

# Supervised Learning

## Introduction & Framework

CMPUT 261: Introduction to Artificial Intelligence

P&M §7.1-7.3

# Logistics & Assignment #2

- **Assignment #2 will be released today**  
See eClass
- Due **Tuesday, October 25** at 11:59pm
  - Submissions past the deadline will have late penalty applied
  - Leave yourself some margin for error when submitting!

# Recap: Uncertainty

- We represent uncertainty about the world by **probabilities**
  - We update our knowledge by **conditioning** on **observations**
  - Observations = learning the value of a **random variable**
- Full, **unstructured joint distributions** are intractable to reason about
- **Conditional independence** is a kind of **structure** that is:
  1. widespread
  2. easy to reason about
  3. allows tractable **inference** (computing distribution of unobserved variables)
- **Belief networks** let us compactly represent joint distributions with a lot of conditional independence
  - **Variable elimination** is an algorithm for efficient inference on belief networks
- **Causal models** allow us to infer the consequences of **interventions**

# Supervised Learning, informally

- In the uncertainty section, we took the probability distribution as **given**
  - Our only problem was to represent and derive distributions
- **Question:** Where do these probabilities **come from**?
- **Supervised learning** is a way to learn probabilities from **examples**
  - Probability of a **target** feature (or **label**) given **input features**
  - i.e., **condition** on input features to get probability of target
- Basic idea:
  - Take a bunch of inputs (e.g., images) and "correct" outputs
  - Learn a model that correctly maps inputs to outputs

# Lecture Outline

1. Recap & Logistics
2. Supervised Learning Problem
3. Measuring Prediction Quality

*After this lecture, you should be able to:*

- define supervised learning task, classification, regression, loss function
- represent categorical target values in multiple ways (indicator variables, indexes)
- define generalization performance
- identify an appropriate loss function for different tasks
- explain why a separate test set estimates generalization performance
- define 0/1 error, absolute error, (log-)likelihood loss, mean squared error, worst-case error
- define trivial predictors and explain why they are useful

# Supervised Learning

**Definition:** A *supervised learning task* consists of

- A set of **input features**  $X_1, \dots, X_n$
- A set of **target features**  $Y_1, \dots, Y_k$
- A set of **training examples**, sampled randomly from some population, for which *both* input and target features are given
- A set of **test examples**, sampled from the **same population**, for which *only* the input features are given

The goal is to **predict** the values of the **target features** given the **input features**; i.e., **learn** a function  $h(x)$  that will map features  $X$  to a prediction of  $Y$

- **Classification:**  $Y_i$  are **discrete**
- **Regression:**  $Y_i$  are **real-valued**

# Supervised Learning Examples

1. **Computational vision:** Given example images and labels representing objects, output a label for the main object in the image
  - *Input features:* Pixel values of the image
  - *Target features:* One feature for each label (e.g., **dog**, **plane**, etc.)
2. **Precision medicine:** Given examples of symptoms, test results, and treatments, output an estimate of recovery time
  - Input features: symptoms, treatment indicators, test results, demographic information
  - Target features: recovery time, survival time, etc.
3. **Natural language processing:** Given example sentences and labels representing "sentiment", output how positive or negative the sentence is
  - Input features: binary indicators for words or characters (\*\*!)
  - Target features: One feature per label (e.g., **positive**, **negative**)

# Regression Example

- Aim is to predict the value of **target**  $Y$  based on **features**  $X$
- Both  $X$  and  $Y$  are **real-valued**
  - **Exact values** of both targets and features may not have been in the training set
  - $e_8$  is an **interpolation** problem:  $X$  is **within the range** of the training examples' values
  - $e_9$  is an **extrapolation** problem:  $X$  is **outside** the range of the training examples' values

Ex.	X	Y
$e_1$	0.7	1.7
$e_2$	1.1	2.4
$e_3$	1.3	2.5
$e_4$	1.9	1.7
$e_5$	2.6	2.1
$e_6$	3.1	2.3
$e_7$	3.9	7

$e_8$	2.9	?
$e_9$	5.0	?



# Data Representation

- For **real-valued** features, we typically just record the feature values
- For **discrete** features, there are multiple options:
  - **Binary features:** Can code  $\{false, true\}$  as  $\{0, 1\}$  or  $\{-1, +1\}$
  - Can record **numeric** values for each possible value
    - **Cardinal values:** **Differences** are meaningful (e.g., 1, 2, 7)
    - **Ordinal values:** **Order** is meaningful (e.g., *Good, Fair, Poor*)
    - **Categorical values:** **Neither** differences nor order meaningful (e.g., *Red, Green, Blue*)
  - Vector of **indicator variables:** One per feature value, exactly one is true (sometimes called a "one-hot" encoding) (e.g., *Red* as (1, 0, 0), *Green* as (0, 1, 0), etc.)

# Classification Example: Holiday Preferences

- An agent wants to learn a person's preference for the **length** of holidays
- Holiday can be for 1,2,3,4,5, or 6 days
- Two possible representations:

Ex.	Y
$e_1$	1
$e_2$	6
$e_3$	6
$e_4$	2
$e_5$	1

Ex.	Y <sub>1</sub>	Y <sub>2</sub>	Y <sub>3</sub>	Y <sub>4</sub>	Y <sub>5</sub>	Y <sub>6</sub>
$e_1$	1	0	0	0	0	0
$e_2$	0	0	0	0	0	1
$e_3$	0	0	0	0	0	1
$e_4$	0	1	0	0	0	0
$e_5$	1	0	0	0	0	0

## Question:

What are the advantages/  
disadvantages of  
each representation?

# Generalization

- **Question:** What does it mean for a trained model to **perform well**?
- We want to be able to make correct predictions on **unseen** data, not just the training examples
  - We are even willing to sacrifice some **training** accuracy to achieve this
  - We want our learners to **generalize**: to go beyond the given training examples to classify **new examples** well
  - **Problem:** We can't measure performance on unobserved examples!
- We can **estimate** generalization performance by evaluating performance on the **test set (Why?)**
  - The learning algorithm doesn't have access to the test data, but we do

# Generalization Example

**Example:** Consider binary two classifiers, **P** and **N**

- **P** classifies all the **positive examples** from the training data as **true**, and all others as **false**
- **N** classifies all of the **negative examples** from the training data as **false**, and all others as **true**

**Question:** Which classifier performs better on the **training data**?

**Question:** Which classifier **generalizes** better?

# Bias

- The **hypothesis** is the function  $h(X)$  that we learn
- The **hypothesis space** is the set of **possible hypotheses**
  - "Training a model" =  
"Choosing a hypothesis from the hypothesis space based on data"
- A preference for one hypothesis over another is called **bias**
  - Bias is not a bad thing in this context!
  - Preference for "simple" models is a bias
  - Which bias works best for **generalization** is an **empirical** question

# Learning as Search

- Given **training data**, a **hypothesis space**, an **error measurement**, and a **bias**, learning can be reduced to **local search**
- Learning searches the hypothesis space trying to find the hypothesis that best fits the data given the bias
  - Search space is prohibitively **large** (typically infinite)
  - Almost all supervised learning methods are versions of **local search**

# Measuring Prediction Error

- We choose our hypothesis partly by measuring its **performance** on training data
  - **Question:** What is the other consideration?
- This is usually described as **minimizing** some quantitative measurement of **error** (or **loss**)
  - **Question:** What might error mean?



# 0/1 Error

## Definition:

The **0/1 error** for a dataset  $E$  of examples and hypothesis  $\hat{Y}$  is the number of examples for which the prediction was not correct:

$$\sum_{e \in E} 1 [Y(e) \neq \hat{Y}(e)]$$

- Not appropriate for **real-valued** target features (**why?**)
- Does not take into account **how wrong** the answer is
  - e.g.,  $1 [2 \neq 1] = 1 [6 \neq 1]$
- Most appropriate for **binary** or **categorical** target features

$1 [ \cdot ]$  is **indicator function**:  
value is 1 if the expression  
in brackets is TRUE, else 0



# Absolute Error

## Definition:

The **absolute error** for a dataset  $E$  of examples and hypothesis  $\hat{Y}$  is the sum of absolute distances between the predicted target value and the actual target value:

$$\sum_{e \in E} |Y(e) - \hat{Y}(e)|$$

- Meaningless for **categorical** variables
- Takes account of **how wrong** the predictions are
- Most appropriate for **cardinal** or *possibly* **ordinal** values

# Squared Error

## Definition:

The **squared error** (or sum of squares error or mean squared error) for a dataset  $E$  of examples and hypothesis  $\hat{Y}$  is the sum of squared distances between the predicted target value and the actual target value:

$$\sum_{e \in E} \left( Y(e) - \hat{Y}(e) \right)^2$$

- Meaningless for **categorical** variables
- Takes account of **how wrong** the predictions are
  - **Large** errors are much more important than **small** errors
- Most appropriate for **cardinal** values

# Worst-Case Error

## Definition:

The **worst-case error** for a dataset  $E$  of examples and hypothesis  $\hat{Y}$  is the maximum absolute difference between the predicted target value and the actual target value:

$$\max_{e \in E} |Y(e) - \hat{Y}(e)|$$

- Meaningless for **categorical** variables
- Takes account of **how wrong** the predictions are
  - but only on **one example**  
(the one whose prediction is furthest from the true target)
- Most appropriate for **cardinal** values

# Probabilistic Predictors

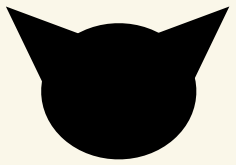
- Rather than predicting **exactly** what a target value will be, many common algorithms predict a **probability distribution** over possible values
  - Especially for **classification** tasks
- Vectors of **indicator variables** are the most common data representation for this scheme:
  - Target features of **training** examples have a single 1 for the **true** value
  - **Predicted** target values are **probabilities** that sum to 1

# Probabilistic Predictions Example

Training examples

<b>X</b>	<b>Y<sub>cat</sub></b>	<b>Y<sub>dog</sub></b>	<b>Y<sub>panda</sub></b>
	1	0	0
	0	1	0

Output on test example

<b>X</b>	<b><math>\hat{Y}_{\text{cat}}</math></b>	<b><math>\hat{Y}_{\text{dog}}</math></b>	<b><math>\hat{Y}_{\text{panda}}</math></b>
	0.5	0.45	0.05

# Likelihood

- For **probabilistic** predictions, we can use **likelihood** to measure the performance of a learning algorithm

## Definition:

The **likelihood** for a dataset  $E$  of examples and hypothesis  $\hat{Y}$  is the **probability** of independently observing the examples according to the probabilities assigned by the **hypothesis**:

$$\Pr(E \mid \hat{Y}) = \prod_{e \in E} \hat{Y}_{Y(e)}(e).$$

- This has a clear Bayesian interpretation
- **Numerical stability issues:** product of probabilities shrinks **exponentially!**
  - *Example:* Probability of **any** sequence of 5000 coin tosses has probability  $2^{-5000}$ !
  - Floating point underflows almost immediately  
(double-precision floating point can't represent anything smaller than  $2^{-1021}$ )

# Log-Likelihood

## Definition:

The **log-likelihood** for a dataset  $E$  of examples and hypothesis  $\hat{Y}$  is the **log-probability** of independently observing the examples according to the probabilities assigned by the hypothesis:

$$\begin{aligned}\log \Pr(E | \hat{Y}) &= \log \prod_{e \in E} \hat{Y}_{Y(e)}(e) \\ &= \sum_{e \in E} \log \hat{Y}_{Y(e)}(e)\end{aligned}$$

- Taking log of the likelihood fixes the underflow issue (**why?**)
- The log function grows **monotonically**, so maximizing log-likelihood is the **same thing** as maximizing likelihood:

$$\left( \Pr(E | \hat{Y}_1) > \Pr(E | \hat{Y}_2) \right) \iff \left( \log \Pr(E | \hat{Y}_1) > \log \Pr(E | \hat{Y}_2) \right)$$

# Trivial Predictors

- The simplest possible predictor **ignores all input features** and just predicts the **same value**  $v$  for any example
- **Question:** Why would we every want to think about these?



# Optimal Trivial Predictors for Binary Data

- Suppose we are predicting a **binary** target
- $n_0$  **negative** examples
- $n_1$  **positive** examples
- **Question:** What is the optimal single prediction?

Measure	Optimal Prediction
0/1 error	0 if $n_0 > n_1$ else 1
absolute error	0 if $n_0 > n_1$ else 1
squared error	$\frac{n_1}{n_0 + n_1}$
worst case	$\begin{cases} 0 & \text{if } n_1 = 0 \\ 1 & \text{if } n_0 = 0 \\ 0.5 & \text{otherwise} \end{cases}$
likelihood	$\frac{n_1}{n_0 + n_1}$
log-likelihood	$\frac{n_1}{n_0 + n_1}$

# Optimal Trivial Predictor Derivations

0/1 error

0 if  $n_0 > n_1$  else 1

$$L(v) = vn_1 + (1 - v)n_0$$

log-likelihood

$$\frac{n_1}{n_0 + n_1}$$

$$L(v) = n_1 \log v + n_0 \log(1 - v)$$

$$\frac{d}{dv}L(v) = 0$$

$$0 = \frac{n_1}{v} - \frac{n_0}{1 - v}$$

$$\frac{n_0}{1 - v} = \frac{n_1}{v}$$

$$\frac{v}{1 - v} = \frac{n_1}{n_0} \wedge (0 < v < 1) \implies v = \frac{n_1}{n_0 + n_1}$$

# Summary

- **Supervised learning** is learning a **hypothesis** function from training examples
  - Maps from **input** features to **target** features
  - **Classification: Discrete** target features
  - **Regression: Real-valued** target features
- **Preferences** among hypotheses are called **bias**
- Choice of **error measurement (loss)** is an important design decision
- Different losses have different optimal trivial predictors
  - Trivial predictors are a baseline: your real model better outperform the trivial predictor