

The Framework of Statistical Learning

- ▷ $X \subset \mathbb{R}^k$: domain set (e.g. space of images with fixed nr. of pixels)
- ▷ $Y \subset \mathbb{R}^q$: target set / set of label (e.g. $Y = \{0, 1\}$ binary classification)
- ▷ Let \mathcal{D} be a probability distribution on $X \times Y$ $Y = \mathbb{R}^q$: regression
- ▷ Assume we are given training set $S = \{(x_i, y_i)\}_{i=1}^n$ i.i.d \mathcal{D}
- ▷ Goal: Find a predictor function $h: X \rightarrow Y$ that minimizes the expected risk $L_{\mathcal{D}}(h) := \mathbb{E}_{\mathcal{D}}[l(Y, h(X))]$ where $l: Y \times Y \rightarrow \mathbb{Z}$ is a given loss/error function

Learning Algorithm:

Specify algorithm that maps S to a specific member function $\hat{h}_n \in \mathcal{F}$ of a hypothesis space $\mathcal{F} \subset \{h: X \rightarrow Y\}$ based on the information of S .

In many cases: $\hat{h}_n = \underset{h \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n l(y_i, h(x_i))$

empirical risk

Empirical Risk Minimization

Ex: \triangleright Linear Regression: $\triangleright l(y, z) = \frac{1}{2}(y - z)^2$
 $\triangleright \mathcal{F} = \left\{ x \mapsto \beta_0 + \langle x, \beta \rangle, \beta_0 \in \mathbb{R}, \beta \in \mathbb{R}^K \right\}$

Bias - Variance Tradeoff:

$$L_2(\hat{h}_n) - \min_h L_2(h) = \underbrace{(L_2(\hat{h}_n) - L_2(h_F))}_{\text{generalization error}} + \underbrace{(L_2(h_F) - \min_h L_2(h))}_{\text{approximation error}}$$

depends on training set S

$h_F = \underset{h \in \mathcal{F}}{\operatorname{argmin}} L_D(h)$

\mathcal{F} hypothesis space

\triangleright smaller if \mathcal{F} is large

$$\hat{h}_n = \underset{h \in \mathcal{F}}{\operatorname{argmin}} \left(\underbrace{\frac{1}{n} \sum_{i=1}^n l(y_i, h(x_i))}_{\text{empirical risk}} + \lambda \cdot R(h) \right)$$

$\lambda > 0$

{ a term that quantifies "complexity" of $h \in \mathcal{F}$ }

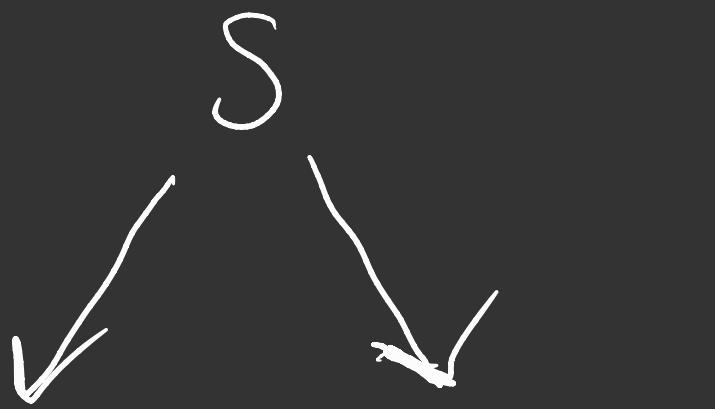
"Regularized loss minimization"
(RLM)

Ridge Regression: Linear Regression with $R(h) = \|\beta\|_{l_2}^2$

$$\hat{\beta}_n = \underset{\beta \in \mathbb{R}^K}{\operatorname{argmin}} \left\{ \frac{1}{n} \|X\beta - y\|_2^2 + \lambda \|\beta\|_2^2 \right\}$$

- ▷ If $\lambda = 0$: \rightarrow Linear regression
- ▷ If $\lambda \rightarrow +\infty$: coefficients $\hat{\beta}_i$ "shrinked to 0"
- ▷ In between: balances fit of linear model to S & size of coefficient vector β

Cross validation



80%

S_{train}

S_{validation}

20%

k-fold Cross Validation

Preprocessing:

If domain set s.t. $X \subset \mathbb{R}^k$, we can define a feature map $\phi: X \rightarrow \tilde{X} \subset \mathbb{R}^\ell$ (often with $k < \ell$) such that $S = (\phi(x_i), y_i)_{i=1}^n$ is used as training set.

Ex: Polynomial features. E.g., if $k=1$, $\ell=5$:

$$\phi(x) = (x, x^2, x^3, x^4, x^5).$$

→ often improves expressive power of a learning model!

Sparse Regression

$$= \sum_{i=1}^n I_{\{\beta_i \neq 0\}} \{\beta\}$$

$$\mathcal{F}_k^{\text{sparse}} = \left\{ h: \mathbb{R}^d \rightarrow \mathbb{R} : h(x) = \beta_0 + \langle \beta, x \rangle; \text{ s.t. } \|\beta\|_0 \leq k \right\}$$

Idea: Only $k \ll d$ coefficients of linear model will suffice to explain/predict the outcome $k \ll d$

Problem: ERM on $\mathcal{F}_k^{\text{sparse}}$ is NP-hard!

↳ computational challenges!

▷ Possible approach: Lasso Regression:

$$\hat{\beta}_n = \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \frac{1}{n} \|A\beta - y\|_2^2 + \lambda \|\beta\|_1 \right\}$$

Computationally, Lasso Regression is

▷ more expensive than Ridge or LR, but

▷ still convex optimization problem, so polynomial (and efficient) regularization

algorithms exist.

▷ parameter

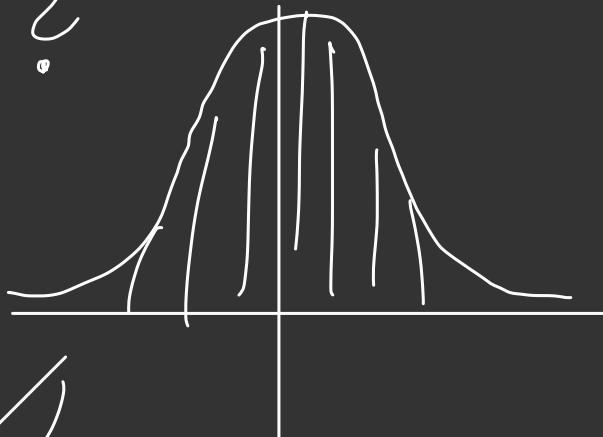
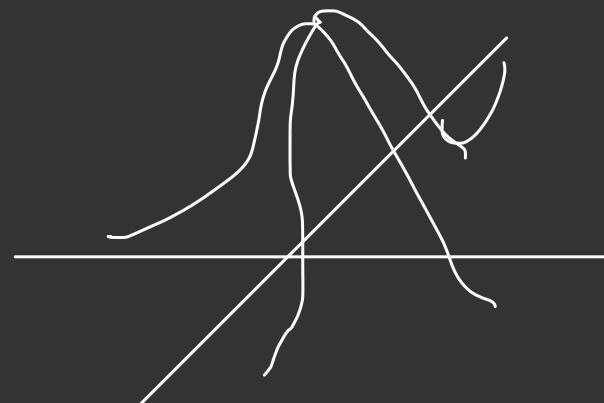
▷ Since Lasso Regression only approximates solution of ERM on $\mathcal{F}_k^{\text{sparse}}$: In addition to estimation and approximation errors, we now also have an optimization error.

High Dimensional Geometry

Q: How do $x_1, x_2 \in \mathbb{R}^d$, $d > 1$ relate to each other when
"generic"?

• 1-dim Gaussian:

$$d=1$$



• d-dim Gaussian:

$$\|x\| \sim \sqrt{d} + O(1)$$

\Rightarrow All very far from origin!

• If $x_1, x_2 \in \mathbb{R}^d$ indep. d-dim Gaussian

$$\Rightarrow \|x_1 - x_2\|_2 \stackrel{\text{close to}}{\sim} \text{const.} \quad \& \text{large.}$$