# CSX415
## Data Science Principals and Practice

Christopher Brown

U.C. Berkeley / Decision Patterns LLC

Spring 2018

Berkeley
UNIVERSITY OF CALIFORNIA

# What is Machine Learning?

**Today's Question:**

**What three things do all Machine Learning Algorithms have in common?**

MACHINE INTELLIGENCE 3.0

# What is **machine learning?**

# A formal *process* for building a *model*

**Other names for ML:**
Machine Learning
Artificial Intelligence
Statistical Learning
Pattern Recognition
Data Mining
Predictive Analytics
Knowledge Discovery
Predictive Modeling
Model Induction
…

# What is a **model?**

a ***function (f)***

... that *estimates* a ***response*** $(\hat{y})$
... associated with (a set of) known ***predictors*** $(x)$

# What is a function?

$$\hat{y} = f(x_1, x_2, \dots, x_n)$$

$$\hat{y} = f(\vec{x})$$

$$\vec{x} \xrightarrow{f} \hat{y}$$

"maps"

Independent
variables, covariates
predictors, attribute,
descriptor, **feature**
...

Prediction,
Forecast,
Estimate,
...

# How do we find $f$?

# Model Training: Finding $f$



**Model Training**

1. Frame Problem
2. Collect and Shape Data
3. Exploratory Data Analysis (EDA)
4. Feature Generation
5. Train Model
6. Validate Model

Iterate

Feature Specification

Model Specification

# How do we use $f$?

There are **two major** ways to distinguish ML problems ... both are determined by $y$

i.e. by what we are trying to *learn*

1.  **Based on availability of $y$**

**Are there previous/historical observations to learn from?**

Yes → **SUPERVISED LEARNING**

No → **UNSUPERVISED LEARNING**

# Not Necessarily Binary …
## There are *special (edge?)* cases

**SPECIAL CASE 1**

Only some $y$'s are known

- and/or -

$y$'s are not directly known inferable

**SEMI-SUPERVISED Learning**

# SPECIAL CASE 2

$y$ 's change during training/scoring

-and/or-

$y$ 's become available during training/scoring

**ADAPTIVE REINFORCEMENT Learning**

\* Less commonly, more frequently "adversary learning"

**2. Based on the type of $y$**

**What values can $y$ assume ?**

Continuous → **Regression**
(predict an count or amount)

Categorical* → **Classification**
(predict a class or category)

*Binary classification is an important special case

# Not Necessarily Binary …
# There are *special (edge?)* cases

# SPECIAL CASE 1

ORDINAL RESPONSE

**Use Either Regression or Classification**

# SPECIAL CASE 2

Date

**Use Either Regression or Classification**

**-or-**

**Special Techniques**

**(forecast|survival)**

$$y$$

Dependent variable,
Target (variable),
Outcome, Response,
**Class (classification)**

Known

# Data Uses

**Dependent variable**,
Target (variable),
Outcome, **Response**,
**Class (classification)**

Independent variables, covariates
predictors, attribute,
descriptor, **feature**
…

Unit of
observation,
Cases,
Instance,
Data Point,
Sample

| Y | $X_1$ | $X_2$ | $X_3$ | … $X_n$ |
|---|---|---|---|---|
| | data, data, data, data, data, data, data, data | | | |
| | data, data, data, data, data, data, data, data | | | |
| | data, data, data, data, data, data, data, data | | | |
| | data, data, data, data, data, data, data, data | | | |
| | data, data, data, data, data, data, data, data | | | |
| | data, data, data, data, data, data, data, data | | | |
| | data, data, data, data, data, data, data, data | | | |
| | data, data, data, data, data, data, data, data | | | |
| | data, data, data, data, data, data, data, data | | | |

*A major limitation of ML is:*

*(nearly) every ML algorithm expects data in a tabular form.*

# Task Breakdown



Other (Misc)

Visualization
Presentation
Formating
Deployment &
Delivery

Data Retrieval,
Management
and
Organization

Statistical
Operations

# Now what about ... $f$?

# How do we find $f$ ?

# Well what properties should $f$ have?

# Desirable Properties of $f$?

- Takes a one or more inputs
- Yields a single output value for each input
- Should be easy* to evaluate
- Outputs, $\hat{y}$, should be "close to" observed values, $y$:

$$\hat{y} \sim y$$

* Computational cheap/efficient

# What do we mean by "Close to"?

**qualitative measure of "close to"?**

…

$$f(\widehat{y}, y)$$

$$\mathcal{L}(\widehat{y}, y)$$

**How do we calculate $\mathcal{L}(\widehat{y}, y)$?**

**Depends on whether we are doing regression or classification**

# qualitative measure of "close to"?

...

**Depends on whether we are doing regression or classification**

# Regression

...

$$\mathcal{L}(\widehat{y}, y) = y - \widehat{y}$$

$$(y - \widehat{y}) = 0$$

...

**That's just one observation**

**We need to evaluate**
$$\mathcal{L}(\widehat{y}, y) = y - \widehat{y}$$
**for all pairs**

And arrive at a single value, we need:
$$L\big(\mathcal{L}(\widehat{y}, y)\big) = (L \ o \ \mathcal{L})(\widehat{y}, y)$$

# Our Model

"Naïve" model

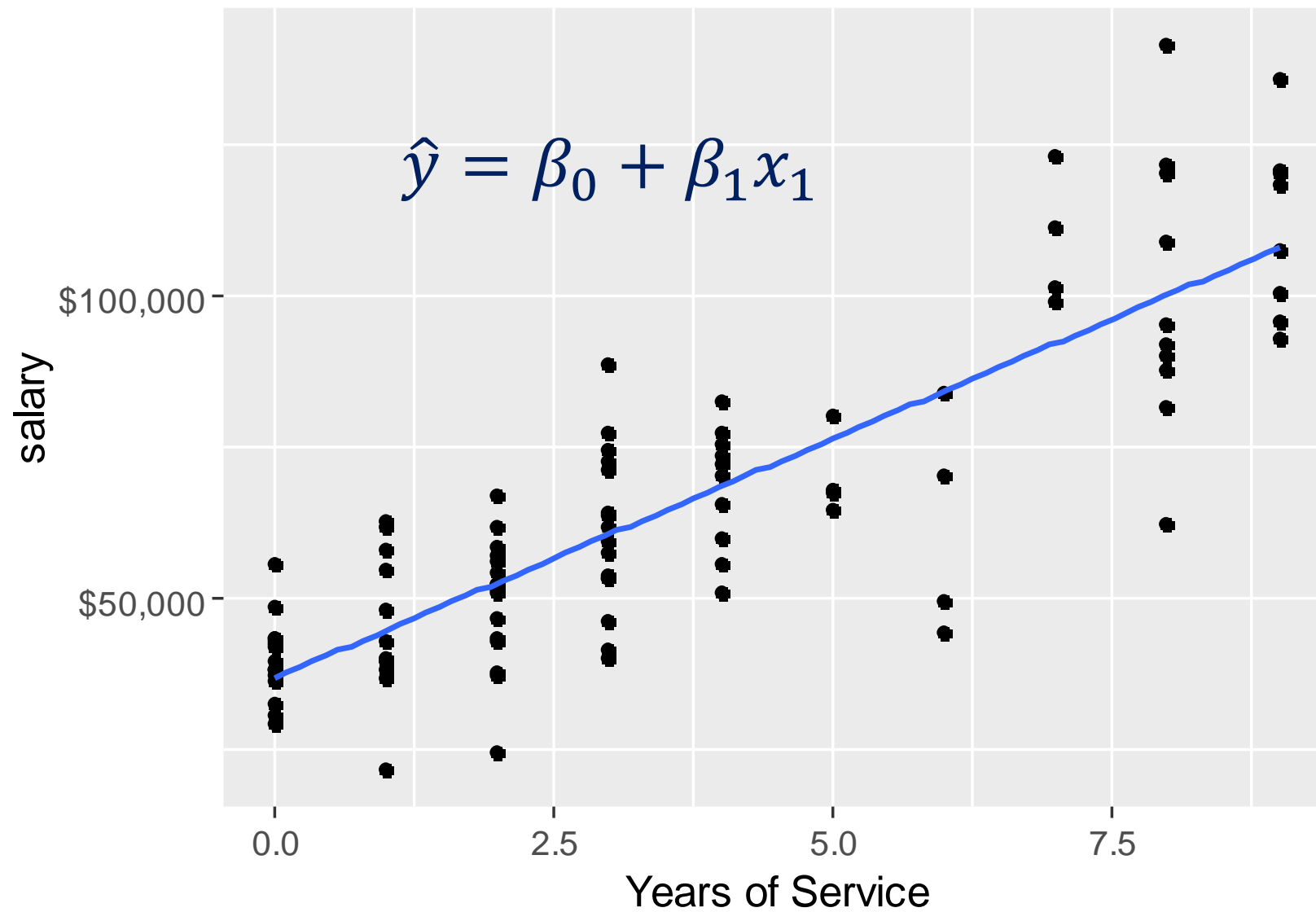$$\hat{y} = mean(y)$$

Linear functions (of one variable)

$$\hat{y} = b + mx$$

$$\hat{y} = \beta_0 + \beta_1 x_1$$

$$\hat{y} = \beta_0 + \beta_1 x_1$$

# Classification

$$\mathcal{L}(\widehat{y}, y) = \begin{Bmatrix} 0 & |y = \widehat{y} \\ 1 & |y \neq \widehat{y} \end{Bmatrix}$$

$$L\big(\mathcal{L}(\widehat{y}, y)\big) = (L \ o \ \mathcal{L})(\widehat{y}, y)$$

# What functions $f$ can be used?

$$\infty$$

# Search / Optimization

Find the **parameters** ($\beta$) that minimize that minimize the loss function …

<span style="color:red">SOLVE:</span>

$$\hat{y} = \beta_0 + \beta_1 x_1$$

$$argmin_\beta\ L(\boldsymbol{y}, \widehat{\boldsymbol{y}})$$

$$argmin_\beta \sum (\boldsymbol{y} - \widehat{\boldsymbol{y}})^2 \text{ (SSE)}$$

**Solution Methods**

• Direct Solution (special case)

• Numerical optimization; recursive goal seeking

# 3 Requirement for ML Algorithm

- A method for evaluating how well the algorithm performs (**ERRORS**)

- A restricted class of functions (**MODEL**)

- A process for proceeding through the restricted class of functions to identify the functions (**SEARCH/OPTIMIZATION**)
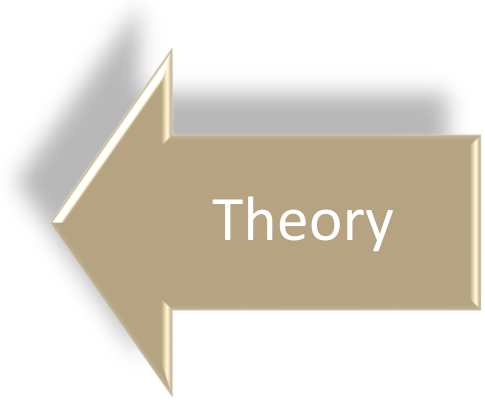
# How to understand Algorithms

1. Errors

2. Model

3. Search Optimization

* Strengths / Limitations

**Frame** problems to make the suitable for solution via machine learning

**Distinguish** fundamental aspects of machine learning algorithms →**know** what algorithms are appropriate for which problems

**Measures/evaluate** model performance

Know how to **improve** a model **and** determine when the model is good enough


Theory

**Execution** is more than building/training models:

**Deploying** machine learning models to operations

**Generating** high quality, graphical and textual results regarding model behavior

**Collaborating** in a group using tools for collaborative/social programming