

Inmas Data Science Workshop

Instructor: Christian Kemmerle
TAs: Daniel Fuentes-Keuthan
Patrick Martin

Dates & Times:

- Fr, March 19: 2-5 pm ET
- Sat, March 20: ▷ 9-12 am ET
▷ 2-5 pm ET
- Sun, March 21: ▷ 9-12 am ET

Goals:

- ▷ 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 statistical

- "reliance upon the test of experience as ultimate standard of validity"

Our focus: Prediction instead of Inference

"machine learning" "artificial intelligence (AI)"

Common Task Framework:

- ▷ Public "training" dataset: list of observations with labels
- ▷ Competitors with common task to infer label prediction rate from training data
- ▷ submit to Referees, report accuracy of prediction rate applied to (hidden) testing dataset.

Ex: • KDD Cup (2006 - 2008)

• ImageNet

Other aspect:

• Available computer hardware

• Better software framework

Goal: Learn from data.

- a) Ex.:
- Predict salary of professors based on discipline, employment length
 - 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 all 64×64 RGB pictures)
- $Y \subset \mathbb{R}^q$: target set / set of labels

▷ Let \mathcal{D} be a probability distribution $X \times Y$.

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 algorithm:

• Specific algorithm that maps S to a specific function $\hat{h}_n \in \mathcal{F}$ of a hypothesis space $\mathcal{F} \subseteq \{h: X \rightarrow Y\}$ based on the information of training set S .

Many learning algorithms amount to: Empirical Risk Minimization

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

Example: Linear Regression ($q=1$)

o $\ell(y, z) := (y - z)^2$

$+ \frac{1}{2} \langle \gamma, M \gamma \rangle$

o $\mathcal{F} := \left\{ x \mapsto \beta_0 + \langle x, \beta \rangle, \beta_0 \in \mathbb{R}, \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} \in \mathbb{R}^k \right\}$

Using ERM, we obtain

$$(\hat{\beta}_0, \hat{\beta}) = \underset{\beta_0, \beta}{\operatorname{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n (\beta_0 + \langle x_i, \beta \rangle - y_i)^2 \right\}$$

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

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

$$X = \begin{pmatrix} -x_1 & - \\ \vdots & - \\ -x_n & - \end{pmatrix}$$

exercise

$$= (X^T X)^{-1} (X^T y)$$

Idea: $L_S(\hat{h}_n) \approx L_D(\hat{h}_n) \approx \text{small}$
if S represents D well

Bias - Complexity Trade-off

Q: How to choose \mathcal{F} ?

$$L_{\mathcal{D}}(\hat{h}_n) - \min_h L_{\mathcal{D}}(h) = \underbrace{\left(L_{\mathcal{D}}(\hat{h}_n) - L_{\mathcal{D}}(h_{\mathcal{F}}) \right)}_{\text{estimation error}} + \underbrace{\left(L_{\mathcal{D}}(h_{\mathcal{F}}) - \min_h L_{\mathcal{D}}(h) \right)}_{\text{approximation error}}$$

generalization error

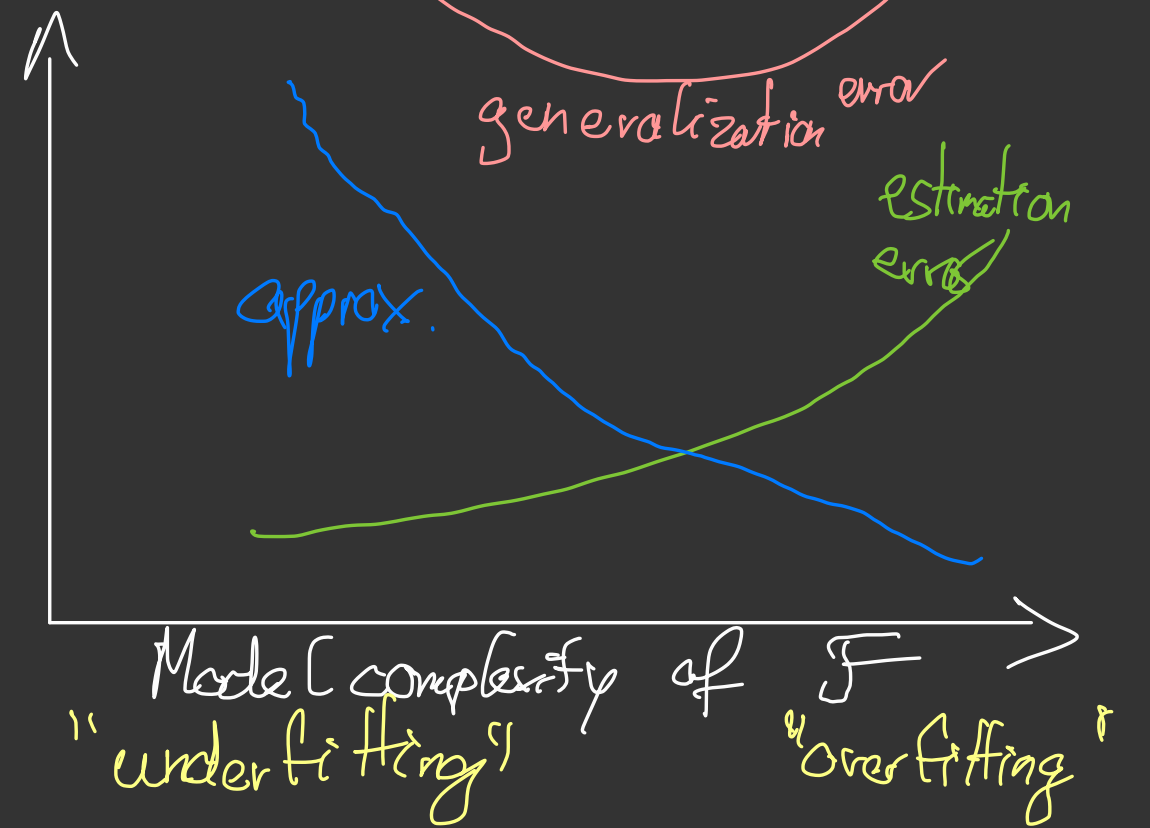
where $h_{\mathcal{F}} = \underset{h \in \mathcal{F}}{\operatorname{argmin}} L_{\mathcal{D}}(h)$

minimizer of expected risk over hypothesis space \mathcal{F}

estimation error

approximation error

error



Model complexity of \mathcal{F}
 "underfitting" "overfitting"

$\frac{1}{\gamma_0}$ $\frac{1}{\gamma}$

Preprocessing: If domain set s.t. $X \subset \mathbb{R}^k$, we can

Define a feature map $\phi: X \rightarrow \tilde{X} \subset \mathbb{R}^l$

(with $k \ll l$) consider

$\tilde{S} = (\phi(x_i), y_i)_{i=1}^n$ instead of $S = (x_i, y_i)_{i=1}^n$

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

$$\Phi(x) = (1, x^1, x^2, x^3, x^4)$$

→ Often improves expressive power of a (learning model).

• $\Phi(x) = \log(x)$

Controlling Model Complexity via Regularization

▷ Modify learning algorithm for some hypothesis space \mathcal{F} :

For $\lambda > 0$,

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

$\forall h \in \mathcal{F}$

$0 \leq R(h)$ is a term that "quantifies complexity" of $h \in \mathcal{F}$.

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

Q: How to choose λ ?

1. Ridge Regression

λ : "regularization parameter"

Choose \mathcal{F} space of linear functions

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

$$\begin{aligned} \Rightarrow \hat{\beta}_n &= \beta(\hat{h}_n) = \underset{\tilde{\beta} \in \mathbb{R}^{k+1}}{\text{argmin}} \left\{ \frac{1}{n} \|X\tilde{\beta} - y\|_2^2 + \lambda \|\tilde{\beta}\|_2^2 \right\} \\ &= (X^T X + \lambda I)^{-1} X^T(y) \end{aligned}$$

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

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

\triangleright λ in between: balancing fit of the linear model and size of coefficients.

Strong connection to hypothesis set: $\mathcal{F}_+ = \left\{ h: \mathbb{R}^{k+1} \rightarrow \mathbb{R} : h(x) = \langle \tilde{\beta}, x \rangle \right.$
 $\left. \text{s.t. } \|\tilde{\beta}\|_2 \leq t \right\}$

2. Sparse Regression

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

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

where $\|\beta\|_0 = \sum_{i=1}^n \mathbb{1}_{\{\beta_i \neq 0\}}$ is nr. of non-zero coefficients of $\beta \in \mathbb{R}^k$.

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

↳ computational challenges!

- Possible approach: Lasso Regression:

$$\hat{\beta}_n = \underset{\beta \in \mathbb{R}^{k+1}}{\text{argmin}} \left\{ \frac{1}{n} \|A\hat{\beta} - y\|_2^2 + \lambda \cdot \|\beta\|_1 \right\} \quad (*)$$

- (*) has no closed form solution, but is a convex optimization problem.

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

Generalization error = Optimization error + estimation error + approximation error

▷ Alternative algorithm: Orthogonal Matching Pursuit

Orthogonal matching pursuit (OMP)

Input: measurement matrix \mathbf{A} , measurement vector \mathbf{y} .

Initialization: $S^0 = \emptyset$, $\mathbf{x}^0 = \mathbf{0}$.

Iteration: repeat until a stopping criterion is met at $n = \bar{n}$:

$$S^{n+1} = S^n \cup \{j_{n+1}\}, \quad j_{n+1} := \operatorname{argmax}_{j \in [N]} \{|(\mathbf{A}^*(\mathbf{y} - \mathbf{A}\mathbf{x}^n))_j|\}, \quad (\text{OMP}_1)$$

$$\mathbf{x}^{n+1} = \operatorname{argmin}_{\mathbf{z} \in \mathbb{C}^N} \{\|\mathbf{y} - \mathbf{A}\mathbf{z}\|_2, \operatorname{supp}(\mathbf{z}) \subset S^{n+1}\}. \quad (\text{OMP}_2)$$

Output: the \bar{n} -sparse vector $\mathbf{x}^\# = \mathbf{x}^{\bar{n}}$.

▷ "Greedy" algorithm.

- Sparse regression: Strongly related to problems in signal and image processing:
 - ▷ Speed-up of MRI measurements
 - ▷ Interpolation of geophysical data