

100 Days Of Machine Learning Code

Sergio-Feliciano Mendoza-Barrera

CONTENTS

I	100DaysOfMLCodeSFMB	1
II	Day 1 [Thu Jul 5 20:58:13 CDT 2018]	1
II-A	Siraj Raval Tweet	1
II-B	Challenge	1
II-C	Following my own path using Julia, R and Python	1
II-D	Install Julia	1
II-E	Install python	1
II-F	Install Emacs	1
III	Day 2 [Fri Jul 6 19:21:03 CDT 2018]	1
III-A	Course description	2
III-B	Prerequisites	2
III-C	Grading	2
III-D	Problem Sets	2
III-E	Projects	2
III-F	Syllabus	2
III-G	Class 1. Course at a Glance	3
IV	Day 3 [Sat Jul 7 13:12:57 CDT 2018]	3
V	Day 4 [Sun Jul 8 12:59:19 CDT 2018]	3
VI	Day 5, 6, 7 [Init Mon Jul 9 16:12:01 CDT 2018]	3
VI-A	Math camp	3
VI-B	Functional and Operators (Matrices)	8
VI-C	Probability Space	12
VII	Day 6 (Pending)	13
VII-A	9.520/6.860, Class 02	13
VIII	References	14

According to the challenge we are going to explore a project and maybe the possibility to monetize a solution of an industry problem using ML.

URL

"Who's ready to take the 100 days of ML code challenge? That means coding machine learning for at least an hour everyday for the next 100 days. Pledge with the #100DaysOfMLCode hashtag, I'll give the first few winners a shoutout!"

I. 100DAYSOFMLCODESFMB

After read this tweet can be a great way to master something, do not know what is, yet.

In day 2 I decided to explore the MIT course **9.520/6.860: Statistical Learning Theory and Applications** Fall 2017.

II. DAY 1 [THU JUL 5 20:58:13 CDT 2018]

A. Siraj Raval Tweet

Who's ready to take the 100 days of ML code challenge? That means coding machine learning for at least an hour everyday for the next 100 days. Pledge with the #100DaysOfMLCode hashtag, I'll give the first few winners a shoutout!

[Link to Tweet.](#)

B. Challenge

Pick an industry that sounds exciting, find a problem they have, think about how AI could be applied to that problem, locate a relevant dataset, apply AI to the dataset, monetize the solution.

C. Following my own path using Julia, R and Python

In that order, trying to go as far as I can using the first one syntax.

D. Install Julia

```
_ _ _(_)_ | A fresh approach to technical computing
(_)_ | (_)(_) | Documentation: https://docs.julialang.org
-- _|_ _ _ | Type "?help" for help.
| | | | | /_`| |
| | | | | (| | | Version 0.6.3 (2018-05-28 20:20 UTC)
/_| \_`|_|_| \_`| | Official http://julialang.org/ release
|_|/ | x86_64-pc-linux-gnu
```

E. Install python

```
Python version :: sys.version_info(major=3, minor=6, micro=5, releaselevel='final', serial=0)
=====
OpenCV version :: 3.4.0
=====
Tensorflow version :: 1.4.1
```

F. Install Emacs

I am joking, Emacs is always installed ;-).

III. DAY 2 [FRI JUL 6 19:21:03 CDT 2018]

9.520/6.860: Statistical Learning Theory and Applications, Fall 2017

Instructors: Tomaso Poggio (TP), Lorenzo Rosasco (LR), Georgios Evangelopoulos (GE)

A. Course description

The course covers foundations and recent advances of Machine Learning from the point of view of Statistical Learning and Regularization Theory.

Understanding intelligence and how to replicate it in machines is arguably one of the greatest problems in science. Learning, its principles and computational implementations, is at the very core of intelligence. During the last decade, for the first time, we have been able to develop artificial intelligence systems that can solve complex tasks, until recently the exclusive domain of biological organisms, such as computer vision, speech recognition or natural language understanding: cameras recognize faces, smart phones understand voice commands, smart speakers/assistants answer questions and cars can see and avoid obstacles.

The machine learning algorithms that are at the roots of these success stories are trained with labeled examples rather than programmed to solve a task. Among the approaches in modern machine learning, the course focuses on regularization techniques, that provide a theoretical foundation to high-dimensional supervised learning. Besides classic approaches such as Support Vector Machines, the course covers state of the art techniques using sparsity or data geometry (aka manifold learning), a variety of algorithms for supervised learning (batch and online), feature selection, structured prediction, and multitask learning and principles for designing or learning data representations. Concepts from optimization theory useful for machine learning are covered in some detail (first order methods, proximal/splitting techniques, . . .).

The final part of the course will focus on deep learning networks. It will introduce an emerging theory formalizing three key areas for the rigorous characterization of deep learning: approximation theory – which functions can be represented efficiently?; optimization theory – how easy is it to minimize the training error?; and generalization properties – is classical learning theory sufficient for deep learning? It will also outline a theory of hierarchical architectures that aims to explain how to build machine that learn using cortex principles and similar to how children learn: from few labeled and many more unlabeled data.

The goal of the course is to provide students with the theoretical knowledge and the basic intuitions needed to use and develop effective machine learning solutions to challenging problems.

B. Prerequisites

We will make extensive use of basic notions of calculus, linear algebra and probability. The essentials are covered in class and in the math camp material. We will introduce a few concepts in functional/convex analysis and optimization. Note that this is an advanced graduate course and some exposure on introductory Machine Learning concepts or courses is expected. Students are also expected to have basic familiarity with MATLAB/Octave.

C. Grading

Pset and project tentative dates: (slides).

D. Problem Sets

Problem Set 1 Problem Set 2 Problem Set 3 Problem Set 4

Submission instructions: Follow the instructions included with the problem set. Use the latex template for the report (there is a maximum page limit). Submit your report online through stellar.mit by the due date/time and a printout in the first class after the due date.

E. Projects

Reports are 1-page, extended abstracts using NIPS style files.

1) *Projects archive*: List of Wikipedia entries, created or edited as part of projects during previous course offerings.

F. Syllabus

URL.

Follow the link for each class to find a detailed description, suggested readings, and class slides. Some of the later classes may be subject to reordering or rescheduling.

Class	Date	Title	Instructor(s)
Class 01	Wed Sep 06	The Course at a Glance	TP
Class 02	Mon Sep 11	The Learning Problem and Regularization	LR
Class 03	Wed Sep 13	Reproducing Kernel Hilbert Spaces	LR
Class 04	Mon Sep 18	Positive Definite Functions, Feature Maps and Mercer Theorem	LR
Class 05	Wed Sep 20	Tikhonov Regularization and the Representer Theorem	LR
Class 06	Mon Sep 25	Logistic Regression and Support Vector Machines	LR
Class 07	Wed Sep 27	Regularized Least Squares	LR
Class 08	Mon Oct 02	Iterative Regularization via Early Stopping	LR
Class 09	Wed Oct 04	Learning with Stochastic Gradients	LR
Class 10	Wed Oct 11	Large Scale Kernel Methods	LR
Class 11	Mon Oct 16	Sparsity Based Regularization	LR
Class 12	Wed Oct 18	Convex Relaxation and Proximal Gradient	LR
Class 13	Mon Oct 23	Structured Sparsity Regularization	LR
Class 14	Wed Oct 25	Multiple Kernel Learning	LR
Class 15	Mon Oct 30	Learning Theory	LR
Class 16	Wed Nov 01	Generalization Error and Stability	LR
Class 17	Mon Nov 06	Online Learning II	Sasha Rakhlin
Class 18	Wed Nov 08	Online Learning II	Sasha Rakhlin
Class 19	Mon Nov 13	Data Representation by Design	GE
Class 20	Wed Nov 15	Learning Data Representation: Dictionary Learning	GE
Class 21	Mon Nov 20	Learning Data Representation: Neural Networks	GE
Class 22	Wed Nov 22	Deep Learning Theory: Approximation	TP
Class 23	Mon Nov 27	Deep Learning Theory: Optimization	TP
Class 24	Wed Nov 29	Deep Learning Theory: Generalization	TP
Class 25	Mon Dec 04	Learning Data Representation: Invariance and Selectivity	TP
Class 26	Wed Dec 06	Deep Networks and Visual Cortex	TP
Class 27	Mon Dec 11	Poster presentations (2 sessions)	

G. Class 1. Course at a Glance

1) *Description*: We introduce and motivate the main theme of much of the course, setting the problem of supervised learning from examples as the ill-posed problem of approximating a multivariate function from sparse data. We present an overview of the theoretical part of the course and sketch the connection between classical Regularization Theory with its RKHS-based algorithms and Learning Theory. We briefly describe several different applications ranging from vision to computer graphics, to finance and neuroscience. The last third of the course will be on data representations for learning and deep learning. It will introduce recent theoretical developments towards a) understanding why deep learning works and b) a new phase in machine learning, beyond classical supervised learning: how to learn in an unsupervised way representations that significantly decrease the sample complexity of a supervised learning.

2) *Slides*:

- Slides for this lecture: PDF.

a) *Youtube video class 2015.* : href6AWZS4Ho2Z8video

[Link here](#)

b) *2017 Course – Center for Brains, Minds and Machines (CBMM)*: hrefQ5itLKscYTAvideo

[Link here](#)

3) *Relevant Reading:*

- Mnih et. al. (Deep Mind), Human-level control through deep reinforcement learning, *Nature* 518, pp. 529–533, 2015.
- Nature Insights, Machine Intelligence (with review article on Deep Learning), *Nature*, Vol. 521 No. 7553, pp. 435–482, 2015.

IV. DAY 3 [SAT JUL 7 13:12:57 CDT 2018]

1) Class 1 video [14.51]

2) Slide [26]

V. DAY 4 [SUN JUL 8 12:59:19 CDT 2018]

1) Class 1 Done

VI. DAY 5, 6, 7 [INIT MON JUL 9 16:12:01 CDT 2018]

A. Math camp

Math camp extra class, optional for those interested: Tue. 09/12, 4:00 pm – 5:30 pm, Singleton auditorium (46-3002).

1) *Description:* We review the basic prerequisites for the course on functional analysis, linear algebra, probability theory and concentration of measure.

2) *Class Reference Material:*

a) *Youtube video:* hrefAsogCoscZgEvideo

[Link here](#)

b) *Local video:* Video file.

c) *Slides:* Slides: PDF. Original URL. Notes/Book appendix: PDF. Original URL.

3) *Some concept testing with data:* We like \mathbb{R}^D because we can

Addition

v = [1, 2, 3];

w = [4, 5, 6];

println(v + w)

[5, 7, 9]

Multiply by numbers

println(3 * v)

[3, 6, 9]

Scalar product

println(v)

println(w)

dot(vec(v), vec(w))

dot(v, w)

[1, 2, 3]

[4, 5, 6]

32

32

Norm

sqrt(dot(vec(v'), vec(v)))

vecnorm(v)

norm(v)

3.7416573867739413

3.7416573867739413

3.7416573867739413

Distances between vectors

```

vecnorm(v - w)
norm(v - w)

5.196152422706632
5.196152422706632

RMS value
norm(v) / sqrt(length(v))

2.160246899469287

Standard deviation
Important note: Julia do not use this definition.
norm(v - mean(v))/sqrt(length(v))

0.8164965809277261

Julia's way:
std(v)

1.0

Angle between two vectors
acos(dot(v, w)/(norm(v) * norm(w)))

0.2257261285527342

```

This what we called "Euclidean" structure. We want to do the samething with $D = \infty$

Vector Space

► A **vector space** is a set V with binary operations

$$+: V \times V \rightarrow V \quad \text{and} \quad \cdot : \mathbb{R} \times V \rightarrow V$$

such that for all $a, b \in \mathbb{R}$ and $v, w, x \in V$:

1. $v + w = w + v$
2. $(v + w) + x = v + (w + x)$
3. There exists $0 \in V$ such that $v + 0 = v$ for all $v \in V$
4. For every $v \in V$ there exists $-v \in V$ such that $v + (-v) = 0$
5. $a(bv) = (ab)v$
6. $1v = v$
7. $(a + b)v = av + bv$
8. $a(v + w) = av + aw$

Figure 1. Vector Space

4) *Vector Space*: Example: \mathbb{R}^D , space of polynomials, space of functions.

5) *Inner Product*:

6) *Cauchy-Schwarz inequality*: $\langle v, w \rangle \leq \langle v, v \rangle^{\frac{1}{2}} \langle w, w \rangle^{\frac{1}{2}}$.

```

println(":: v and w inner product ::")
dot(v, w)
println(":: must be less or equal to ::")
sqrt(dot(v, v)) * sqrt(dot(w, w))

:: v and w inner product :::
32
:: must be less or equal to :::
32.83291031876401

```

7) *Norm*: Can define norm from inner product:

$$\|v\| = \langle v, v \rangle^{\frac{1}{2}}$$

8) *Metric*:

Inner Product

- ▶ An **inner product** is a function $\langle \cdot, \cdot \rangle: V \times V \rightarrow \mathbb{R}$ such that for all $a, b \in \mathbb{R}$ and $v, w, x \in V$:
 1. $\langle v, w \rangle = \langle w, v \rangle$
 2. $\langle av + bw, x \rangle = a\langle v, x \rangle + b\langle w, x \rangle$
 3. $\langle v, v \rangle \geq 0$ and $\langle v, v \rangle = 0$ if and only if $v = 0$.
- ▶ $v, w \in V$ are orthogonal if $\langle v, w \rangle = 0$.
- ▶ Given $W \subseteq V$, we have $V = W \oplus W^\perp$, where $W^\perp = \{v \in V \mid \langle v, w \rangle = 0 \text{ for all } w \in W\}$.
- ▶ Cauchy-Schwarz inequality: $\langle v, w \rangle \leq \langle v, v \rangle^{1/2} \langle w, w \rangle^{1/2}$.

Figure 2. Inner Product

- ▶ A **norm** is a function $\|\cdot\|: V \rightarrow \mathbb{R}$ such that for all $a \in \mathbb{R}$ and $v, w \in V$:
 1. $\|v\| \geq 0$, and $\|v\| = 0$ if and only if $v = 0$
 2. $\|av\| = |a| \|v\|$
 3. $\|v + w\| \leq \|v\| + \|w\|$
- ▶ Can define norm from inner product: $\|v\| = \sqrt{\langle v, v \rangle}$.

Figure 3. Norm definition

- ▶ A **metric** is a function $d: V \times V \rightarrow \mathbb{R}$ such that for all $v, w, x \in V$:
 1. $d(v, w) \geq 0$, and $d(v, w) = 0$ if and only if $v = w$
 2. $d(v, w) = d(w, v)$
 3. $d(v, w) \leq d(v, x) + d(x, w)$
- ▶ Can define metric from norm: $d(v, w) = \|v - w\|$.

Figure 4. Distance

- $B = \{v_1, \dots, v_n\}$ is a **basis** of V if every $v \in V$ can be uniquely decomposed as

$$v = a_1 v_1 + \dots + a_n v_n$$
 for some $a_1, \dots, a_n \in \mathbb{R}$.
- An orthonormal basis is a basis that is orthogonal ($\langle v_i, v_j \rangle = 0$ for $i \neq j$) and normalized ($\|v_i\| = 1$).

Figure 5. Basis

9) Basis:

10) Hilbert Space, overview: Goal: to understand Hilbert spaces (complete inner product spaces) and to make sense of the expression

$$f = \sum_{i=1}^{\infty} \langle f, \phi_i \rangle \phi_i, \quad f \in \mathcal{H}$$

Need to talk about

- 1) Cauchy sequence
- 2) Completeness
- 3) Density
- 4) Separability

- Recall: $\lim_{n \rightarrow \infty} x_n = x$ if for every $\epsilon > 0$ there exists $N \in \mathbb{N}$ such that $\|x - x_n\| < \epsilon$ whenever $n \geq N$.
- $(x_n)_{n \in \mathbb{N}}$ is a **Cauchy sequence** if for every $\epsilon > 0$ there exists $N \in \mathbb{N}$ such that $\|x_m - x_n\| < \epsilon$ whenever $m, n \geq N$.
- Every convergent sequence is a Cauchy sequence (why?)

Figure 6. Cauchy Sequence

11) Cauchy sequence: See definition and examples in the video:

<https://www.youtube.com/watch?v=d190jhAifI>

Link here

```
function nonCauchy(n::Int64)
  for i = 1:n
    print((1 + ((-1)^i)), ", ")
  end
end

n = 12;
nonCauchy(n)

nonCauchy (generic function with 1 method)
```

0, 2, 0, 2, 0, 2, 0, 2, 0, 2,

12) Completeness:

13) Hilbert Space:

14) Orthonormal Basis:

- A Hilbert space has a countable orthonormal basis if and only if it is separable.
- Can write:

- ▶ A normed vector space V is **complete** if every Cauchy sequence converges.
- ▶ Examples:
 1. \mathbb{Q} is not complete.
 2. \mathbb{R} is complete (axiom).
 3. \mathbb{R}^n is complete.
 4. Every finite dimensional normed vector space (over \mathbb{R}) is complete.

Figure 7. Completeness

- ▶ A **Hilbert space** is a complete inner product space.
- ▶ Examples:
 1. \mathbb{R}^n
 2. Every finite dimensional inner product space.
 3. $\ell_2 = \{(a_n)_{n=1}^{\infty} \mid a_n \in \mathbb{R}, \sum_{n=1}^{\infty} a_n^2 < \infty\}$
 4. $L_2([0, 1]) = \{f: [0, 1] \rightarrow \mathbb{R} \mid \int_0^1 f(x)^2 dx < \infty\}$

Figure 8. Hilbert Space

$$f = \sum_{i=1}^{\infty} \langle f, \phi_i \rangle \phi_i, \text{ for all } f \in \mathcal{H}$$

- ▶ A Hilbert space has a countable orthonormal basis if and only if it is separable.
- ▶ Can write:

$$f = \sum_{i=1}^{\infty} \langle f, \phi_i \rangle \phi_i \text{ for all } f \in \mathcal{H}.$$
- ▶ Examples:
 1. Basis of ℓ_2 is $(1, 0, \dots), (0, 1, 0, \dots), (0, 0, 1, 0, \dots), \dots$
 2. Basis of $L_2([0, 1])$ is $1, 2 \sin 2\pi nx, 2 \cos 2\pi nx$ for $n \in \mathbb{N}$

Figure 9. Orthonormal Basis

B. Functional and Operators (Matrices)

- 1) Maps:
- 2) Representation of Continuous Functionals:
- 3) Matrix:
 - A is symmetric if $A^T = A$

Next we are going to review basic properties of maps on a Hilbert space.

- ▶ functionals: $\Psi : \mathcal{H} \rightarrow \mathbb{R}$
- ▶ linear operators $A : \mathcal{H} \rightarrow \mathcal{H}$, such that

$$A(af + bg) = aAf + bAg, \text{ with } a, b \in \mathbb{R} \text{ and } f, g \in \mathcal{H}.$$

Figure 10. Maps

Let \mathcal{H} be a Hilbert space and $g \in \mathcal{H}$, then

$$\Psi_g(f) = \langle f, g \rangle, \quad f \in \mathcal{H}$$

is a continuous linear functional.

Riesz representation theorem

The theorem states that every continuous linear functional Ψ can be written uniquely in the form,

$$\Psi(f) = \langle f, g \rangle$$

for some appropriate element $g \in \mathcal{H}$.

Figure 11. Representation of continuous functionals

4) *Eigenvalues and Eigenvectors:* Video introduction to eigenvalues and eigenvectors
[hrefG4N8vJpf7hMvideo](#)

[Link here](#)

```
A = [3 2; 3 -2]
eig(A)

2x2 Array{Int64,2}:
3  2
3 -2
([4.0, -3.0], [0.894427 -0.316228; 0.447214 0.948683])
```

[MIT lecture](#)

[hrefDzqE7tj7eIMvideo](#)

[Link here](#)

- ▶ Every linear operator $L : \mathbb{R}^m \rightarrow \mathbb{R}^n$ can be represented by an $m \times n$ matrix A .
- ▶ If $A \in \mathbb{R}^{m \times n}$, the transpose of A is $A^\top \in \mathbb{R}^{n \times m}$ satisfying
$$\langle Ax, y \rangle_{\mathbb{R}^m} = (Ax)^\top y = x^\top A^\top y = \langle x, A^\top y \rangle_{\mathbb{R}^n}$$
for every $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^m$.

Figure 12. Matrix

- ▶ Let $A \in \mathbb{R}^{n \times n}$. A nonzero vector $v \in \mathbb{R}^n$ is an eigenvector of A with corresponding eigenvalue $\lambda \in \mathbb{R}$ if $Av = \lambda v$.
- ▶ Symmetric matrices have real eigenvalues.
- ▶ **Spectral Theorem:** Let A be a symmetric $n \times n$ matrix. Then there is an orthonormal basis of \mathbb{R}^n consisting of the eigenvectors of A .
- ▶ Eigendecomposition: $A = V\Lambda V^\top$, or equivalently,

$$A = \sum_{i=1}^n \lambda_i v_i v_i^\top.$$

Figure 13. Eigenvalues and Eigenvectors

```

A = [5 1; 3 3]
eig(A)
2x2 Array{Int64,2}:
5 1
3 3
([6.0, 2.0], [0.707107 -0.316228; 0.707107 0.948683])

println(:: Eigenvectors ::)
x1 = [1; 1]
x2 = [1; -3]
println(:: ===== ::)

println(:: First member A x1 ::)
A * x1
println(:: Second member lambda1 x1 ::)
lambda1 = 6.0
lambda1 * x1

println(:: ===== ::)
println(:: First member A x2 ::)
A * x2

println(:: Second member lambda2 x2 ::)
lambda2 = 2.0
lambda2 * x2

:: Eigenvectors :::
2-element Array{Int64,1}:
1
1
2-element Array{Int64,1}:
1
-3
:: ===== ::

:: First member A x1 :::
2-element Array{Int64,1}:
6
6
:: Second member lambda1 x1 :::
6.0
2-element Array{Float64,1}:
6.0

```

```

6.0

:: ===== ::

:: First member A x2 ::

2-element Array{Int64,1}:
 2
 -6

:: Second member lambda2 x2 ::

2.0

2-element Array{Float64,1}:
 2.0
 -6.0

A2 = [1 5; 3 3]
eig(A2)

2×2 Array{Int64,2}:
 1 5
 3 3
([-2.0, 6.0], [-0.857493 -0.707107; 0.514496 -0.707107])

println(:: Eigenvalues ::)
x1 = [-0.857493; 0.514496]
x2 = [-0.707107; -0.707107]
println(:: ===== ::)

println(:: First member A x1 ::)
A2 * x1
println(:: Second member lambda1 x1 ::)
lambda1 = -2.0
lambda1 * x1
println(:: ===== ::)

println(:: First member A x2 ::)
A2 * x2
println(:: Second member lambda2 x2 ::)
lambda2 = 6.0
lambda2 * x2

:: Eigenvectors ::

2-element Array{Float64,1}:
 -0.857493
 0.514496

2-element Array{Float64,1}:
 -0.707107
 -0.707107
:: ===== ::

:: First member A x1 ::

2-element Array{Float64,1}:
 1.71499
 -1.02899

:: Second member lambda1 x1 ::

-2.0

2-element Array{Float64,1}:
 1.71499
 -1.02899
:: ===== ::

:: First member A x2 ::

2-element Array{Float64,1}:
 -4.24264
 -4.24264

:: Second member lambda2 x2 ::
```

```
6.0
2-element Array{Float64,1}:
-4.24264
-4.24264
```

The **Eigendecomposition** is very useful in Machine Learning algorithms.

- ▶ Every $A \in \mathbb{R}^{m \times n}$ can be written as

$$A = U\Sigma V^\top,$$

where $U \in \mathbb{R}^{m \times m}$ is orthogonal, $\Sigma \in \mathbb{R}^{m \times n}$ is diagonal, and $V \in \mathbb{R}^{n \times n}$ is orthogonal.

- ▶ Singular system:

$$\begin{aligned} Av_i &= \sigma_i u_i & AA^\top u_i &= \sigma_i^2 u_i \\ A^\top u_i &= \sigma_i v_i & A^\top A v_i &= \sigma_i^2 v_i \end{aligned}$$

Figure 14. Singular Value Decomposition

5) *Singular Value Decomposition:*

- ▶ The spectral norm of $A \in \mathbb{R}^{m \times n}$ is

$$\|A\|_{\text{spec}} = \sigma_{\max}(A) = \sqrt{\lambda_{\max}(AA^\top)} = \sqrt{\lambda_{\max}(A^\top A)}.$$

- ▶ The Frobenius norm of $A \in \mathbb{R}^{m \times n}$ is

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2} = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2}.$$

Figure 15. Matrix Norm

6) *Matrix Norm:*

A real symmetric matrix $A \in \mathbb{R}^{m \times m}$ is positive definite if

$$x^t Ax > 0, \quad \forall x \in \mathbb{R}^m.$$

A positive definite matrix has positive eigenvalues.

Note: for positive semi-definite matrices $>$ is replaced by \geqslant .

Figure 16. Positive Definite Matrix

7) *Positive Definite Matrix:* [hrefojUQkGNQbQvideo](#)
[Link here](#)

- ▶ The adjoint of a bounded linear operator $L: \mathcal{H}_1 \rightarrow \mathcal{H}_2$ is a bounded linear operator $L^*: \mathcal{H}_2 \rightarrow \mathcal{H}_1$ satisfying

$$\langle Lf, g \rangle_{\mathcal{H}_2} = \langle f, L^*g \rangle_{\mathcal{H}_1} \text{ for all } f \in \mathcal{H}_1, g \in \mathcal{H}_2.$$
- ▶ L is self-adjoint if $L^* = L$. Self-adjoint operators have real eigenvalues.
- ▶ A bounded linear operator $L: \mathcal{H}_1 \rightarrow \mathcal{H}_2$ is compact if the image of the unit ball in \mathcal{H}_1 has compact closure in \mathcal{H}_2 .

Figure 17. Adjoint and Compactness

8) *Adjoint and Compactness:*

- ▶ Let $L: \mathcal{H} \rightarrow \mathcal{H}$ be a compact self-adjoint operator. Then there exists an orthonormal basis of \mathcal{H} consisting of the eigenfunctions of L ,

$$L\phi_i = \lambda_i \phi_i$$

and the only possible limit point of λ_i as $i \rightarrow \infty$ is 0.

- ▶ Eigendecomposition:

$$L = \sum_{i=1}^{\infty} \lambda_i \langle \phi_i, \cdot \rangle \phi_i.$$

Figure 18. Spectral Theorem for Compact Self-Adjoint Operator

9) *Spectral Theorem for Compact Self-Adjoint Operator:*

C. Probability Space

Typo: $P(X) = \infty$ must be $P(X) = 1$

1) *Real Random Variables:* [ZvWZKLaZVjovideo](#)

Link here

2) *Convergence of Random Variables:*

3) *Law of Large Numbers:*

4) *Concentration Inequalities:*

a) *Chebyshev's Inequality:*

$$P(|X - \mu| \geq \epsilon) \leq \frac{Var(X)}{\epsilon^2}$$

VII. DAY 6 (PENDING)

A. 9.520/6.860, Class 02

1) *Description:* We formalize the problem of learning from examples in the framework of statistical learning theory and introduce key terms and concepts such as loss functions, empirical and excess risk, generalization error and consistency. We briefly describe foundational results and introduce the concepts of hypothesis space and regularization.

2) *Class Reference Material:* L. Rosasco, T. Poggio, *Machine Learning: a Regularization Approach*, MIT-9.520 Lectures Notes, Manuscript, Dec. 2017.

3) *Chapter 1 – Statistical Learning Theory:* Note: The course notes, in the form of the circulated book draft is the reference material for this class. Related and older material can be accessed through previous year offerings of the course.

A triple (Ω, \mathcal{A}, P) , where Ω is a set,

\mathcal{A} a Sigma Algebra, i.e. a family of subsets of Ω s.t.

- ▶ $\mathcal{X}, \emptyset \in \mathcal{A}$,
- ▶ $A \in \mathcal{A} \Rightarrow \Omega \setminus A \in \mathcal{A}$,
- ▶ $A_i \in \mathcal{A}, i = 1, 2, \dots \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{A}$.

P a probability measure, i.e a function $P : \mathcal{A} \rightarrow [0, 1]$

- ▶ $P(\mathcal{X}) = \infty$ (hence and $P(\emptyset) = 0$),
- ▶ Sigma additivity: If $A_i \in \mathcal{A}, i = 1, 2, \dots$ are disjoint, then

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$$

Figure 19. Probability Space

A measurable function $X : \Omega \rightarrow \mathbb{R}$, i.e. mapping elements of the sigma algebra in open subsets of \mathbb{R} .

- ▶ Law of a random variable: probability measure on \mathbb{R} defined as

$$\rho(I) = P(X^{-1}(I))$$

for all open subsets $I \subset \mathbb{R}$.

- ▶ Probability density function of a probability measure ρ on X : a function $p : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$\int_I d\rho(x) = \int_I p(x)dx$$

for open subsets $I \subset \mathbb{R}$.

Figure 20. Real Random Variables

4) Further Reading:

- F. Cucker and S. Smale, On the mathematical foundations of learning, Bulletin of the American Mathematical Society, 2002.
- T. Evgeniou, M. Pontil and T. Poggio, Regularization networks and support vector machines, Advances in Computational Mathematics, 2000.
- S. Villa, L. Rosasco and T. Poggio, On learnability, complexity and stability, "Empirical Inference, Festschrift in Honor of Vladimir N. Vapnik." Springer-Verlag, Chapter 7, 2013.
- V. Vapnik, An overview of statistical learning theory, IEEE Trans. on Neural Networks, 10(5), 1999.

5) Video: [SFxypsvhhMQvideo](#)

Link here

VIII. REFERENCES

- 1) Introductory Machine Learning Notes. Lorenzo Rosasco, MIT, 2017. Original URL.
- 2) Gilbert Strang MIT web page.

$X_i, i = 1, 2, \dots$, a sequence of random variables.

► Convergence in probability:

$$\forall \epsilon \in (0, \infty), \quad \lim_{i \rightarrow \infty} \mathbb{P}(|X_i - X| > \epsilon) = 0.$$

► Almost Sure Convergence

$$\mathbb{P}\left(\lim_{i \rightarrow \infty} X_i = X\right) = 1.$$

Figure 21. Convergence of Random Variables

$X_i, i = 1, 2, \dots$, sequence of independent copies of a random variable X

Weak Law of Large Numbers:

$$\forall \epsilon \in (0, \infty), \quad \lim_{n \rightarrow \infty} P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}[X]\right| > \epsilon\right) = 0.$$

Strong Law of Large Numbers:

$$P\left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \mathbb{E}[X]\right) = 1.$$

Figure 22. Law of Large Numbers

3) Multivariate Calculus best book.

println(":: Update! ::")

:: Update! ::

$X_i, i = 1, 2, \dots$, sequence of independent copies of a random variable $X, \forall \epsilon \in (0, \infty)$

- ▶ Markov inequality

$$P \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| > \epsilon \right) \leq \frac{\mathbb{E}[X]}{\epsilon}$$

- ▶ Chebysev Inequality

$$P \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| > \epsilon \right) \leq \frac{\text{Var}[X]}{\epsilon^2}$$

- ▶ Höeffding Inequality: If $|X_i| \leq c$

$$P \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| > \epsilon \right) \leq 2e^{-\frac{\epsilon^2 n}{c^2}}$$

Figure 23. Concentration Inequalities (ERRONEOUS)