

Inmas Machine Learning Workshop

Instructor: Christian Kiemmerle

TAs:

- Lanlan Ji
- Vittorio Laprinzo
- Salma Tarmoun
- Sichen Yang

- Goals:
- ▷ Learn to use computational tools to learn from data
 - ↳ Intuition for their power & challenges
 - ▷ Gain intuition of geometry of high-dimensional spaces
 - ▷ Learn to use Python to apply these techniques

What is Data Science?

[Tukey '62]: "Data Analysis" as an empirical science:

- Procedures for gathering data, interpret data
- Uses mathematical statistics
- "reliance upon the test of experience as ultimate standard of validity"

Our focus: Prediction instead of Inference ← focus of "statistics"

Last 15-20 years: Led to technological advances in

pattern / image recognition, machine translation, targeted advertisement,
(semi-) autonomous cars

"machine learning", "artificial intelligence (AI)"

Common Task Framework: [Lieberman '10, Donoho]

- ▷ Public "training" dataset: List of observations with labels
- ▷ Competitors with common task to infer prediction rule from training data
↳ ^{label}
- ▷ Submit to Referee, reports accuracy of prediction rule applied to (hidden) testing dataset

Ex: · \$1M Netflix prize (2006-2008)
· ImageNet

Other aspects: · Available computer hardware
· Software / communication frameworks facilitating reproducibility

Goal: Learn from data.

- a) Ex.:
- Predict salary of professors from employment data
 - Detect spam e-mails based on large set of spam/non-spam e-mails.

Supervised Learning

b) "Learning without teacher", find meaningful data representation / summary

- Ex.:
- Find categories among pictures on phone
 - Visualize complex genetics data to be interpreted by humans

Unsupervised Learning

The Framework of Statistical Learning

- ▷ $X \subset \mathbb{R}^k$: domain set (e.g. space of ^{images with fixed} pixel number)
- ▷ $Y \subset \mathbb{R}^q$: target set / set of labels (e.g. $Y = \{0, 1\}$ if spam/
non-spam classification)
- ▷ Let \mathcal{D} be a probability distribution on $X \times Y$

Assume we are given a training set $S := (x_i, y_i)_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$

- ▷ Goal: Find a predictor/classifier 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 Algorithms:

Specific 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 .

Many learning algorithms:

Empirical Risk Minimization

$$\hat{h}_n =$$

$$\underset{h \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(x_i))$$

Ex: Linear Regression: $\triangleright \ell(y, z) = (y - z)^2$

$$\triangleright \mathcal{F} := \left\{ x \mapsto \beta_0 + \alpha x, \beta_0, \beta_1 \in \mathbb{R}, \beta = \begin{pmatrix} \beta_1 \\ \beta_0 \end{pmatrix} \in \mathbb{R}^2 \right\}$$

Bias - Variance Tradeoff:

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

generalization error

$$h_{\mathcal{F}} = \underset{h \in \mathcal{F}}{\operatorname{argmin}} L_2(h)$$

estimation error

\triangleright induced by choice / size of S

\triangleright For fixed sample size $|S|$,
Larger if \mathcal{F} large

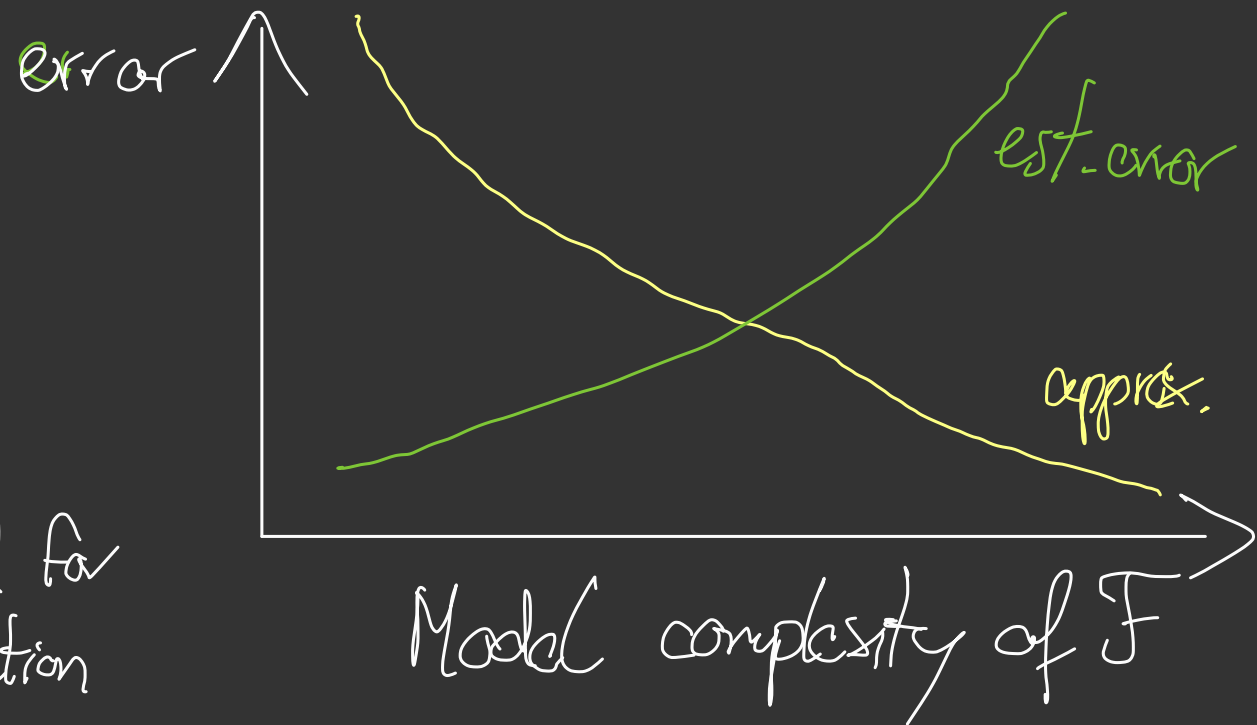
approximation error

\triangleright smaller if \mathcal{F} large

Cross Validation.

Goal: Reduce bias induced by training set S .

▷ Fair comparison only if: S_{train}, S_{test} $\sim \mathcal{D}$
 used in learning method used for evaluation



▷ K-Fold Cross validation



$R(h)$ a term that quantifies "complexity" of $h \in \mathcal{F}$
 $\lambda > 0$

Controlling Model Complexity via Regularization

"Regularized loss minimization" (RLM)

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

Ridge Regression:

Linear Regression with

$$R(h) = \|\beta\|_2^2$$

Choose \mathcal{F} space of linear functions

$$R(h) := \|\beta(h)\|_2^2$$

$$\Rightarrow \hat{\beta}_n = \beta(\hat{h}_n) = \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{n} \|X\beta - y\|_2^2 + \lambda \|\beta\|_2^2$$

$$= (X^T X + \lambda I)^{-1} X^T (y)$$

▷ If $\lambda = 0$: \rightarrow linear regression

▷ If $\lambda \rightarrow \infty$: coefficients $\hat{\beta}$ "shrunk to 0".

▷ λ in between: balancing fit of linear model and size of coefficients.

• Strong connection to hypothesis set $\mathcal{F}_+ = \{h: \mathbb{R} \rightarrow \mathbb{R} : h(x) = \langle \beta, x \rangle \text{ s.t. } \|\beta\|_2 \leq t\}$

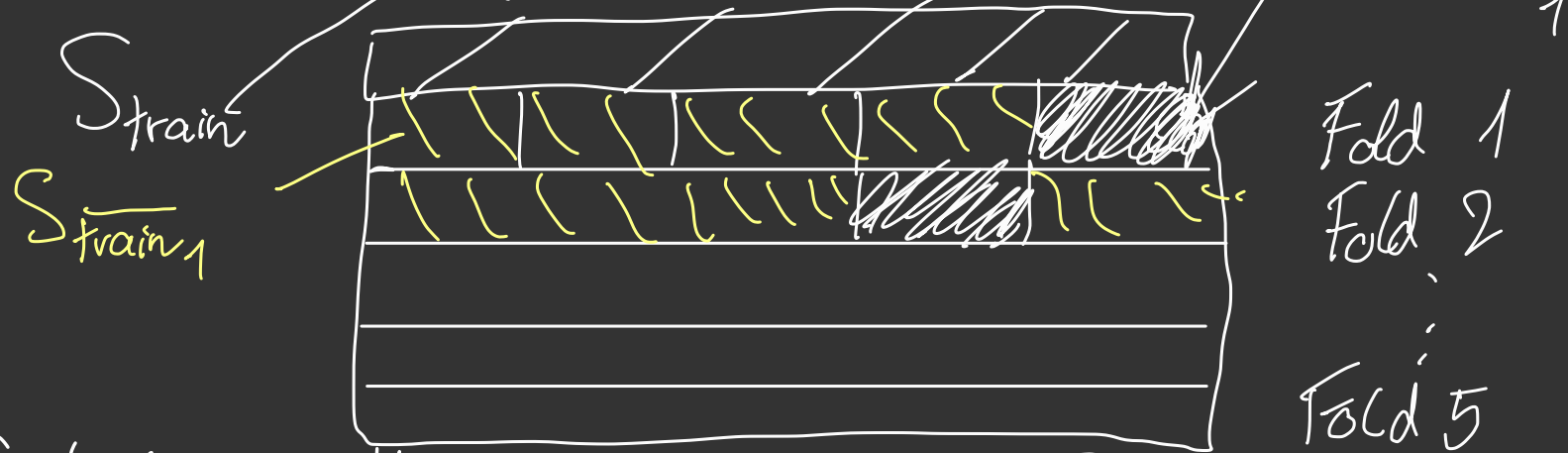
bound on coefficients

Cross Validation: Idea: Reduce bias induced by training set S .

▷ Fair comparisons only if:



▷ K-Fold Crossvalidation



Used for
"Model Selection"

Controlling Model Complexity via Regularization

▷ Modify learning algorithm for some hypothesis space \mathcal{F} : For some $\lambda > 0$,

$$\hat{h}_n := \underset{h \in \mathcal{F}}{\operatorname{argmin}} \left\{ \underbrace{L_S(h)}_{\text{loss term}} + \underbrace{\lambda R(h)}_{\text{regularization term}} \right\}$$

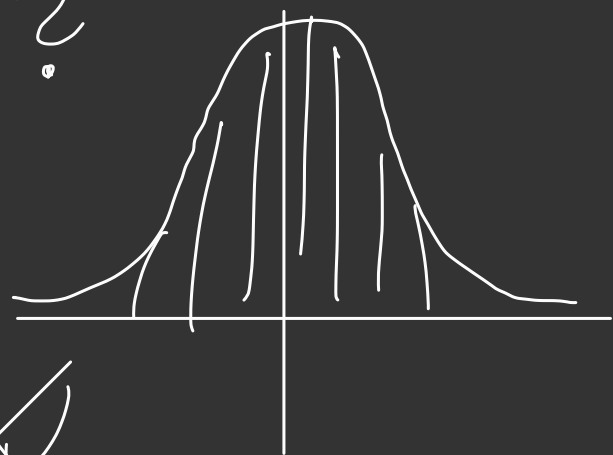
$R(h)$ a term that quantifies "complexity" of $h \in \mathcal{F}$.

"regularized loss minimization" (RLM),

▷ Compared to ERM, term $\lambda R(h)$ penalizes too complex instances of \mathcal{F} .

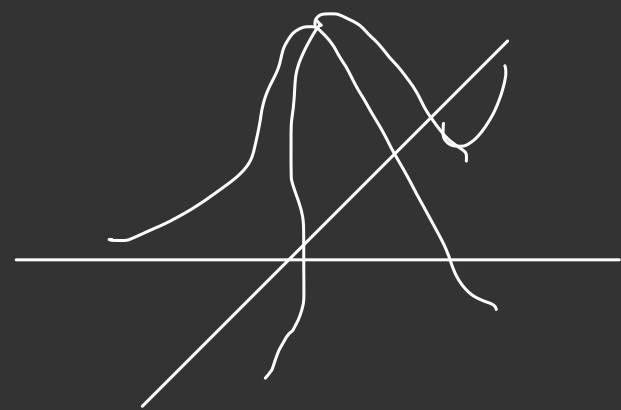
High Dimensional Geometry

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



• 1-dim Gaussian:

$d=2$



• d -dim Gaussian:

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

\Rightarrow All very far from origin!

• $\forall x_1, x_2 \in \mathbb{R}^d$ indep. d -dim Gaussian
 $\Rightarrow \|x_1 - x_2\|_2 \sim$ ^{close to} const. & large.

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

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

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

→ Often improves expressive power of a learning model!

2. Sparse Regression

- If features are designed to "explain" the target variable as a linear combination of few features (e.g., $k \ll n$), we can use

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

where $\|\beta\|_0 = \sum_{i=1}^n \mathbb{1}_{\{\beta_i \neq 0\}}$ is the number of non-zero coefficients of β .

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

↳ computational challenges!

- Possible approach: Lasso Regression:

$$\hat{\beta}_n = \arg \min_{\beta \in \mathbb{R}^2} \left\| A\beta - y \right\|_2^2 + \lambda \|\beta\|_1 \quad (*)$$

- Unlike linear/ridge regression, (*) has no closed form solution, but convex optimization problem ← well-established theory/methods exist.

- With respect to original class $\mathcal{F}_h^{\text{sparse}}$:

$$\text{Generalization error} = \underline{\text{Optimization error}} + \text{estimation error} + \text{approximation error}$$