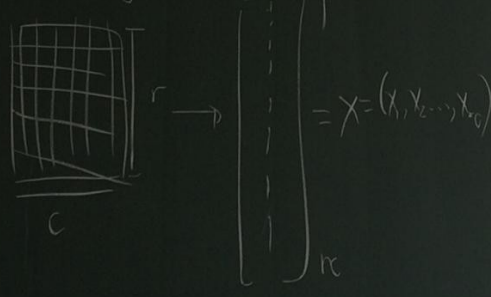


digit recognition

$y = \text{Output}$
 $= \{0, 1, 2, \dots, 9\}$
 $= \text{label of the number}$

$X = \text{Input space}$
 $= \text{Vectorized discretized image of the digit}$



$f: (x_1, x_2, \dots, x_{rc}) \mapsto \{0, 1, 2, \dots, 9\}$

$P\{Y=j|X\}$

(I) Shape of the data (n, p) . $\frac{n}{p} < 1 \Rightarrow \text{Wide}$
 (II) Content of data (Type) $\frac{n}{p} > 1 \Rightarrow \text{Thin}$

Consider the i^{th} obs.
 $X_i = (X_{i1}, X_{i2}, \dots, X_{ip})^T$ p -dim input
 What is X_{ij} ?

(i) $X_{ij} \in \{0, 1\}$ Binary input
 $X_i \in \{0, 1\}^p = X$

(ii) $X_{ij} \in [0, 1]$ Unitary Input
 Neural Networks

$(X_i, Y_i) \text{ iid } p_{XY}(x, y)$

$f: X \rightarrow Y$
 $x \mapsto f(x)$

X : Input space
 Y : Output space


Def. f is a function.
 ref to as a learning machine

f encodes the pattern of relationship (association) b/w X and Y .

$f \in \mathcal{H}$

eg: $Y^X = \text{All possible mappings from } X \text{ to } Y$

\mathcal{H} = function space Hypothesis space



Note: (a) \mathcal{H} may be finite and countable
 $\mathcal{H} = \{f_1, f_2, \dots, f_m\}$

(b) $\mathcal{H} \equiv \text{countably infinite}$

(c) $\mathcal{H} = \text{Infinite}$
 Choosing \mathcal{H} is an act of approximation
eg $\mathcal{H} = \text{Polynomials of degree } m$

ML

Likelihood function and Maximum Likelihood Estimation (MLE)

Refresher: if $Y_i \sim F_Y(y; \theta)$
 then the pdf of Y is $p_Y(y; \theta)$.
 The likelihood of θ is

$$L(\theta; Y) = p(y_1, y_2, \dots, y_n; \theta)$$

$$\stackrel{iid}{=} \prod_{i=1}^n p_Y(y_i; \theta)$$

Thanks to the convexity of $\log(\cdot)$ and given that the logarithm of $L(\theta; Y)$ is easier to manipulate, we use

$$\ell(\theta; Y) = \log L(\theta; Y)$$

$$= \log \text{likelihood of } \theta$$

MLE
 The maximum likelihood estimator of θ is

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in \Theta} \{L(\theta; Y)\}$$

\Rightarrow capped $\theta = \arg \max_{\theta \in \Theta} \{\log L(\theta; Y)\}$

Ex 1: if $\eta(X) = \beta^T X$, then our param is simply β

$$L(\beta; Y) = \prod_{i=1}^n p(y_i | x_i, \beta)$$

where $p(y_i | x_i, \beta) = \frac{1}{1 + e^{-\beta^T x_i}} \left(\frac{e^{-\beta^T x_i}}{1 + e^{-\beta^T x_i}} \right)^{y_i}$

Likelihood function and Maximum Likelihood Estimation (MLE)

Refresher: if $Y_i \sim F_Y(y; \theta)$
 then the pdf of Y is $p_Y(y; \theta)$.
 The likelihood of θ is

$$L(\theta; Y) = p(y_1, y_2, \dots, y_n; \theta)$$

$$\stackrel{iid}{=} \prod_{i=1}^n p_Y(y_i; \theta)$$

Thanks to the convexity of $\log(\cdot)$ and given that the logarithm of $L(\theta; Y)$ is easier to manipulate, we use

$$\ell(\theta; Y) = \log L(\theta; Y)$$

$$= \log \text{likelihood of } \theta$$

MLE
 The maximum likelihood estimator of θ is

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in \Theta} \{L(\theta; Y)\}$$

\Rightarrow capped $\theta = \arg \max_{\theta \in \Theta} \{\log L(\theta; Y)\}$

Ex 1: if $\eta(X) = \beta^T X$, then our param is simply β

$$L(\beta; Y) = \prod_{i=1}^n p(y_i | x_i, \beta)$$

where $p(y_i | x_i, \beta) = \frac{1}{1 + e^{-\beta^T x_i}} \left(\frac{e^{-\beta^T x_i}}{1 + e^{-\beta^T x_i}} \right)^{y_i}$

Focus on $p(y|x)$ posterior distribution of Y given $X=x$

Central to AI and ML
b/c we are typically given $X=x$
and we seek to predict Y given $X=x$

AI: Artificial Intelligence
ML: Machine Learning -

Eg: Digit Rec

- Take a picture of a given digit
- Extract the grayscale 28×28 matrix
- Vectorizer $\begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \rightarrow \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$

Majority rule

$$\hat{f}(x) = \begin{cases} 0 & \text{if } P(Y=0|x) \text{ is largest} \\ 1 & \text{if } P(Y=1|x) \text{ is largest} \\ \vdots & \vdots \\ 9 & \text{if } P(Y=9|x) \text{ is largest} \end{cases}$$

$$\hat{f}(x) = \arg \max_{j=0,1,\dots,9} P(Y=j|x)$$

$\hat{f}(x) = \text{label of } x$

$\tilde{x}_{ij} \rightarrow \frac{x_{ij} - x_{(i)}^j}{x_{(i)}^j - x_{(i)}^1}$

$\tilde{x}_{ij} = \tilde{x}_{ij} (x_{(i)}^j - x_{(i)}^1) + x_{(i)}^1$

(III) Distribution
 $p_{XY}(x,y)$

(IV) Missing Ness
Pattern and incidence of missing values

(V) Outlier and Novelty Selection

Understanding $p_{XY}(x,y)$

- Focus on $p(x)$ ref to density estimation in X key to Anomaly analysis
- Focus on $p(x|y)$ ref estimation of the so-called class conditional densities
Pillar of generative pattern recog via discriminant analysis
- Focus of $p(y)$ is usually the easiest in PR this is simply discrete prob estm

Eg.

Binary Classification

$$y = \{0, 1\}$$

$$X \subset \mathbb{R}^p \quad x = (x_1, \dots, x_p)^T = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} \quad x_j \in \mathbb{R}$$

Hypothesized

Logistic Regr

$$\text{Assume } P\{Y=1|x\} = \frac{1}{1+e^{-\eta(x)}}$$

where $\eta: X \rightarrow \mathbb{R}$

$$\eta(x) = \beta^T x$$

$\eta(x)$ - nonlinear

$$\begin{aligned} P(Y=0|x) &= 1 - P(Y=1|x) \\ &= 1 - \pi(x) \end{aligned}$$

$$Y=0 \text{ or } Y=1$$

$$P(A^c) = 1 - P(A)$$

$$\begin{aligned} p(y|x) &= \pi(x)^y (1-\pi(x))^{1-y} \\ &= \left(\frac{1}{1+e^{\eta(x)}}\right)^y \left(\frac{e^{\eta(x)}}{1+e^{\eta(x)}}\right)^{1-y} \end{aligned}$$

Goal: Estimate $\eta(x)$

(iii) $x_{ij} \in \{a_1, a_2, \dots, a_m\}$

(iv) $x_{ij} \in [a, b] \subset \mathbb{R}$

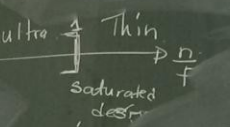
(v) $x_{ij} \in \mathbb{R}^{p \times 1}$

$x_{ij} \equiv$ Frequency of word j in document i
 \equiv Text Mining with Bag of Words assumption.

(vi) $x_{ij} \in \{\text{strongly disagree, disagree, neutral, agree, strongly agree}\}$ 2 values
 \Rightarrow Ordinal Input

x_i	x_j	x_i'	x_j'
100	10	100	10
100	10	100	10
100	10	100	10
100	10	100	10
100	10	100	10
100	10	100	10
100	10	100	10
100	10	100	10
100	10	100	10
100	10	100	10

$$x_{ij} \rightarrow \frac{x_{ij} - \min_j(x_j)}{\text{Range}(x_j)}$$



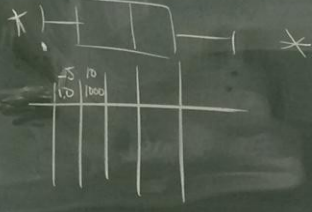
Cubicle

$$x_i \in [a, b]^p$$

$$x_i \in [0, 1]^p$$

$$x_{ij} \rightarrow \frac{x_{ij} - \min_j(x_j)}{\max_j(x_j) - \min_j(x_j)}$$

Outlier: $x^* > Q_3 + 3 IQR$



$p(x, y) = p(x)p(y|x)$
 $= p(y)p(x|y)$

For Thomas: $p(y|x) = \frac{p(y)p(x|y)}{p(x)}$

Posterior evidence \Rightarrow given Normalizer
 Partition f

$E[Y|x] = \int y p(y|x) dy$
 $=$ Condition expect of response Y given predictor X

Regression

Dealing with $p(y|x)$ in PR
 in $E[Y|x]$ in Reg.

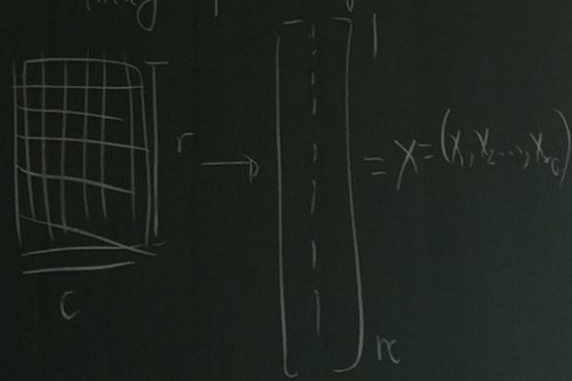
Bayes in Learning

(2) Bayes for estimation θ or $\eta(X)$
 instead of $\hat{\theta}_{MLE} = \arg\max_{\theta} \{ \log L(\theta) \}$
 Construct $p(\theta|D) = \frac{p(\theta)p(D|\theta)}{p(D)}$

Ex: digit recognition

$Y =$ Output
 $= \{0, 1, 2, \dots, 9\}$
 $=$ label of the number

$X =$ Input space
 $=$ Vectorized discretized image of the digit



$f: (x_1, x_2, \dots, x_{100}) \rightarrow \{0, 1, 2, \dots, 9\}$
 $P\{Y=j|X\}$

(I) Shape of the data (n, p) . $\frac{n}{p} < 1 \Rightarrow$ Ultra Thin
 (II) Content of data (Type) $\frac{n}{p} > 1$ Thin

Consider the i^{th} obs.
 $X_i = (x_{i1}, x_{i2}, \dots, x_{ip})$ p -dim input
 What is x_{ij} ?

(i) $x_{ij} \in \{0, 1\}$ Binary input
 $x_i \in \{0, 1\}^p = X$

(ii) $x_{ij} \in [0, 1]$ Unitary Input
 Neural Networks

$p_{XY}(x,y) = p(x)p(y|x)$
 $= p(y)p(x|y)$

For Thomas: $p(y|x) = \frac{p(y)p(x|y)}{p(x)}$

Dealing with $p(y|x)$ in PR ①
 or $E[Y|x]$ in Reg.

Bayes in Learning

Bayes for estimation θ or $\eta(X)$

instead of $\hat{\theta}_{MLE} = \arg \max_{\theta} \{ \log L(\theta; Y) \}$
 Construct $p(\theta | D) = \frac{p(\theta)p(D|\theta)}{p(D)}$

$E[Y|x] = \int y p(y|x) dy$
 = Condition expect of response Y given predictor X

Regression

Posterior Evidence
 \Rightarrow given Normalizer
 $X=x$ Partition

Likelihood function and Maximum Likelihood Estimation (MLE)

Refresher: if $Y \sim F_Y(y|\theta)$
 then the pdf of Y is $p(y|\theta)$
 The likelihood of θ is

$L(\theta; Y) = p(y_1, y_2, \dots, y_n | \theta)$
 $\stackrel{iid}{=} \prod_{i=1}^n p(y_i | \theta)$

Thanks to the convexity of $\log(\cdot)$
 and given that the log-likelihood of $L(\theta)$ is easier to manipulate, we use

$l(\theta; Y) = \log L(\theta; Y)$
 $= \log \text{likelihood of } \theta$

MLE
 The maximum likelihood estimator of θ is

$\hat{\theta}_{MLE} = \arg \max_{\theta \in \Theta} \{ L(\theta; Y) \}$
 $\stackrel{H}{\Rightarrow} \text{optimal } \hat{\theta} = \arg \max_{\theta \in \Theta} \{ \log L(\theta; Y) \}$

[Box] if $\eta(x) = \beta^T x$
 then our param is $\eta(\beta)$

$L(\beta; Y) = \prod_{i=1}^n p(y_i | x_i; \beta)$
 then $p(y_i | x_i; \beta) = \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}}$

Eg.

Binary Classification

$$y = \{0, 1\}$$

$$X \subset \mathbb{R}^p \quad x = (x_1, \dots, x_p)^T \quad x_j \in \mathbb{R}$$

Hypothesized

Logistic Regr

$$\text{Assume } P\{Y=1|x\} = \frac{1}{1+e^{-\eta(x)}}$$

where $\eta: X \rightarrow \mathbb{R}$

$$\eta(x) = \beta^T x \quad \checkmark$$

$\eta(x)$ - nonlinear

$$P(Y=0|x) = 1 - P(Y=1|x) = 1 - \pi(x)$$

$$Y=0 \text{ or } Y=1$$

$$P(A^c) = 1 - P(A)$$

$$p(y|x) = \pi(x)^y (1-\pi(x))^{1-y} = \left(\frac{1}{1+e^{-\eta(x)}}\right)^y \left(\frac{e^{-\eta(x)}}{1+e^{-\eta(x)}}\right)^{1-y}$$

Goal: Estimate $\eta(x)$

$$p(x,y) = p(x)p(y|x) = p(y)p(x|y)$$

$$\text{Bayes' Theorem: } p(y|x) = \frac{p(y)p(x|y)}{p(x)}$$

Posterior
Evidence
Normalizer
Partition

$$E[Y|x] = \int y p(y|x) dy = \text{Condition expect of respn given prediction } x$$

Regression

Dealing with $p(y|x)$ in PR
or $E[Y|x]$ in Reg

Bayes in Learning

(2) Bayes for estimation θ or $\eta(x)$

$$\text{Construct } p(\theta|D) = \frac{p(\theta)p(D|\theta)}{p(D)}$$

instead of $\hat{\theta}_{MLE} = \arg\max_{\theta} \log L(\theta)$

Likelihood function and Maximum Likelihood Estimation (MLE)

Refresher: if $Y_i \sim F_Y(y; \theta)$
 then the pdf of Y is $p(y|\theta)$.
 The likelihood of θ is

$$L(\theta; Y) = \prod_{i=1}^n p(y_i; \theta)$$

Thanks to the convexity of $\log(\cdot)$ and given that the logarithm of $L(\theta; Y)$ is easier to manipulate, we use

$$l(\theta; Y) = \log L(\theta; Y) = \log \text{likelihood of } \theta$$

MLE
 The maximum likelihood estimator of θ is

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in \Theta} \{L(\theta; Y)\}$$

\Rightarrow capped $\theta = \arg \max_{\theta \in \Theta} \{\log L(\theta; Y)\}$

Eg. if $\eta(x) = \beta^T X$, then our param is simply β

$$L(\beta; Y) = \prod_{i=1}^n p(y_i | x_i, \beta)$$

where $p(y_i | x_i, \beta) = \frac{1}{1 + e^{-\beta^T x_i}} \left(\frac{e^{-\beta^T x_i}}{1 + e^{-\beta^T x_i}} \right)^{y_i}$

Focus on $p(y|x)$ posterior density of Y given $X=x$

Central to AI and ML
 b/c we are typically given $X=x$ and we seek to predict Y given $X=x$

AI: Artificial Intelligence
 ML: Machine Learning

Eg. Digit Rec

- Take a picture of a given digit
- Extract the grayscale $r \times c$ matrix
- Vectorizer $\begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \rightarrow \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$

Majority rule

$$\hat{f}(x) = \begin{cases} 0 & \text{if } P(Y=0|x) > \log 2 \\ 1 & \text{if } P(Y=1|x) > \log 2 \\ 9 & \text{if } P(Y=9) > \log 2 \end{cases}$$

$$\hat{f}(x) = \arg \max_{j=0,1,9} P(Y=j|x)$$

$(x_1, y_1) \dots (x_n, y_n)$
 n obs

Confusion matrix

	+1	-1	
+1	TP	FN	
-1	FP	TN	
			n

TN = True negative count
 out of n by f
 TP = True Positive
 count of n by f
 FN = False Neg count

	+1	-1	
+1	0	b	
-1	a	0	

in illness
 $b > a$

When the impact of FP
 is very different from that
 of FN it is advisable
 to a non symmetric loss

Perfect f has $FP = FN = 0$
 $TN + TP = n$

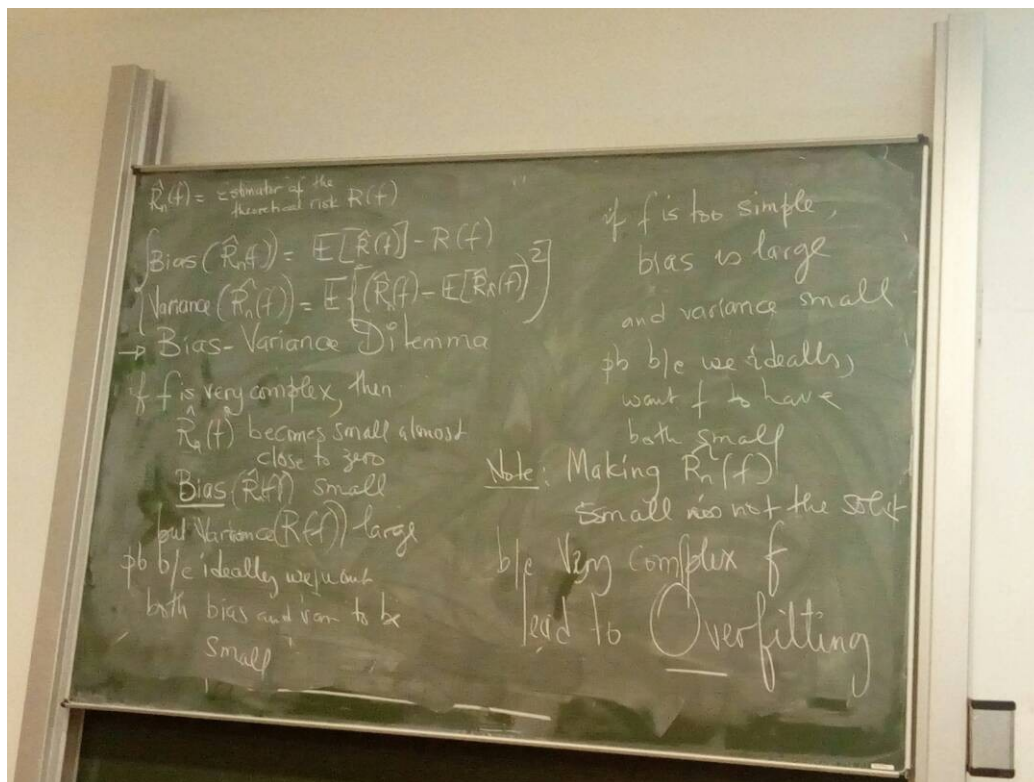
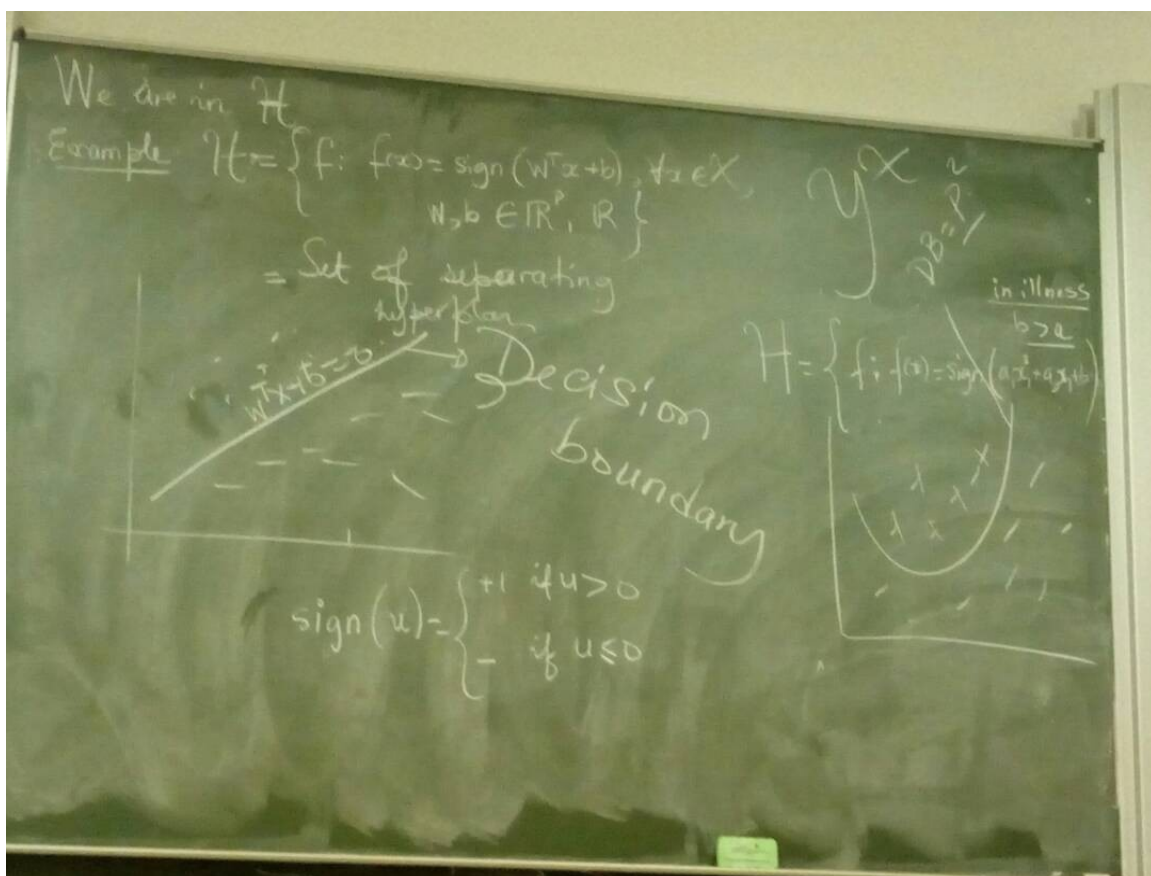
$R_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$
 Empirical Risk
 = Training Error
 = Random Variable
 = Statistic
 = Estimator of
 the theoretical

risk $R(f) = E[\ell(Y, f(X))]$
 $= \int \ell(y, f(x)) p_{xy}(x, y) dx dy$

Note In practice
 $R(f)$ is never known
 b/c $p_{xy}(x, y)$ is unknown

$P(|x - E[x]| < \varepsilon) \geq 1 - \frac{V(x)}{\varepsilon^2}$
 Chebyshev's theorem

$\hat{\theta}$ is a consistent est
 of θ if
 $\lim_{n \rightarrow \infty} P\{\|\hat{\theta} - \theta\| < \varepsilon\} = 1$
 $\hat{\theta} \xrightarrow{P} \theta$



Girls (vs) Ladies