

## An Introduction to Statistical Learning

### What this is

- the following are notes from Ch1 from ISLRv2
- **Chapter 1: Introduction**
- pp 1-6
- `**/Spring_2026/DATA780/Unit1/ISLRv2_corrected_June_2023.pdf`

### An Overview of Statistical Learning

**Statistical learning** is a collection of **methods** used to **understand**, discover **patterns**, and make **predictions** in data and is widely used in business, medicine, science [sic], and public policy applications. `

### What is Statistical Learning?

|                              |  |
|------------------------------|--|
| <b>Supervised Learning</b>   | <ul style="list-style-type: none"> <li>• used to <b>predict or estimate</b> the <b>output</b> (Y) using <b>inputs</b> (X)</li> <li>• Common tasks include predicting <b>numerical</b> and <b>categorical</b> values and classifications</li> </ul>   |
| <b>Unsupervised Learning</b> | <ul style="list-style-type: none"> <li>• <b>only inputs</b> are observed</li> <li>• there is <b>no output</b> variable</li> <li>• used to <b>discover structure or relationships</b> in the data</li> <li>• common tasks include <b>clustering</b> and <b>dimension reduction</b></li> </ul> |

### Regression Example: Wage Data

|                         |   |
|-------------------------|---|
| <b>Goal</b>             | <ul style="list-style-type: none"> <li>• <b>predict</b> a man's <b>wage</b> using information such as age, education level, and calendar year</li> </ul>  |
| <b>Key Observations</b> | <ul style="list-style-type: none"> <li>• wages <b>increase with age</b> up to age 60, <b>then decreases</b></li> <li>• wages <b>increase over time</b></li> <li>• <b>higher education</b> is associated with <b>higher wages</b></li> <li>• there is <b>substantial variability</b> in wages</li> </ul> |
| <b>Insights</b>         | <ul style="list-style-type: none"> <li>• no single variable predicts wage well on its own</li> <li>• <b>combining multiple variables improves prediction</b></li> <li>• relationships can be non-linear (especially with age)</li> </ul>  |

### Classification Example: Stock Market Data

|                  |  |
|------------------|--|
| <b>Goal</b>      | <ul style="list-style-type: none"> <li>• <b>predict trends</b> in stock market</li> </ul>  |
| <b>Data Used</b> | <ul style="list-style-type: none"> <li>• daily percentage changes in the S&amp;P 500</li> <li>• <b>past 5 days of returns</b></li> </ul> |

|   |  |
|---|--|
| <b>Key Considerations</b>                             | <ul style="list-style-type: none"> <li>the <b>output</b> is <b>categorical</b>, not numerical</li> <li>past returns show very weak predictive power</li> </ul>   |
| <b>Results</b>  | <ul style="list-style-type: none"> <li>prediction accuracy of ~60% using statistical learning methods</li> </ul>   |
| <b>Unsupervised Learning Example: Gene Expression</b> |  |
| <b>Goal</b>   | <ul style="list-style-type: none"> <li>identify groups of similar cancer cell lines using gene expression data</li> </ul>  |
| <b>Data Used</b>                                      | <ul style="list-style-type: none"> <li>64 cancer cell lines</li> <li>6,830 gene measurements per cell line</li> <li>no output variable</li> </ul>  |
| <b>Approach</b>                                       | <ul style="list-style-type: none"> <li>Reduce thousands of variables to a small number (dimension reduction)</li> <li>visualize the data in two dimensions</li> <li>look for natural groupings/clusters</li> </ul>   |
| <b>Result</b>   | <ul style="list-style-type: none"> <li>cell lines form visible clusters</li> <li>clusters tend to correspond to actual cancer types</li> <li>cancer type information was not used to form the clusters</li> <li>evidence provided that unsupervised learning can uncover real structure</li> </ul> |
| <b>Supervised Learning: A Brief History</b>           |  |
| <b>Early Developments</b>                             | <ul style="list-style-type: none"> <li>1800s: least squares &gt;&gt; linear regression</li> <li>1936: linear discriminant analysis</li> <li>1940s: logistic regression</li> <li>1970s: generalized linear models</li> </ul>  |
| <b>Major Advancements</b>                             | <ul style="list-style-type: none"> <li>early methods were restricted by computational limits</li> <li>1980s onward computing power enabled non-linear methods</li> </ul>   |
| <b>Modern Era</b>                                     | <ul style="list-style-type: none"> <li>statistical learning is a distinct field</li> <li>methods have become widely available through tools such as R</li> <li>expanded far beyond statistics and computer science</li> </ul>  |

## NumPy: Broadcasting

### What this is

- the following are notes from the NumPy v2.4 manual
- <https://numpy.org/doc/stable/user/basics.broadcasting.html>

### An Overview of NumPy Broadcasting

#### What Broadcasting is

- a rule NumPy uses to let you do **arithmetic** on **arrays** with **different shapes**

#### What Broadcasting does

- when NumPy sees two arrays that do **not match shape** exactly it attempts to **expand** the **smaller array**

#### Result

- NumPy **avoids loops** and unnecessary data copying, making array math **fast** and **efficient**

### Why Broadcasting Matters

#### NumPy without Broadcasting

- normally NumPy does element-by-element operations which traditionally only work with array with the same shape

```
a = np.array([1, 2, 3])
b = np.array([2, 2, 2])
a * b          # Works
```

#### NumPy with Broadcasting

- let's NumPy **relax** the **constraints**
- facilitates **scalar** and **non-conformable** array arithmetic

### Scalar Multiplication Example

#### Problem

- a problem requires **multiplying** a 3 x 1 array by a **scalar**

```
a = np.array([1, 2, 3])
b = 2
a * b
```

#### Solution

- NumPy conceptually **“stretches”** ‘b’ to **match** the shape of a
- this is achieved **without duplicating** the scalar in **memory**
- essentially makes the operation:

```
[1 * 2, 2 * 2, 3 * 2]
```

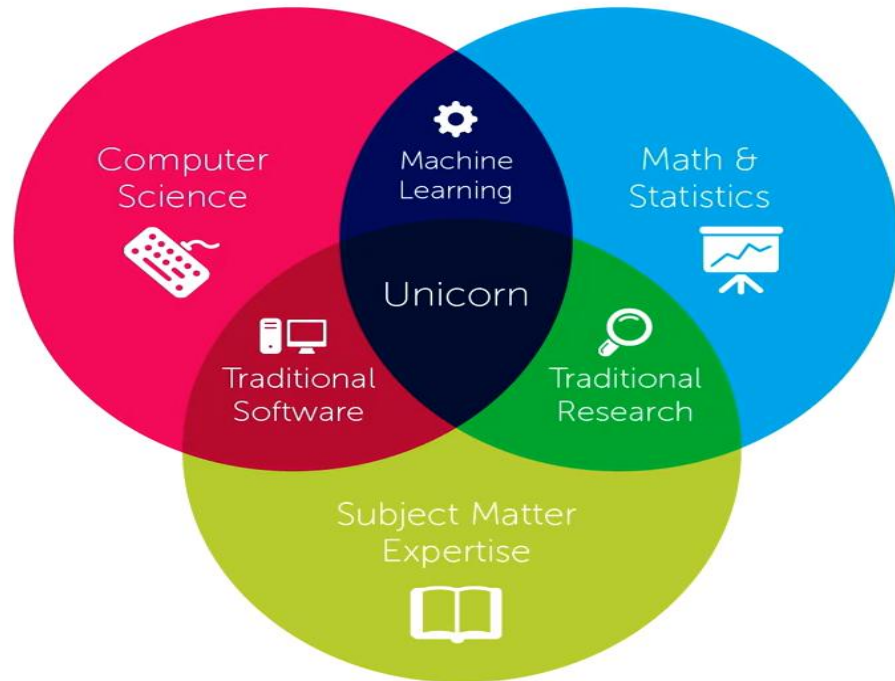
| How Broadcasting Works |   |
|------------------------|---|
| Rules                  | <ul style="list-style-type: none"> <li>compare shapes from the <b>trailing</b> (rightmost) <b>dimension</b></li> <li>two dimensions are <b>compatible</b> if               <ul style="list-style-type: none"> <li>they are <b>equal</b></li> <li>one of them <b>equals '1'</b></li> </ul> </li> <li>if <b>neither</b> condition is <b>true</b> <ul style="list-style-type: none"> <li>NumPy raises <b>'ValueError'</b></li> </ul> </li> </ul>   |
| What it does           | <ul style="list-style-type: none"> <li>If <b>one array</b> has <b>fewer dimensions</b>, NumPy <b>treats</b> missing dimension sizes as <b>'1'</b></li> <li>If a <b>dimension</b> is <b>'1'</b> in one array but <b>larger</b> in the <b>other</b>, NumPy <b>expands</b> that dimension</li> </ul>   |
| Resulting Array        | <ul style="list-style-type: none"> <li>the resulting shape has the <b>maximum</b> number of <b>dimensions</b> from the <b>input arrays</b></li> <li>each <b>axis</b> is the <b>size</b> of the <b>larger</b> of the two</li> <li>if shapes <b>cannot be aligned</b> NumPy raises an <b>error</b></li> </ul>   |
| Examples               | <ul style="list-style-type: none"> <li>NumPy <b>expands</b> the dimensions of <b>'B'</b> so it is compatible with <b>'A'</b> <pre> Array A shape: (8, 1, 6, 1) Array B shape:  (7, 1, 5) Result shape: (8, 7, 6, 5) </pre> </li> <li>when operation is with scalar and array <b>'b'</b> is <b>broadcasted</b> into <b>'a'</b> <pre> a = np.array([1, 2, 3]) b = 2 a * b           # Works by broadcasting b </pre> </li> <li>when operation is with 1-D and 2-D arrays <pre> a = np.array([[0, 0, 0],                [10,10,10],                [20,20,20],                [30,30,30]])  b = np.array([1,2,3])  a + b           # b stretches across each row of a </pre> </li> </ul> |

|                                |  |
|--------------------------------|--|
|                                | <ul style="list-style-type: none"><li>• a set of arrays are <b>broadcast-able</b> if all can be expanded to a <b>common shape</b> that <b>meets the rules</b></li></ul> <pre>Shapes: (5,1), (1,6), (6,), () → all can broadcast to (5,6)</pre> <ul style="list-style-type: none"><li>• scalar shape '()' acts like a shape of <b>ones</b></li><li>• the <b>1-D array</b> '(6,)' is treated like <b>(1, 6)</b></li><li>• all <b>arrays</b> are conceptually <b>expanded</b> to <b>(5, 6)</b></li><li>• this is all achieved with <b>minimal memory</b> use by NumPy</li></ul> |
| <b>When Broadcasting Fails</b> | <ul style="list-style-type: none"><li>• when two <b>arrays</b> do <b>not</b> satisfy the <b>rules</b> and broadcasting is <b>not possible</b> NumPy shows a <b>broadcasting error</b></li></ul> <pre>Array shapes: (4,3) and (4,) → Error: cannot broadcast because trailing dimensions don't match</pre> <ul style="list-style-type: none"><li>•</li></ul>  |
| <b>Broadcasting Summary</b>    | <ul style="list-style-type: none"><li>• let's NumPy <b>combine</b> arrays of <b>different</b> shapes avoiding loops and memory bloat</li><li>• <b>compares</b> shapes from <b>trailing</b> dimensions</li><li>• dimensions must either <b>match</b> or be <b>'1'</b></li><li>• if conditions <b>not met</b> NumPy raises an <b>error</b></li></ul>   |

## Async Materials

### Machine Learning

# Data Science



## Machine Learning Terminology

### Inputs to Models

- **features** – the individual **measurable properties** in a model
- **covariates** – input variables that are **statistically related** to the output
- **dimensions** – **number** of input variables or **axes** in a feature space
- **parameters** – learnable **variables needed** to compute the model's **output**

### Contents of a Dataset

- **instances** – a single **unit of observation** in a dataset
- **samples** – a **single** observed **data point**
- **examples** – a **training** or **testing** data point

### Quick Reference

| Concept Type | Term             | Emphasis                               |
|--------------|------------------|--|
| inputs       | <b>feature</b>   | <b>measurable</b> input variable       |
|              | <b>covariate</b> | <b>input</b> related to <b>outcome</b> |
|              | <b>dimension</b> | <b>number</b> of input <b>axes</b>     |
| dataset unit | <b>instance</b>  | one <b>entity</b> or <b>case</b>       |
|              | <b>sample</b>    | <b>observed</b> data point             |
|              | <b>example</b>   | labeled <b>training/test</b> case      |

|                                  |  |
|----------------------------------|--|
| <b>Training</b>                  | <ul style="list-style-type: none"> <li>• <b>inference</b> – using understood relationships to draw conclusions</li> <li>• <b>prediction</b> – using data, machine learning, and statistical modeling to forecast future outcomes</li> <li>• <b>classification label</b> – predefined category or class assigned to a data point correlating to output in supervised machine learning</li> <li>• <b>types of learning</b> <ul style="list-style-type: none"> <li>○ <b>supervised</b> – mapping input to output</li> <li>○ <b>unsupervised</b> – learn data characteristics</li> <li>○ <b>reinforcement</b> – learn how to interact with an environment through sequential ‘stacked’ learning</li> </ul> </li> </ul> |
| <b>Types of Machine Learning</b> |  |
| <b>Supervised Learning</b>       | <ul style="list-style-type: none"> <li>• <b>inputs &gt; outputs</b> <ul style="list-style-type: none"> <li>○ {example1, output1}</li> <li>○ {example2, output2}</li> </ul> </li> <li>• speech &gt; text</li> <li>• <b>regression</b></li> </ul> <div data-bbox="532 993 1490 1312"> </div>   |
| <b>Unsupervised Learning</b>     | <ul style="list-style-type: none"> <li>• does <b>not</b> make use of <b>prespecified/annotated examples</b></li> <li>• {example1, example2, example3...}</li> <li>• <math>\{X_i\}_{i=1}^n</math></li> <li>• <b>clustering</b> <ul style="list-style-type: none"> <li>○ attempting to <b>discover</b> the salient <b>groups</b> of the data</li> </ul> </li> <li>• <b>dimensionality reduction</b> <ul style="list-style-type: none"> <li>○ <b>identifies</b> the degrees of freedom or <b>core descriptors</b> of the data</li> </ul> </li> <li>• <b>generative models</b> <ul style="list-style-type: none"> <li>○ given training data, generate new samples from same distribution</li> </ul> </li> </ul>        |
| <b>Reinforcement Learning</b>    | <ul style="list-style-type: none"> <li>• interact with an environment and assess the environment state to then output actions that change that state</li> </ul>  |

|  |   |
|--|---|
|  | <ul style="list-style-type: none"> <li>• predictions build on one another to maximize ‘rewards’ through iterative application of the assessment, action, change state, assess, ...</li> </ul>   |
| <b>Keys to Success with Machine Learning</b>             |   |
| <b>Three Pillars of Machine Learning</b>                 | <ul style="list-style-type: none"> <li>• big <b>datasets</b></li> <li>• <b>fast</b> processing</li> <li>• <b>innovative</b> methods</li> </ul>  |
| <b>Typical ML Model Lifecycle</b>                        | <ul style="list-style-type: none"> <li>• <b>select</b> loss/<b>model-type</b> for data task</li> <li>• <b>optimize</b> model with training <b>data</b></li> <li>• <b>evaluate</b> on held-out data to <b>validate</b> choices</li> <li>• use with unseen <b>future</b> data</li> </ul>  |
| <b>Three Pillars to YOUR Success in Machine Learning</b> | <ul style="list-style-type: none"> <li>• <b>statistics</b> and <b>mathematics</b></li> <li>• computer <b>engineering</b></li> <li>• <b>data analytics</b> acumen</li> </ul>   |
| <b>Philosophy: Intelligence Can Grow</b>                 | <p><b>A Growth Mindset Drives Motivation and Achievement</b></p> <pre> graph LR     A((I can get smarter)) --&gt; B((Learning is my goal))     A --&gt; C((Effort makes me stronger))     B --&gt; D((I'd spend more time and work harder))     C --&gt; D     D --&gt; E((Higher Achievement))   </pre> <p>Blackwell, Trzesniewski &amp; Dweck (2007) <i>Child Development</i></p> |



## Mathematics Review

### Differentiation

#### taking derivatives of functions

$$\frac{d}{dx} f(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$$

### Basic Facts

- *Constant rule:* if  $f(x)$  is constant, then

$$f'(x) = 0.$$

- *Sum rule:*

$$(\alpha f + \beta g)' = \alpha f' + \beta g' \text{ for all functions } f \text{ and } g \text{ and all real numbers } \alpha \text{ and } \beta.$$

- *Product rule:*

$(fg)' = f'g + fg'$  for all functions  $f$  and  $g$ . As a special case, this rule includes the fact  $(\alpha f)' = \alpha f'$  whenever  $\alpha$  is a constant, because  $\alpha' f = 0 \cdot f = 0$  by the constant rule.

- *Quotient rule:*

$$\left(\frac{f}{g}\right)' = \frac{f'g - fg'}{g^2} \text{ for all functions } f \text{ and } g \text{ at all inputs where } g \neq 0.$$

- *Chain rule* for composite functions: If  $f(x) = h(g(x))$ , then

$$f'(x) = h'(g(x)) \cdot g'(x).$$

[Derivative](#)

### Gradients

#### derivatives that take in multiple variables and produce a real value output

### Gradients

$$f : \mathbb{R}^n \mapsto \mathbb{R}$$

$$\nabla f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_n}(p) \end{bmatrix}$$

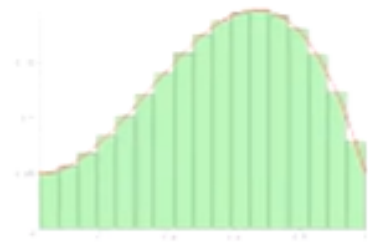
**Example: Gradient**

$$f(x_1, x_2, x_3) = x_1 x_2^2 + \log(x_3)$$

$$\nabla f(a) = \begin{bmatrix} a_2^2 \\ 2a_1 a_2 \\ \frac{1}{a_3} \end{bmatrix}$$

**Derivatives as  
Approximators****Derivatives as Approximators**

$$f(x + h) \approx f(x) + \nabla f(x)^T h$$

**Integration****Integration as a Limit of Sums**

$$\sum_{i=1}^n f(t_i) \Delta_i$$

**Integration in Multiple Dimensions**

$$\int_A f(x, y) \boxed{d(x, y)}$$

$$\int_A f(x) dx, x \in \mathbb{R}^n$$

Multiple Integral



Double integral as volume under a surface  $z = 10 - \left(\frac{x^2}{8} + \frac{y^2}{8}\right)$ . The rectangular region at the bottom of the body is the *domain* of integration, while the **surface** is the graph of the two-variable function to be integrated.

**Antiderivative**

integration and differentiation are inverse operations

**Antiderivative**

$$\frac{d}{dx} f(x) = g(x), \text{ then } \int g(x) dx = f(x) + C .$$

|                               |   |
|-------------------------------|---|
| <b>Live Session</b>           |   |
| <b>Introduction</b>           | <b>05 Jan 2026</b>  |
| <b>Instructor</b>             | Rei Sanchez-Arias   |
| <b>Email</b>                  | <a href="mailto:reisanar@unc.edu">reisanar@unc.edu</a>  |
| <b>Website</b>                | <a href="https://www.reisanar.com/">https://www.reisanar.com/</a>   |
| <b>Office Hours</b>           | Mondays 12:00 pm to 1:00 pm   |
| <b>What to expect</b>         | <b>DATA780</b>  |
| <b>Meeting Time</b>           | Monday 6:00 pm to 7:30 pm   |
| <b>Final Project</b>          | <p>project deliverables:</p> <ul style="list-style-type: none"> <li>• project writeup: ~8 pages NeruIPS format</li> <li>• open-source repository with executable code for methods developed</li> <li>• the github repo can be private or <b>public</b></li> <li>• spotlight presentation (last live session) <ul style="list-style-type: none"> <li>○ 6 to 7 minutes for presentation</li> <li>○ 2 minutes for Q&amp;A</li> <li>○ prepared slide deck</li> <li>○ optional live demo</li> </ul> </li> <li>• final project can be combined for DATA780 and DATA740 (if it applies to both)</li> <li>• your project may innovate some new LM methodology, or make an improvement to an existing methodology</li> </ul> |
| <b>Final Project Proposal</b> | <ul style="list-style-type: none"> <li>• state the task, goals</li> <li>• what methods do you plan to use?</li> <li>• what datasets will you consider?</li> <li>• evaluations metrics you plan to use</li> </ul>  |
| <b>Assessments</b>            | Quizzes will be weekly  |
| <b>Syllabus</b>               | <a href="https://digitalcampus.instructure.com/courses/55373/assignments/syllabus">https://digitalcampus.instructure.com/courses/55373/assignments/syllabus</a>   |
| <b>Live Session Summary</b>   | <p>Meeting Notes found at the top of the Modules page on canvas</p> <p><a href="https://digitalcampus.instructure.com/courses/55373/pages/meeting-notes?module_item_id=9295727">https://digitalcampus.instructure.com/courses/55373/pages/meeting-notes?module_item_id=9295727</a></p>  |

| What to expect  |   | DATA780 |
|---|---|---------|
| <b>Resources<br/>(free pdf<br/>downloads<br/>available)</b> | <ul style="list-style-type: none"> <li>• <b>Mathematics for Machine Learning</b> <ul style="list-style-type: none"> <li>○ <a href="https://mml-book.github.io/book/mml-book.pdf">https://mml-book.github.io/book/mml-book.pdf</a></li> </ul> </li> <li>• <b>Data-Driven Science and Engineering</b> <ul style="list-style-type: none"> <li>○ <a href="https://datasciencebook.com/">https://datasciencebook.com/</a></li> </ul> </li> <li>• <b>Introduction to Statistical Learning</b> <ul style="list-style-type: none"> <li>○ <a href="https://www.statlearning.com/">https://www.statlearning.com/</a></li> </ul> </li> </ul> |         |
| <b>Blogs/Websites</b>                                       | <b>3Blue1Brown Videos (by Grant Sanderson)</b> <ul style="list-style-type: none"> <li>• Linear Algebra Series           <ul style="list-style-type: none"> <li>○ <a href="https://www.3blue1brown.com/topics/linear-algebra">https://www.3blue1brown.com/topics/linear-algebra</a></li> </ul> </li> <li>• Calculus Series           <ul style="list-style-type: none"> <li>○ <a href="https://www.3blue1brown.com/topics/calculus">https://www.3blue1brown.com/topics/calculus</a></li> </ul> </li> </ul>   |         |
| <b>TODO</b>   |   |         |
| <b>High Priority</b>  | <ul style="list-style-type: none"> <li>• review linear algebra terminology</li> <li>• <b>finish CW-Data780-unit_01.ipynb</b></li> </ul>   |         |
| <b>Low Priority</b>   | <ul style="list-style-type: none"> <li>• review assignments           <ul style="list-style-type: none"> <li>○ <b>Homework 1</b></li> <li>○ <b>Final Project</b></li> </ul> </li> </ul>   |         |