▶ Recap of shrinkage by ridge regression
- ▶ Geometric view of compression by LASSO
- ▶ Generalization of LASSO

Last time

- ▶ Recap of shrinkage by ridge regression
- ▶ Geometric view of compression by LASSO
- ▶ Generalization of LASSO
- ▶ ℓ^0 regularization and compressed sensing

Today

- ▶ Classification ERM

Today

- ▶ Classification ERM
- ▶ Logistic regression, Bayesian view

Today

- ▶ Classification ERM
- ▶ Logistic regression, Bayesian view
- ▶ Perceptron algorithm and convergence proof

Today

- ▶ Classification ERM
- ▶ Logistic regression, Bayesian view
- ▶ Perceptron algorithm and convergence proof
- ▶ Support vector regression or maximum margin classification

Today

- ▶ Classification ERM
- ▶ Logistic regression, Bayesian view
- ▶ Perceptron algorithm and convergence proof
- ▶ Support vector regression or maximum margin classification
- ▶ Convex optimization

Today

- ▶ Classification ERM
- ▶ Logistic regression, Bayesian view
- ▶ Perceptron algorithm and convergence proof
- ▶ Support vector regression or maximum margin classification
- ▶ Convex optimization
- ▶ Bias-complexity tradeoff

Linear predictors

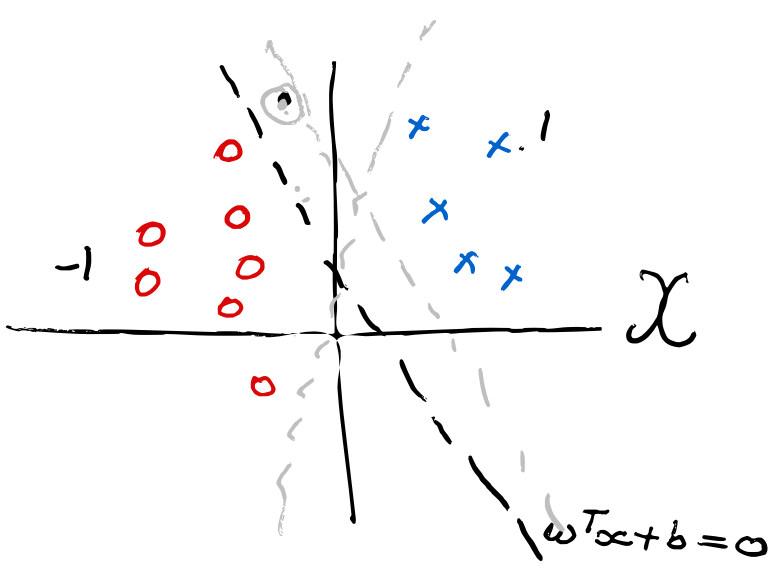
Classification ERM

$$\mathcal{H} = \{ h^{lin}(\cdot; w, b) : w \in \mathbb{R}^d, b \in \mathbb{R} \}$$

$$\text{sgn}(x) = \begin{cases} 1 & x > 0 \\ -1 & x < 0 \end{cases}$$

$\mathcal{HS} =$

$$\{ \text{sgn} \circ h^{lin}(\cdot, w, b) : w \in \mathbb{R}^d, b \in \mathbb{R} \}$$



$$R(h) = \mathbb{E}_{z \sim D} \mathbb{1}_{\{(x,y) : h(x) \neq y\}}$$

ERM

$$S: \{(x_i, y_i) : 1 \leq i \leq m\}$$

iid from \mathcal{D} .

$$R_S(h) = \frac{1}{m} \sum_{(x,y) \in S} \mathbb{1}_{\{(x,y) \in S : h(x) \neq y\}}$$

$$h(\cdot, w, b) \in \mathcal{HS}$$

Realizability: $\exists h(\cdot, w, b) \in \mathcal{HS}$

s.t. $h(x, w, b) = y$ for almost every $(x, y) \sim \mathcal{D}$.

$$h^{ERM} = \underset{h \in \mathcal{HS}}{\text{argmin}} \frac{1}{m} \sum_{(x,y) \in S} \mathbb{1}_{\{(x,y) \in S : h(x) \neq y\}}$$

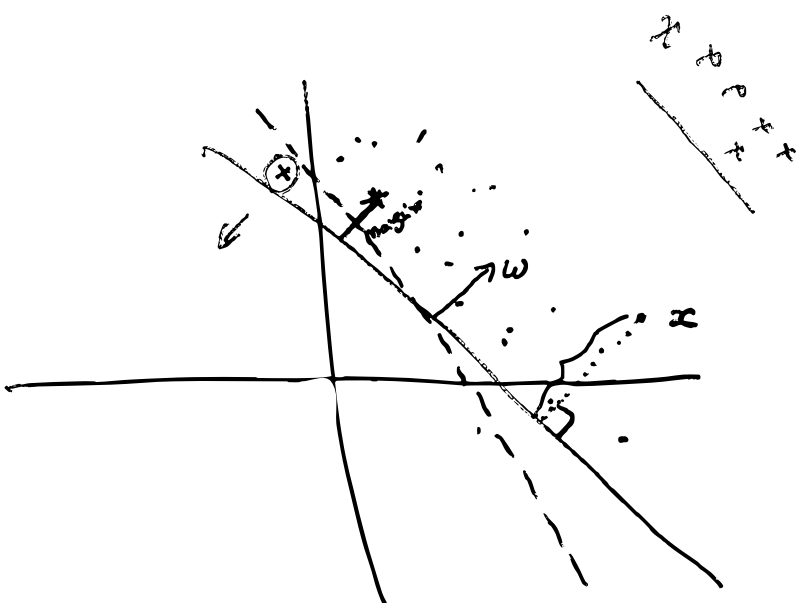
any h^{ERM} will be s.t.

$$h^{ERM}(x_i) = y_i \quad \forall i \in [m] \quad (1, 2, \dots, m)$$

Margin

Given a sample set S and a classifier $h \in \mathcal{HS}$,

$$\text{Margin}(h, S) = \min_{(x,y) \in S} d_h(x, y)$$



$$\max_{w, b} \text{Margin}(h(\cdot, w, b), S)$$

subject to

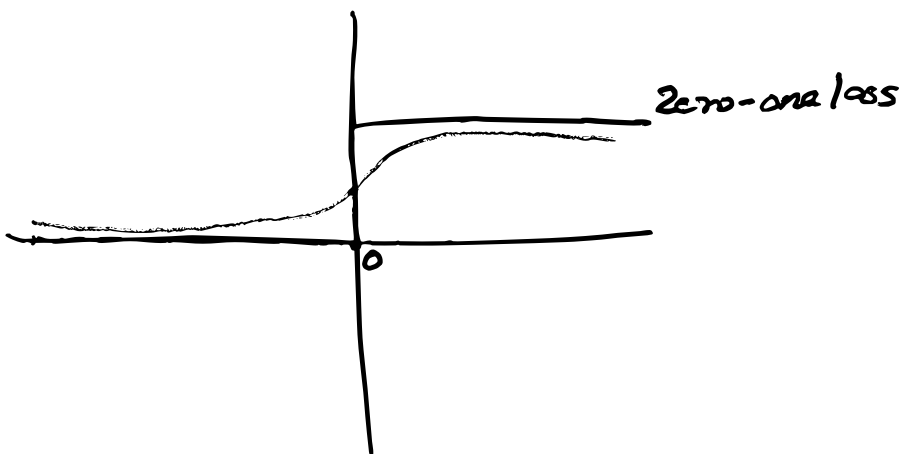
$$\text{SVM} \quad \rightarrow h(x_i, w, b) = y_i \quad \forall i \in [m]$$

logistic regression

$$p(x) = \frac{1}{1 + e^{-x}}$$

ERM :

$$\operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i \langle w, x_i \rangle})$$



Likelihood function:

$$P(x, y | w, b) = \frac{1}{1 + e^{-y(\langle w, x \rangle + b)}}$$

Bayes rule

$$P(\underline{w, b} | S) \propto P(S) \underline{P(S | w, b)}$$

MLE :

$$\underline{P(S | w, b)} = \prod_{i=1}^m \frac{1}{1 + e^{-y_i(\langle w, x_i \rangle + b)}}$$

$$\log P(S | w, b) = - \sum_{i=1}^m \log(1 + e^{-y_i(\langle w, x_i \rangle + b)})$$

$$\max_{w, b} \underline{\log \text{likelihood}}(w, b, S)$$

Convex optimization

Halfspaces

Linear programming:

$$\max_{w \in \mathbb{R}^d} \langle u, w \rangle$$

$$\text{subject to } Aw \geq v$$

$$y_i \langle w, x_i \rangle > 0 \quad \forall i \in [m]$$

$$y_i (\langle w, x_i \rangle + b) > 0$$

Recall realizability. $\exists (w^*, b^*)$

Define $S := \min_i y_i (\langle w^*, x_i \rangle + b^*)$

$$\frac{y_i (\langle w^*, x_i \rangle + b^*)}{S} \geq 1$$

$$y_i (\langle w, x_i \rangle + b) \geq 1$$

Perceptron algorithm

$$\omega^{(t+1)} = \omega^{(t)} + y_i x_i, \text{ where}$$

i is such that x_i is mis-classified by $p_{\omega^{(t)}}$

$$h(x_i, w, b) = y_i \quad \forall i \in [m]$$

$$\omega^{(0)} = 0 \in \mathbb{R}^{d+1}$$

$$\omega^{(t+1)} = \omega^{(t)} + y_i x_i$$

$$\omega^{(T)}$$

Convergence proof

Assumptions: \rightarrow realizability

$$\text{Given: } B = \min \{ \|\omega\| \cdot \forall i \in [m] \text{ s.t. } y_i \langle \omega, x_i \rangle \geq 1 \}$$

$$\rightarrow R = \max_i \|x_i\|$$

Thm: The perceptron algorithm stops after at most $(RB)^2$ steps

Given:

$$\|\omega^*\| = B$$

$$\omega^{(1)} = 0$$

$$\omega^{(t+1)} = \omega^{(t)} + y_i x_i$$

$$\omega^{(t+1)} - \omega^{(t)} = y_i x_i$$

$$\underbrace{y_i \langle \omega^{(t+1)}, x_i \rangle}_{\text{"}} > y_i \langle \omega^{(t)}, x_i \rangle$$

$$y_i \langle \omega^{(t)} + x_i y_i, x_i \rangle$$

$$\underbrace{y_i \langle \omega^{(t)}, x_i \rangle}_{\text{"}} + \underline{\|x_i\|^2}$$

Proof:

$$\omega^{(t+1)} = \omega^{(t)} + y_i x_i$$

$$\|\omega^{(t+1)}\|^2 = \|\omega^{(t)}\|^2 + \|x_i\|^2 + 2y_i \langle \omega^{(t)}, x_i \rangle$$

$$\text{Since } y_i \langle \omega^{(t)}, x_i \rangle \leq 0$$

$$\|\omega^{(t+1)}\|^2 \leq \|\omega^{(t)}\|^2 + \|x_i\|^2$$

$$\leq \|\omega^{(t)}\|^2 + R^2$$

Telescoping,

$$\|\omega^{(T+1)}\|^2 \leq \|\omega^{(1)}\|^2 + R^2$$

$$\leq \|\omega^{(T-1)}\|^2 + 2R^2$$

\vdots

$$\textcircled{+} \leq TR^2 \quad (\omega^{(1)} = 0)$$

$$\begin{aligned} \langle \omega^{(t+1)}, \omega^* \rangle &= \langle \omega^{(t)} + y_i x_i, \omega^* \rangle \\ &= \langle \omega^{(t)}, \omega^* \rangle + y_i \langle x_i, \omega^* \rangle \end{aligned}$$

$$\Rightarrow \langle \omega^{(t+1)} - \omega^{(t)}, \omega^* \rangle = y_i \langle x_i, \omega^* \rangle$$

$$(\because y_i \langle x_i, \omega^* \rangle \geq 1 \text{ since } \omega^* \text{ is a correct classifier})$$

$$\Rightarrow \langle \omega^{(T+1)}, \omega^* \rangle = \sum_{t=1}^T \langle \omega^{(t+1)} - \omega^{(t)}, \omega^* \rangle$$

$$(\because \omega^{(1)} = 0)$$

$$\geq T \quad \textcircled{*}$$

Using Cauchy-Schwarz,

$$|\langle \omega^{(T+1)}, \omega^* \rangle| \leq \|\omega^{(T+1)}\| \|\omega^*\|$$

$$\leq \|\omega^{(T+1)}\| B$$

$\Rightarrow \textcircled{*}$

$$T \leq \|\omega^{(T+1)}\| B$$

Applying $\textcircled{+}$,

$$T \leq \|\omega^{(T+1)}\| B$$

$$\leq \sqrt{T} R B$$

$$\Rightarrow T^2 \leq R B.$$