

## Day 8 - Statistical Learning Framework

Agenda:

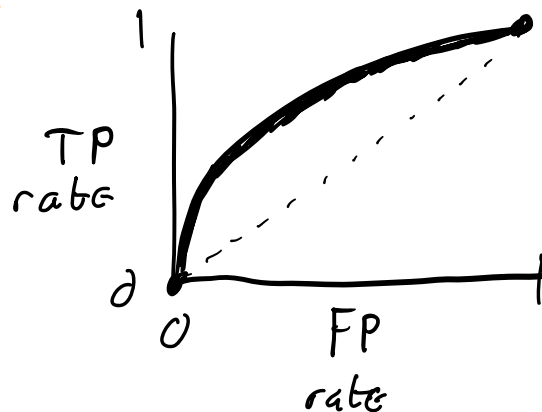
- Statistical learning framework
- Derivation of square loss for regression
- Derivation of log loss / cross-entropy loss for classification
- Terms related to the statistical learning framework
- Bias variance tradeoff

### ROC CURVE

TP rate = recall = sensitivity

= proportion of positive class correctly classified

$$= \frac{TP}{TP+FN}$$



FP rate =  $\frac{FP}{TN+FP}$  = 1 - Specificity

= proportion of negative class incorrectly classified

Precision =  $\frac{TP}{TP+FP}$

= proportion of predicted positives that are correct

		Predicted	
		+	-
True	+	TP	FN
	-	FP	TN

# CS 6140: Machine Learning — Fall 2021 — Paul Hand

HW 3

Due: Wednesday October 6, 2021 at 2:30 PM Eastern time via [Gradescope](#).

Names: [Put Your Name(s) Here]

You can submit this homework either by yourself or in a group of 2. You may consult any and all resources. Make sure to justify your answers. If you are working alone, you may either write your responses in LaTeX or you may write them by hand and take a photograph of them. If you are working in a group of 2, you must type your responses in LaTeX. You are encouraged to use [Overleaf](#). Create a new project and replace the tex code with the tex file of this document, which you can find on the [course website](#). To share the document with your partner, click Share > Turn on link sharing, and send the link to your partner. When you upload your solutions to Gradescope, make sure to take each problem with the correct page or image.

**Question 1.** *Linear regression with multivariate responses.*

Consider training data  $\{(x^{(i)}, y^{(i)})\}_{i=1 \dots n}$ , where  $x^{(i)} \in \mathbb{R}^d$  and  $y^{(i)} \in \mathbb{R}^k$ . Consider a model  $y = Ax$ , where  $A \in \mathbb{R}^{k \times d}$  is unknown. Estimate  $A$  by solving least squares linear regression

$$Y = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{pmatrix} = A \begin{pmatrix} x^{(1)} \\ \vdots \\ x^{(n)} \end{pmatrix} \quad \min_A \sum_{i=1}^n \|y^{(i)} - Ax^{(i)}\|^2 \quad \min \|Y - AX^t\|^2$$

$k \times n$        $k \times d$        $d \times n$

(a) Find  $A$  in the case of training data  $\left\{ \left( \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right), \left( \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right), \left( \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 2 \\ 3 \end{pmatrix} \right) \right\}$ . You may use a com-

puter to perform linear algebra. Hint: the problem can be simplified by observing that each output dimension can be computed separately from the others. If you use this fact, justify it in your response.

$$Y = AX^t$$

$$y = Ax$$

$$\min_{x_1, x_2} f(x_1) + g(x_2) \quad \min f(x_1) \quad \min g(x_2)$$

**Response:**

(b) Consider the case of generic training data. Let  $Y$  be the  $k \times n$  matrix such that  $Y_{ji} = y_j^{(i)}$ . Let  $X$  be the  $n \times d$  matrix where  $X_{ij} = x_j^{(i)}$ . Provide a formula for the least squares estimate of  $A$ . Make sure to check that the matrix dimensions match in any matrix products that appear in your answer. Use the same hint as in part (a).

**Response:**

(c) Show that any prediction under this learned model is a linear combination of the response values  $(y^{(1)}, \dots, y^{(n)})$ . That is, for the  $A$  in part (b), show that  $Ax \in \text{span}(y^{(1)}, \dots, y^{(n)})$  for any  $x$ . You may assume that  $X$  is rank  $d$ .

**Response:**

**Question 2.** *Logistic Regression*

Consider training data  $\{(x_i, y_i)\}_{i=1 \dots n}$ , where  $x_i \in \mathbb{R}^d$  and  $y_i \in \{0, 1\}$ . Consider the logistic data model  $\hat{y} = \sigma(\theta \cdot x)$ , where  $x \in \mathbb{R}^d$ ,  $\theta \in \mathbb{R}^d$ , and  $\sigma$  is the logistic function  $\sigma(z) = e^z / (e^z + 1)$ .

(a) Show that  $\sigma'(z) = \sigma(z)(1 - \sigma(z))$ .

**Response:**

(b) Let  $f(\theta) = \sum_{i=1}^n -y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)$ , where  $\hat{y}_i = \sigma(\theta \cdot x_i)$ . Compute  $\nabla f(\theta)$ . Use the fact in part (a) to simplify your answer.

**Response:**  $d\theta$   $(\times d)$

(c) If  $M = \sum_{i=1}^n x_i x_i^t$ , show that  $z^t M z \geq 0$  for any  $z \in \mathbb{R}^d$ .

**Response:**  $d \times d$

$x_i^t x_i$  inner product  
 $x_i x_i^t$  outer product

(d) Using a summation and vector and/or matrix products, write down a formula for the Hessian,  $H$ , of  $f$  with respect to  $\theta$ . Show that  $z^t H z \geq 0$  for any  $z \in \mathbb{R}^d$ .

**Response:**

$$H(f) = \frac{\partial^2 f}{\partial \theta_i \partial \theta_j}$$

$i, j = 1 \dots d$

$$\nabla f = \left\{ \frac{\partial f}{\partial \theta_i} \right\}_{i=1 \dots d}$$

## Statistical Framework for ML (supervised)

Assume:

- $(x_i, y)$  are sampled from a joint probability distribution
- Training data  $D = \{(x_i, y_i)\}_{i=1 \dots n}$  are iid samples
- Test data are also iid samples OF THE SAME DISTRIBUTION!

Can estimate the model/predictor by maximum likelihood estimation

Results (usually) in an optimization problem

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{H}} \sum_{i=1}^n \ell(f(x_i), y_i) \quad \text{"empirical risk minimization"}$$

where

$\ell$  - loss function eg  $\ell(\hat{y}, y) = |\hat{y} - y|^2$

$\mathcal{H}$  - hypothesis class eg degree  $d$  polynomial

Evaluate performance on test data  $\{(x_i, y_i)\}_{i=1 \dots m}$

$$\frac{1}{m} \sum_{i=1}^m \ell(y_i, \hat{h}(x_i))$$

## Linear Regression and Square Loss

Let  $a \in \mathbb{R}^d$ ,  $x \in \mathbb{R}^d$

Model:  $y_i = x_i^t a + \varepsilon_i$  w/  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$

Data:  $\mathcal{D} = \{(x_i, y_i)\}_{i=1 \dots n}$

Estimate  $a$  by maximum likelihood

pdf of  $\varepsilon_i$  is  $\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{z^2}{2\sigma^2}}$  over  $z \in \mathbb{R}$

likelihood of data (using  $\varepsilon_i = y_i - x_i^t a$ )

$$L(a) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - x_i^t a)^2}{2\sigma^2}} \quad \text{by independence of data}$$

$$\log L(a) = -\sum_{i=1}^n \frac{(y_i - x_i^t a)^2}{2\sigma^2} + \text{terms constant in } a$$

maximizing data likelihood  $\Leftrightarrow$  minimizing square loss

$$\max_a L(a) \quad \Leftrightarrow \quad \min_a \underbrace{\sum_{i=1}^n (x_i^t a - y_i)^2}_{\text{Square loss } \ell(\hat{y}, y) = |\hat{y} - y|^2}$$

# Logistic Regression and Cross Entropy Loss

Model:

Let  $a \in \mathbb{R}^d$

$$P(y=1|x) = \sigma(x^t a)$$

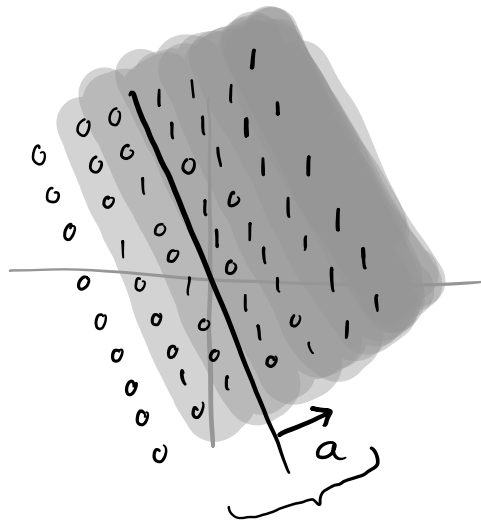
$$P(y=0|x) = 1 - \sigma(x^t a)$$

$$\text{w/ } \sigma(z) = \frac{e^z}{e^z + 1} = \frac{1}{1 + e^{-z}}$$



Data:  $\{(x_i, y_i)\}$

Visually:



width of region of uncertainty  $\approx \frac{1}{\|a\|_2}$

Estimate  $a$  by maximum likelihood

$$L(a) = \prod_{i=1}^n P(y_i=0|x_i)^{1-y_i} P(y_i=1|x_i)^{y_i}$$

$$\log L(a) = \sum_{i=1}^n (1-y_i) \log P(y_i=0|x_i) + y_i \log P(y_i=1|x_i)$$

Cross entropy loss

$$\mathcal{L}_{CE}(P, q) = - \sum_{z \in \mathcal{Z}} P(z) \log q(z) = - \mathbb{E}_P(\log q)$$

discrete  
r.v.s over  $\mathcal{Z}$

$\sum_z [\log q(z)] p(z)$

Maximizing data likelihood  $\Leftrightarrow$  minimizing cross entropy loss

$$\max_a L(a) \Leftrightarrow \min_a - \sum_{i=1}^n \left( y_i \log(\sigma(x_i^T a)) + (1-y_i) \log(1-\sigma(x_i^T a)) \right)$$

$\mathcal{L}_{CE} \left( \begin{pmatrix} y_i \\ 1-y_i \end{pmatrix}, \begin{pmatrix} \sigma(x_i^T a) \\ 1-\sigma(x_i^T a) \end{pmatrix} \right)$

## Formalism for Statistical Framework for ML (supervised)

---

**Domain Set** -  $X$  - arbitrary set of objects/instances that could be labelled

- usually represented as a feature vector in  $\mathbb{R}^d$
- could be infinite dimensional

**Label Set** -  $Y$  - set of possible labels

- eg.  $\mathbb{R}^d$  for regression
- $\{1,0\}$  for binary classification
- Finite set for multiclass classification

**Training data** -  $S = \{(x_i, y_i)\}_{i=1 \dots n}$   
n points in  $X \times Y$

**Predictor/hypothesis** - any function  $h: X \rightarrow Y$  that  
 $x \mapsto y$   
outputs a prediction  $y$  for any instance  $x$

**Hypothesis Class** -  $H$  a set of predictors/hypotheses that are being considered  
eg  $H = \{\text{degree } d \text{ polynomials}\}$



## Data generation model

### Simple version

- Assume  $x \sim D$ , where  $D$  is a <sup>probability</sup> distribution over  $X$
- Each sample is independent
- $y = f(x)$  for a "correct" function  $f$ .

### Realistic version

- Assume  $(x, y) \sim D$ , a joint probability distribution over  $X \times Y$

There is some marginal distribution of  $X$ ,  $P_X$ .

For any  $x$ , there is a conditional distribution over  $y$   $D_{y|x}$

## Loss

- how bad is the prediction of an instance relative to its label

$$\underset{\substack{\text{label} \\ y}}{\mathcal{L}}(y, \underset{\substack{\text{prediction} \\ \hat{y}}}{\hat{y}}) \in \mathbb{R}$$

Examples

- Square loss  $\mathcal{L}(y, \hat{y}) = \|y - \hat{y}\|^2$  if  $y, \hat{y} \in \mathbb{R}^d$

- log loss  $\mathcal{L}(y, \hat{y}) = \sum_{i=1}^k y_i \log \hat{y}_i$  if  $y \in \mathbb{R}^k$  are one-hot encodings &  $\hat{y} \in \mathbb{R}^k$  is a probability dist over  $k$  labels

- 0-1 loss  $\mathcal{L}(y, \hat{y}) = \begin{cases} 1 & \text{if } \hat{y} \neq y \\ 0 & \text{if } \hat{y} = y \end{cases}$

**Risk** - expected loss of a predictor  
for new data samples

$$R(h) = \mathbb{E}_{(x,y) \sim D} \ell(y, h(x))$$

aka "generalization error"  
"error" "test error"  
"population error"

Goal of learning:

To find a  $h$  such that  $R(h)$   
is minimal. Want to solve

$$\min_{h \in H} R_D(h)$$

challenge: We don't know  $D$ . We only  
have samples  $S$

Empirical Risk  
Minimization

— approximation of risk based  
on training data  $S$

$$\hat{h} = \operatorname{argmin}_{f \in \mathcal{H}} \underbrace{\sum_{i=1}^n \ell(y_i, f(x_i))}_{\text{Empirical risk}}$$

Test Error — Use a finite test set  
to assess generalization

$$\frac{1}{m} \sum_{i=1}^m \ell(y_i^{\text{test}}, \hat{h}(x_i^{\text{test}}))$$