**Lecture 2**   concept $\rightarrow$    $\{(x_i, y_i)\}_{i=1}^{n}$ Data $\longrightarrow$ output $\longrightarrow$ $\boxed{\text{ML}}$ $\longrightarrow$ program $\hat{f}_n$  ;   Data (experience) $\longrightarrow$ task $\rightarrow$ Understanding (performance)

-formance Measures $\rightarrow$   X, Y observation

$f(x) = $ y-prediction

1. example

$loss(Y, f(x)) = \begin{cases} (Y - f(x))^2 & \text{— continuous } \textbf{regression} \\ 0, 1 & \text{— discrete } \textbf{classification} \end{cases}$

$x \in \mathcal{X}, (x, y) \text{ any cell}$

$(x, y) \sim P_{XY}$    $\boxed{\text{Risk} \quad R(f) = \mathbb{E}_{XY}[loss(Y, f(x))]} \quad \longrightarrow (1)$

-ayes Optimal Rule $\rightarrow$   Goal: construct prediction rule   $f^*: \mathcal{X} \rightarrow \mathcal{Y}$

Ideal optimal $\rightarrow f^* = \arg\min_{f} \mathbb{E}_{XY}[loss(Y, f(x))] = \{f \mid \min \mathbb{E}_{XY}[loss(Y, f(x))]\} \rightarrow (2)$

large $x$

-ayes Risk $\rightarrow$   $R(f^*) \leq R(f) \quad \forall f$   we have $\hat{f}$   $\left[\frac{1}{n}\sum_{i=1}^{n}[loss(Y_i, f(x_i))]\right]$

-Training Algorithm $\rightarrow$ Data $= \{(x_i, y_i)\}_{i=1}^{n}$ ,   $\hat{f}_n = $ map from $x$ to $y$ , not overfit

-supervised Learning $\rightarrow$   $\cdot \frac{1}{n}\sum_{i=1}^{n}[loss(Y_i, \hat{f}_n(x_i))]$

---

**Lecture 3**

-er formance Revisit $\rightarrow$

data $D_n = \{(x_i, y_i)\}_{i=1}^{n}$

Excess Risk ($n^r \rightarrow 0$)

$n \rightarrow \infty$

Application of ML

EX1 loan

▷ credit score $\rightarrow$ regression
▷ loan decision $\rightarrow$ classification

-L in Application $\rightarrow$

EX2 : chess
▷ Nature of training sample / Exp
   → Game vs Pro (limited, not much control)
   - Pro's games (nearly unlimited, no control)
   - self vs self (unlimited, flexible)

⌐ Task (T)
├ Training Sample / Exp (E)
├ Type of output (cata, %, number) (O)
├ Performance measure / loss Fn (P)
├ Input (image, credit score) (X)
├ Hypothesis space (H)
│   function   $H: X \rightarrow O$
└ optimize (P) using $H: X \rightarrow O$

---

light GBM , XGBoost    • **supervised** $\Rightarrow$   • **unsupervised**

Octave $\Rightarrow$ free open source Sw
└ Matlab $= \checkmark$ good ML invironment

# Lecture 3

## Linear Regression

**Types of supervised Problems** →

① regression    ② classification

**Linear function** →

$f(x;w) = w_0 + w_1 x_1 + \dots w_d x_d = \vec{W} \cdot \vec{x}$
$= \vec{w}^T x_i$

$x \in x \in R$ line
$x \in R^2$ plane
$x \in R^d$ hyper plane

**Notation** →

- $x_i \in \overset{[N \times d]}{X} \in R^d$  ↳ $d=0$ mean const ⇒ $w = [w_0, w_1, \dots w_d]$
  $x_i = [x_{i0}, x_{i1}, \dots x_{id}]$

$X =$ [d+1 by N matrix]

- $y_i \in y$
- $\bar{X} = N \times (d+1)$ data matrix
- $\bar{y} =$ label vector  $\bar{y} = [y_1, \dots y_N]^T$

$f(x;w) = w_0 + \sum_{j=1}^{d} w_j x_{ij} = W \cdot x_i$

- loss function $\ell(\hat{y}, y)$

**Loss function** →

From function $y = f(\bar{x}; \vec{w})$

$x_0, y_0 =$ new data

$R(w) = E_{(x_0, y_0) \sim p[\bar{x}, y]} [\ell(f(x_0; \vec{w}), y_0)]$

goal = minimize to loss $R(w)$ for new data

$w^* = \arg\min_w \sum_{i=1}^{N} (y_i - W \cdot x_i)^2$

**Least Square** →

to do this

we should minimize $L(w, X, y) = L(w) = E[\ell(y, w, x)] \approx \frac{1}{N} \sum_{i=1}^{N} (y_i - W \cdot x_i)^2$

$$\boxed{\frac{\partial L}{\partial w_j} = \nabla_{w_j} L = 0 \quad \forall j}= \left[ \frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_1}, \dots \frac{\partial L}{\partial w_d} \right] \quad \frac{1}{N}(y - Xw)^T (y - Xw)$$

**Condition 1**

$\frac{\partial L}{\partial w_j} = -\frac{2}{N} \sum_{i=1}^{N} (y_i - W \cdot x_i) x_{ij} = 0 \quad \text{——(1)}$

$= \frac{1}{N}(y^T - X^T w^T)(y - Xw)$

**Condition 2**

Error are uncorrelated w/ data and linear fn

$\sum_{i=1}^{N} (y_i - W \cdot x_i) = 0 \quad \text{——(2)}$

**Derivative of loss**

$L(w) = \frac{1}{N}(y^T - x^T w^T)(y - xw)$

★ From $\frac{\partial a^T b}{\partial a}$ : $\left[ \frac{\partial}{\partial a} \right] \left[ a^T \right] [b]$ : $\left[ \frac{\partial}{\partial a} \right] \left( [b^T][a] \right) : \frac{\partial(b^T a)}{\partial a} : [b]$ ✓

★ matrix $\frac{\partial(a^T B a)}{\partial a_i} = \frac{\partial}{\partial a_i}(a^T (Ba))^T = \frac{\partial}{\partial a_i}((Ba)^T (a^T)^T) = \frac{\partial}{\partial a_i}(a^T B^T a) \leftarrow \boxed{B = B^T}$ condition

dummy index → $B_{ij} a_i \to B_{kj} a_k$, $a_i B_{ik} \to a_i B_{ik}$ condition

เทคนิคเผาลงตรงนี้ทุก

$\frac{\partial(a^T B a)}{\partial a_k} = \frac{\partial}{\partial a_k}\left( \sum_{i,j} a_i B_{ij} a_j \right) = \sum_{kj} a_j B_{kj} \left( \frac{\partial}{\partial a_k} a_k \right) + \sum_{ik} \left( \frac{\partial}{\partial a_k} a_k \right) B_{ik} a_i$

$\frac{\partial(a^T B a)}{\partial a_k} = a_j B_{kj} + B_{ik} a_i \Rightarrow (B^T + B) a$

Derivative of loss 2 $\Rightarrow$

$$\frac{\partial L(w)}{\partial w} = \frac{1}{N}\frac{\partial}{\partial w}\left[v^T y - w^T x^T y - \overset{w^T(y^T)^T}{(y^T x w)} + w^T x^T x w\right]$$

$$= \frac{1}{N}\left[-x^T y - (y^T x)^T + (x^T x + (x^T x)^T)w\right]$$

$$= \frac{1}{N}\left[-2x^T y + 2 x^T x w\right]$$

$(x^+)$

Moore-Penrose pseudoinverse of $x$ $\longrightarrow$

$$\frac{\partial L}{\partial w} = -\frac{2}{N}(x^T y - x^T x w) \cdots ③ = 0$$

$$\boxed{w^* = (x^T x)^{-1} x^T y} \qquad \boxed{x^+ \triangleq (x^T x)^{-1} x^T}$$

$$\text{prediction}: \quad \hat{y} = w^* \cdot x_0 = (y^T x^{+T}) x_0 \qquad \overset{\mathbb{R}}{}\overset{\mathbb{R}^N}{}$$

## Lecture 4   Generalization

More training data $\to$ worse fit
$\hookrightarrow$ better prediction


Test error
Empirical loss

$$f^* = \arg\min_{f: x \to \mathbb{R}} \; E^{New}_{(x_0,y_0)\sim P(x,y)}\left[(f(\tilde{x}_0) - y_0)^2\right] \quad —— ①$$

$$\iint p(x,y)\,dx\,dy = \int\left(\int p(y|\bar{x})\,dy\right)P(x)\,dx \quad —— ②$$

rule of probability: $P(\bar{x},y) = P(y|\bar{x})\,p(\bar{x})$

$$② \to def \quad E_{P(y,\bar{x})}[g(y,\bar{x})] = \iint g(y,\bar{x})\,p(y|\bar{x})\,p(\bar{x})\,dy\,d\bar{x} \quad ③$$

$$E_{P(y,x)}[g(y,x)] = \iint\left\{\int g(y,\bar{x})\,p(y|\bar{x})\,dy\right\}p\bar{x}\,d\bar{x} \quad —— ④$$

notation $\longrightarrow (x_0,y_0)\sim P(x,y) \Rightarrow (x_0,y_0)$ has pb distribution of $P(x,y)$   $E_{x_0\sim p(x)}[\quad]$

$$④ \to E_{(x_0,y_0)\sim P(x,y)}\left[(f(x_0)-y_0)^2\right] = \int\left\{E_{y_0\sim p(y|x)}\left[(f(x_0)-y_0)^2 | x_0\right]\right\}p(x_0)\,dx_0 \quad —— ⑤$$

vary $f(x)$ given $x_0$
$\downarrow$

try to minimize it for each $x_0$
conditional expectation

$$⑤ \quad \frac{\delta}{\delta f(x)}E_{P(y|x)}\left[(f(x_0)-y_0)^2 | x_0\right] = 2E_{y_0\sim p(y|x)}\left[2(f(x_0)-y_0)| x_0\right] = 0$$

$$= 2\left(\int f(x_0)\,p(y|x)\,dy - E_{y_0\sim p(y|x)}[y_0|x_0]\right) = 0$$

mean function of $(x_0)$
should return $E_P[y]$ $\leftarrow f^*(x_0) = E_{p(y|x)}[y_0|x_0] = E_{y_0\sim p(y|x)}[y_0|x_0]$

---

① 1บาท ② $T_i^*$ (การส่งดินของ)

③ sportclub - organization/competition/ fund raizing service

Rule of Prob

sum rule $\quad p(x) = \sum_y P(x,y)$

## Iteration

$$w^j = w^j - \alpha\frac{\partial}{\partial\theta_j}\ell(w_0,w_1)$$

generative vs. discriminate approach

$$\hat{y}(x_0) : E_{y\sim p(y|x)}[y|x_0]$$

F = observe
$\hat{F}$ = estimate
$\overset{*}{F}$ = optimal

1) Generative $\longmapsto$ Approach

Estimate $p(x,y)$

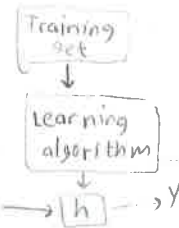$\downarrow$ Normalize $\rightarrow$ find $p(y|x)$

2) discriminate $\longmapsto$ Approach

Estimate $p(y|x)$ from data

# supervised learning
- training from the right data

# Regression Problem
- Predict real-valued output

Training set
$\downarrow$
Learning algorithm
$\downarrow$
$x \longrightarrow \boxed{h} \dashrightarrow y$

Decomposition of error

Expected $\longmapsto$ $\hat{w}$   LSQ estimate from training data
Loss     $w^*$   optimal regression parameter

$\Big\}$ w = parameter vector : $\bar{w}\cdot\bar{x}$

$$y - \hat{w}\cdot x = (y-w^*\cdot x) + (w^*\cdot x - \hat{w}\cdot x)$$

$$E[y-\overset{\varepsilon}{w^*}x]=0 \not< E[w^*x - \hat{w}x]=0$$

$$E_{P(x,y)}\left[(y-\hat{w}\cdot x)^2\right] = E_{P(x,y)}\left[(y-w^*\cdot x)^2\right] + 2E_{P(x,y)}\left[(y-w^*\cdot x)(w^*\cdot x-\hat{w}\cdot x)\right] + E_{P(x,y)}\left[(w^*x-\hat{w}x)^2\right]$$

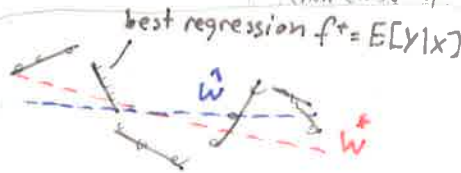prediction error $\propto$ linear fn $w^*\cdot x - \hat{w}\cdot x$

$$\boxed{E_{P(x,y)}\left[(y-\hat{w}\cdot x)^2\right] = E_{P(x,y)}\left[(y-\overset{*}{w}\circ x)^2\right] + E_{P(x,y)}\left[(w^*\cdot x - \hat{w}\cdot x)^2\right]} \quad\text{——— (6)}$$

Approximation Error (Variance)

Estimation Error (Bias)
(how close opt $w^*$ from infinit training $\hat{w}$)

error $\not< N$
$\propto$ hypothesis

best regression $f^* = E[y|x]$



$\hat{w}$   $w^*$

$\lim\limits_{N\to\infty} \hat{w} = w^*$

MSE (Mean Square Error) & bias-variance decomposition

Question for underfitting, overfitting, model capacity, MSE for estimator & predictor

$X \sim D$ distribution
$X = (X_1, X_2, \ldots, X_n)$ ; $\hat{y} = f(x;w)$ ;

- estimated variance : $\hat{\sigma}^2 = \frac{1}{n}\sum_{i}^{n}(x_i - \bar{x})^2$

- Estimator properties: bias $= E[\hat{y}] - \hat{y}_{real}$
               variance $\cdot$ Var$(\hat{y})$

Berkeley sheet

**Adding Noise** ⟶ $y = f(x;w) + v$ ⟶ useful info

$$E_{p(y|x)}[f(x;w) + v \mid x] = f(x;w) + E_{p(v)}[v]$$

**Gaussian Noise Model** ⟶ $y = f(x;w) + v$ , $v \sim N(v; 0, \sigma^2)$

distribution of — mean — Variance

$$p(y|x;w) \xrightarrow[M=0,\sigma^2]{add} p(y|x;w,\sigma) = N(y; f(x;w), \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi}\sigma} e^{-\left(\frac{(y - f(x;w))^2}{2\sigma^2}\right)} \quad \text{Gaussian} \quad (7)$$

**Likelihood** ⟶ likelihood of parameter $w$ given observed data $X = [x_1, ..., x_N]$, $Y = [y_1, ... y_N]^T$

matrix          vector

$$\mathcal{L} = P(Y|X;W,\sigma)$$

prob to observe data $y$ given $X$
under model with parameters ; $W, \sigma$

[I ID] independently, identically, distributed between set data $(x_i, y_i)$

$$\boxed{P(\vec{Y}|\underline{X}; \vec{w}, \sigma) = \prod_{i=1}^{N} p(y_i | \vec{x}_i, W, \sigma)} \quad \text{IID} \quad (8)$$

**Maximum Likelihood** ⟶ $\hat{W}_{ML} = \arg\max_W P(Y|X;W,\sigma) \cdots (7),(8)$

$$\log(\hat{W}_{ML}) = \arg\max_W \sum_{i=1}^{N} \log p(y_i | x_i, W, \sigma)$$

$$= \arg\max_W \sum_{i=1}^{N} \left[ -\frac{(y_i - f(x_i;w))^2}{2\sigma^2} - \log \sigma\sqrt{2\pi} \right]$$

$$= \arg\max_W -\frac{1}{2\sigma} \sum_{i=1}^{N} (y_i - f(x_i;w))^2 - N\log\sigma\sqrt{2\pi} \quad \text{independent to } W$$

max likeli

$$\boxed{L(f(x;w), y) = -\log p(y|x;w,\sigma) = \sum_{i=1}^{N} (y_i - f(x_i;w))^2} \quad \frac{(9)}{\text{for Gaussian Noise}}$$

**General Additive Regression Model** ⟶ $\hat{y} = f(x;w) = W_0 + W_1 \phi_1(x) + W_2\phi_2(x) + ... + W_m\phi_m(x)$

↱ some function (input: vector) $x_i$ = vector

$$\hat{W} = \begin{bmatrix} w_0 \\ \vdots \\ w_m \end{bmatrix} = ? \qquad X = \begin{bmatrix} \phi_0(x_1) & \cdots & \phi_m(x_1) \\ \vdots & & \vdots \\ \phi_0(x_N) & \cdots & \phi_m(x_N) \end{bmatrix}$$

$\hat{y} = W_0 + W_1 x + W_2 x^2 + W_m x^m$

$$X = \begin{bmatrix} 1 & x_1 & & x_1^m \\ 1 & x_2 & & x_2^m \\ \vdots & & & \\ 1 & x_N & & x_N^m \end{bmatrix}$$

$$\hat{w} = (x^T x)^{-1} x^T y$$

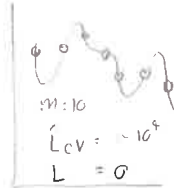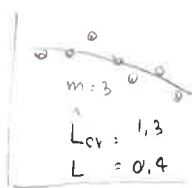Overfitted ⟶ สังเกต: too sensitive, unstable to each data point
Problem

boxed: crossvalidation

Leave-one-out cross-validation

—(10)

$$\hat{L}_{cv} = \frac{1}{N} \sum_{i=1}^{N} (y_i - f(x_i; \hat{w}_{-i}))^2$$ ; $\hat{w}_{-i}$ : fit parameter without $i$-th data



$m = 3$
$\hat{L}_{cv} = 1.3$
$L = 0.4$

$m = 10$
$\hat{L}_{cv} = \sim 10^5$
$L = 0$

---

Lecture 5        boxed: Regulation

Road Map ⟶ ▷ regulation = a tool against overfitting
              ▷ gradient descent

(Penalty) VO

Penalizing Model ⟶ **Ituition:** penalize # of bits required to encode the parameter
complexity

shrinkage
(5.1) Method

$$(4.9) \longrightarrow \boxed{W^* = \arg\max_{W} \left\{ \frac{1}{2} \sum_{i=1}^{N} \log p(data_i ; \dot{w}) - penalty(W) \right\}} \longrightarrow$$

→ given $y_i | x_i$

Loss or
$(E[l(\cdot)])$       minimum difference $y, \hat{y}$

Ridge ⟶
Regression

$$W_{ridge} = \arg\min_{W} \left\{ \sum_{i=1}^{N} (y_i - \underset{\hat{y}_i}{\underbrace{W \cdot x_i}})^2 + \lambda \sum_{j=1}^{m} w_j^2 \right\} \longrightarrow (5.2)$$

$\leftarrow f(x_i, w, -1$

$(5.4) \longrightarrow$ $$\hat{W}_{ridge}^* = (\lambda I + X^T X)^{-1} X^T y$$

$w^2$ → (5.3)

proof
below

$L = E[l(\cdot)]$

Lasso ⟶ $E(w^*) = W_{lasso} = \arg\min_{W} \left\{ \sum_{i=1}^{N} (y_i - W \cdot x_i)^2 + \lambda \sum_{j=1}^{m} |w_j| \right\}$
Regression

$|w_j|$

Problem of lasso → can't $\frac{\partial L}{\partial w}$ → **Need Numerical Opt. tools**

eq 5.5

constrain form

$$\bar{w} : \sum_{j=1}^{m} w_j^2 \leq t \qquad \bar{w} : \sum_{j=1}^{m} |w_j| \leq t$$

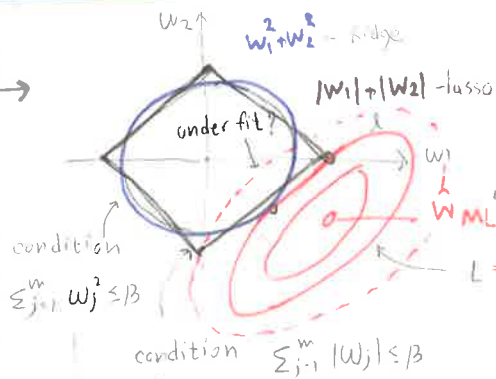---

Prove $\overset{*}{W}_{ridge}$ (from lecture 3)

$L = E[l(\cdot)] \approx \sum (y_i - w \cdot x_i)^2 + \lambda \sum_{j=1}^{m} w_j^2$

$= (y^T - x^T w^T)(y - x_w) + \lambda w^T w$

* increase the size of data
also reduce over-fitting Pb

$W^*$

ometry of $\longrightarrow$
rror Surface

$\hat{w} = \arg\max_{w: \|w\|_q \le \beta} \left\{ -\sum_{i=1}^N (y_i - w \cdot x_i)^2 \right\}$

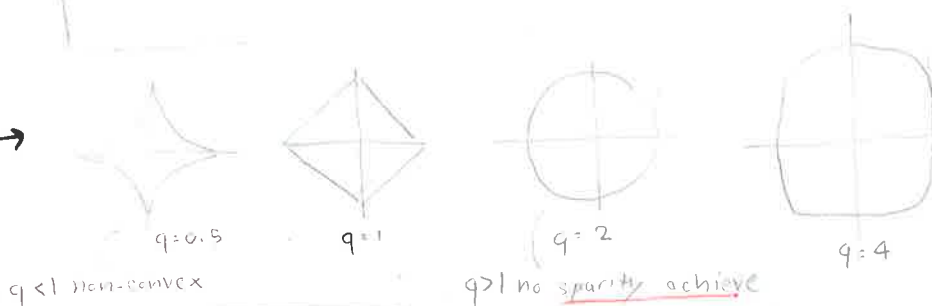$w_1^2 + w_2^2$ – Ridge

$|w_1| + |w_2|$ – lasso

under fit

over fit

optimization (5.5)

condition $\sum_{j=1}^m w_j^2 \le \beta$

$w_1$

$w_{ML}$

$L: \sum_{i=1}^N (y_i - \hat{w} \cdot x_i)^2$

condition $\sum_{j=1}^m |w_j| \le \beta$

choice of $\lambda$ $\longrightarrow$ high cross validation $\hat{L}_{cv}$

iew of $L_q$ $\longrightarrow$
penalty of
function $\|w\|_q$

$q = 0.5$     $q = 1$     $q = 2$     $q = 4$

$q < 1$ non-convex     $q > 1$ no sparity achieve

$$\|w\|_q = \left( \sum_{j=1}^m |w_j|^q \right)^{1/q} \longrightarrow (5.5)$$

$\therefore \|w\|_\infty : \max_j |w_j|$

---

## Lagrange Multiplier 2. <proof (5.c)>

Now we consider 'C' a variable (changable)

$\begin{cases} \vec{x}^* \to \vec{x}^*(c) \\ \lambda^* \to \lambda^*(c) \end{cases}$   $\mathscr{L}(\vec{x}^*(c), \lambda^*(c), c) \quad f(\vec{x}^*(c)) - \lambda^*(c)(g(\vec{x}^*(c)) - \breve{c})) \longrightarrow (5.d)$

$\mathscr{L}(\vec{x}^*(c), \lambda^*(c), c) = \hat{M}(c) - 0$

at optimization $g(\vec{x}^*(c)) = c$

$\dfrac{d\hat{M}(c)}{dc} = \dfrac{d\mathscr{L}}{dc} = \dfrac{\partial\mathscr{L}}{\partial\vec{x}^*}\dfrac{d\vec{x}^*}{dc} + \dfrac{\partial\mathscr{L}}{\partial\lambda^*}\dfrac{d\lambda^*}{dc} + \dfrac{\partial\mathscr{L}}{\partial c}\dfrac{dc}{dc} = \dfrac{\partial f}{\partial\vec{x}^*}\dfrac{d\vec{x}^*}{dc} + \cdots + \dfrac{\partial c\lambda(c)}{\partial c}\dfrac{dc}{dc} \Rightarrow \boxed{\dfrac{dL(*)}{dc} = \dfrac{d\hat{M}^*}{dc} = \lambda^*} \to (5.c)$

(5.a bottom)

(Mathematic Optimization)

(with) ## Lagrange Multiplier 1.

General form

maximize    Lagrange Multiplier constrain function $g(\vec{x}) = x \cdots$

const

$\boxed{\mathscr{L}(\vec{x}, \lambda) = f(\vec{x}) - \lambda(g(\vec{x}) - c)}$ $\to (5.a)$

$\boxed{\nabla\mathscr{L} = 0 \hookrightarrow \nabla f = \lambda\nabla g}$

Ex. labor = $20/h$ = h

steel = $2000/ton = s$

revenue = $R(h,s) = 100 h^{2/3} s^{1/3}$ $\to$ what to optimize

$y(h,s) = $ Budget = $720000 = 20h + 2000s$ $\to$ constrain

$\nabla\mathscr{L} = 0 \Rightarrow$ we get the maximized sol'n $(\vec{x}^*, \lambda^*)$ $\to (5.b)$

$Max = \hat{M} = f(\vec{x}^*) = f(\vec{x}^*(c))$ $\vec{x}^*$ is a function of 'c'

$\lambda^* = \dfrac{d\hat{M}(c)}{dc}$ change of maximization $\to (5.c)$

$\nabla R = \lambda \nabla g$ (Lagrange Multiplier)

$g(h,s):$ $\nabla R$

**Book Ch 1.2** | Probability Theory |

pdf, cdf $\longrightarrow$ $\int_{-\infty}^{\infty} p(x)dx = 1$ ; $\int_{-\infty}^{x} p(x)dx = P(x)$ ; $p(x) \geq 0$ $\longrightarrow$ (1)

pdf transform $\longrightarrow$ from $x$ to $y$ $\quad$ $p_y(y) = p_x(x)\left|\dfrac{dx}{dy}\right|$ $\qquad$ (2)

given $x = g(y)$ $\qquad = p_x(g(y))\,|g'(y)|$

Expectation $\longrightarrow$ $\mathbb{E}[f] = \overset{\text{discrete}}{\underset{x}{\sum} p(x) f(x)} = \overset{\text{continuous}}{\int p(x) f(x) dx}$

$\underline{\text{conditional}}$: $\mathbb{E}_x[f|y] = \sum_x p(x|y) f(x)$

$\underline{\text{var}}$ $\qquad$ $Var[f] = \mathbb{E}\left[(f(x) - \mathbb{E}[f(x)])^2\right] = \mathbb{E}[f(x)^2] - (\mathbb{E}[f(x)])^2$

$\underline{\text{cov}}[x,y]$ $\quad$ $\mathbb{E}_{x,y}\left[\{\vec{x} - \mathbb{E}[\vec{x}]\}\{\vec{y}^T - \mathbb{E}[\vec{y}]^T\}\right] = \mathbb{E}_{x,y}[\vec{x}\vec{y}^T] - \mathbb{E}[\vec{x}]\mathbb{E}[\vec{y}^T]$ } matrix

Bayesian $\longrightarrow$ $\vec{w} \Rightarrow$ polynomial curve fitting $\qquad$ $p(D|w)$ conditional probability
Prob $\qquad$ $p(w) \Rightarrow$ prior probability distribution

$D \Rightarrow \{t_1, \dots, t_N\}$

posterier prob. $\quad \underset{\text{post}}{p(w|D)} = \dfrac{\overset{\text{likelihood} \quad \text{prior}}{p(D|\vec{w})\, p(\vec{w})}}{p(D)}$ $\longrightarrow$ (3)

Gaussian $\longrightarrow$ $\mathcal{N}(x|\mu, \sigma^2) = \dfrac{1}{(2\pi\sigma^2)^{1/2}} e^{\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}}$ $\longrightarrow$ (4)
Distribution

$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu,\sigma^2) x^2 dx = \mu^2 - \sigma^2 \quad \rightarrow Var[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma$

For vector $\vec{x}$ [Nx1], $\Sigma$ [N×N] = covariance matrix , $|\Sigma| = $ det of $\Sigma$

$\boxed{\mathcal{N}(\vec{x}|\vec{\mu}, \Sigma) = \dfrac{1}{(2\pi)^{N/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\vec{x}-\mu)^T \Sigma^{-1}(x-\mu)\right\}}$ $\longrightarrow$ (5)

In case of $\vec{X} = \{x_1, x_2, \dots, x_N\}^T$ independent and identically distributed (i.d.d)

$\hookrightarrow \Sigma = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}_N$

$\boxed{p(\vec{x}|\mu, \sigma^2) = \prod_{n=1}^{N} \mathcal{N}(x_n|\mu, \sigma^2)}$ $\longrightarrow$ (5a)

$$\ln p(\vec{x}\,|\,\mu,\sigma^2) = -\frac{1}{2\sigma^2}\sum_{n=1}^{N}(x_n-\mu)^2 - \frac{N}{2}\ln\sigma^2 - \frac{N}{2}\ln(2\pi)$$

ML : Maximum likelihood $\boxed{\vec{x}=\mu_{ML}\ \frac{1}{N}\sum_{n=1}^{N}x_n}$ $\boxed{S^2=\sigma_{ML}^2 = \frac{1}{N}\sum^{N}(x_n-\bar{x})^2 = \frac{N-1}{N}\sigma^2 = E[\sigma_{ML}^2]}$

(6)

**non-biased**
$$\hat{\sigma}^2 = \left(\frac{N}{N-1}\right)\sigma_{ML}^2 = \frac{1}{N-1}\sum_{n=1}^{N}(x_n-\mu_{ML})^2$$

proof $\longrightarrow$ **proof**

$$Var[x_i] = E[(x_i-\mu)^2] = E[x_i^2] - E[\mu]^2 = \sigma^2 = \sigma_{model}^2 \longrightarrow (6.a)$$

$$Var[\bar{x}] = Var\left[\frac{1}{N}\sum_{n=1}^{N}x_i\right] = \langle Cor[ax,by]=(ab)\cdot cor[x,y]\rangle = \frac{1}{N^2}Var\left[\sum_{n=1}^{N}x_i\right]$$

$$= \frac{1}{N^2}\sum_{n=1}^{N}Var[x_i] = \frac{1}{N}\sigma^2 \longrightarrow (6.b)$$

$$\therefore\ E[\sigma_{ML}^2] = E\left[\frac{1}{N}\sum_{n=1}^{N}(x_n-\bar{x})^2\right] = \frac{1}{N}E\left[\sum(x_n-\mu)^2 - 2\sum(x_n-\mu)(\bar{x}-\mu) + \sum(\bar{x}-\mu)^2\right]$$

(6.a) $\searrow$

$$= \left(\frac{1}{N}\right)\left\{\sum^{N}E[(x_n-\mu)^2] - \sum^{N}E[(\bar{x}-\mu)^2]\right\}$$

(6.b) $\nearrow$

$$E[\sigma_{ML}^2] = \left(\frac{1}{N}\right)\left\{N\sigma^2 - N\left(\frac{1}{N}\sigma^2\right)\right\} = \left(\frac{N-1}{N}\right)\sigma^2 \longrightarrow (6)$$

---

urve-fitting $\longrightarrow$ = error minimalization + regulation

Input $\vec{x} = \{x_1,\dots x_N\}^T$, $\vec{Y}_{ob} = \{Y_1,\dots Y_N\}^T$, $y = y(x,\vec{w})$

$$p(t\,|\,x,\vec{w},\beta) = N(y_{obs}\,|\,y(x,\vec{w}),\beta^{-1})\quad : \beta^{-1}=\sigma^2$$

$$\mathcal{L}\ (likelihood) = \prod_{n=1}^{N}N(y_{obs,n}\,|\,y(x_n,\vec{w}),\beta^{-1}) \longleftarrow \begin{matrix}(5)\\(5a)\end{matrix}$$

$$\log(\mathcal{L}) = -\frac{\beta}{2}\sum_{n=1}^{N}\{y(x_n,\vec{w})-t_n\}^2 + \frac{N}{2}\ln\beta - \frac{N}{2}\ln(2\pi)$$

---

1.2.6
Bayesian $\longrightarrow$ (3) $\rightarrow$ post/predict distri $\quad p(Y_+|x_+,\vec{x},\vec{Y}) = \int p(Y_+|x_+,\vec{w})\,p(\vec{w}|\vec{x},\vec{Y})\,d\vec{w}$
Curve
Fitting $\qquad\qquad y_+,w_+ = New\ data,\quad \vec{x},\vec{y} = trained\ data$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \longmapsto f(\ ) \longmapsto \longmapsto pdf \longmapsto$

1. Book 1.3

## Model Selection

Lecture ⑤.1 → get the highest w/ restriction
⑤.2

cross-validation ⟶ $\left(\frac{S-1}{S}\right)$ "leave-one-out technique"

① leave 1st point, train the rest $2^{nd}$-$N^{th}$

② find prediction for 1st point

③ do it to the rest (take out $2^{nd}$ to $N^{th}$)

$\ln p(D|\vec{w}_{ML}) - M$

ook 1.5

## Decision Theory

▸ training data $(\vec{x}, \vec{y})$ (regression)

$(\vec{x}, C_k)$ (classification) ⟶ Inference →

$p(\vec{x}, \vec{y})$
$p(\vec{x}, C_k)$ } joint distribution (unknown)

⚡ EX: $\vec{x}$ - pixel of images / $C_k = \{$ normal, cancer cells $\}$
$\vec{y}'' = \{0, 1\}$

5.1 Minimize
isclassification
rate

⟶ **Bayes thm**

Book ⑤ → $p(C_k|\vec{x}) = \dfrac{p(\vec{x}|C_k) \, p(C_k)}{p(\vec{x})}$ ——————(7)

↳ $p(\vec{x}, C_k) = $ joint distribution

↳ $p(C_k) = $ prior prob for class $\{C_k\}$

↳ $p(C_k|\vec{x}) = $ proterior prob for class $\{C_k\}$

example of
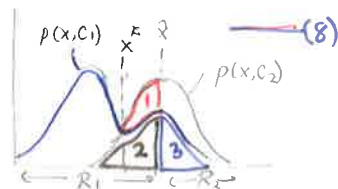$C_0 = $ normal
$C_1 = $ cancer

⟶ How to assign class to $x$?

(Input Space) $\vec{x} \in$ (region) $R_1 \xrightarrow{assign}$ class $C_1$

Ex ↓ $p(\text{mistake}) = p(\vec{x} \in R_1, C_2) + p(\vec{x} \in R_2, C_1)$

$= \int_{R_1} p(\vec{x}, C_2) d\vec{x} + \int_{R_2} p(\vec{x}, C_1) d\vec{x}$

↑ $p(\text{correct}) = \sum_{k=1}^{K} p(\vec{x} \in R_k, C_k)$

$= \sum_{k=1}^{K} \int_{R_k} p(\vec{x}, C_k) d\vec{x}$  find $C_k$ that $p(\vec{x}, C_k) = $ largest



——————(8)

$p(x,C_1)$    $p(x,C_2)$

Error
1: $C_2 \checkmark$ $C_1 \times$
2: $C_2 \times C_1 \checkmark$
Posterior: opt $x = \hat{x}$
where $p(x,C_1), p(x,C_2)$ crosses

$\boxed{\text{Information Theory}}$

- amount of Information $\rightarrow$ degree of surprise to learn $X$

- $h(x)$ $h(y)$ $\rightarrow$ Information content $\Big\}$ $h(x) = -\log_2 p(x) \xrightarrow{\text{bits}}$ (8)

- $p(x,y) = p(x)p(y)$ "Indep" $\qquad$ or $= -\ln p(x)$

ntropy $\longrightarrow$ Entropy $^{(8)}$ $\boxed{H[X] = E[h(x)] = -\sum_{x} p(x_i) \log_2 p(x)} \longrightarrow$ (9)

$\qquad\qquad\qquad$ (uniform $h(x) \rightarrow$ higher $H[X]$)

w/ regulation $\sum_i p(x_i) = 1$ $\quad$ **Lagrange multiplier**

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\nabla \mathcal{L} = 0 : \nabla R - \lambda \nabla G$

(9) $\quad \boxed{\begin{aligned} \tilde{H} &= -\sum_i p(x_i) \ln p(x_i) + \lambda\left(\sum_i p(x_i) - 1\right) \\ &= -\int p(x) \ln(p(x))dx + \lambda\left(\int p(x)dx - 1\right) \end{aligned}}$ $\longrightarrow$ (10)

Example $\longrightarrow$ Goal: find the model of Normal distribution $p(x)$

or mal distri!

Constrains : $\displaystyle\int_{-\infty}^{\infty} p(x)dx = 1$

$\qquad\qquad\qquad \displaystyle\int_{-\infty}^{\infty} x p(x)dx = x$ $\Big\}$ (11)

$\qquad\qquad\qquad \displaystyle\int_{-\infty}^{\infty} (x-\mu)^2 p(x)dx = \sigma^2$

$\begin{aligned}(11)\\(10)\end{aligned}\Rightarrow$ $\tilde{H} = -\displaystyle\int_{-\infty}^{\infty} p(x)\ln p(x)dx + \lambda_1\left(\int_{-\infty}^{\infty} p(x)dx - 1\right)$

$\qquad\qquad + \lambda_2\left(\displaystyle\int_{-\infty}^{\infty} x p(x)dx - \mu\right) + \lambda_3\left(\int_{-\infty}^{\infty} (x-\mu)^2 p(x)dx^2 - \sigma^2\right)$

$\qquad \nabla\tilde{H} = 0 = p(x)\left[-\ln p(x) + \lambda_1 + \lambda_2 x + \lambda_3 (x-\mu)^2\right]$

$\qquad p(x) = \exp\left[\lambda_1 + \lambda_2 x + \lambda_3 (x-\mu)^2\right] = \dots = \dfrac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\dfrac{(x-\mu)^2}{2\sigma^2}\right\}$
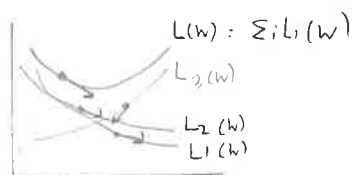
**Stirling Approximation** $\ln N! \cong N \ln N - N$

$\ln N! = \ln 1 + \ln 2 + \dots + \ln N$

$\ln N! - \frac{1}{2}(\ln 1 + \ln N) = (\ln 1 + \dots + \ln N) - \frac{1}{2}(\ln 1 + \ln N)$

# Lecture 9 | Decision Trees

Review
Stochastic
gradient descent

$L(w) = \Sigma_i L_i(w)$

$L_2(w)$

$L_2(w)$
$L_1(w)$

Goal: $\min_w L(w) = \sum_{i=1}^{N} L_i(w)$

---

## Lecture 6 | Gradient descent; bias-Varian tradeoff

$(x^T x)^{-1} x = x^+$ pseudoinverse

Problem
w/ $x^+$

▷ Least square closed form sol$^n$ $\binom{Lee}{5.9}$ → $\hat{w} = (\lambda I + x^T x)^{-1} x^T y$

▷ sometimes $x^+$ too large to compute

Alternative
Numerical Opti:
(gradient descent)

▷ Gradient ascent = 'uphill climbing'

▷ Gradient descent = 'down hill'

[Mech]  $t = 0$, $w^{(t)}$

$$g^{(t)} = \nabla f(\underline{X}, \bar{y}; w^{(t-1)})$$

learning rate  ⎫
⎬ (Lec 6.1)
update model $\boxed{w^{(t)} = w^{(t-1)} - \eta \, g^{(t)}}$  ⎭

[Example]

$\eta < \eta_{opt}$     $\eta = \eta_{opt}$     $\eta > \eta_{opt}$

optimum point

Bias of
an estimator

(Lec 6.2)

$\boxed{bias(\hat{\theta}) \triangleq \mathbb{E}_x[\hat{\theta} - \theta]}$

different $\mathbb{E}$ of predicted from correct

$\mathbb{E}[\hat{\mu}_{ML}] = \frac{1}{N} \Sigma_i^N x_i = \mu$ ✓

(book 6)

$\mathbb{E}[\hat{\sigma}^2_{ML}] = \frac{1}{N} \Sigma_i (x_i - \hat{\mu})^2 = \frac{N-1}{N} \sigma^2$

Consistency of
an estimator

Estimator $\hat{\theta}$ is consistent if $\boxed{\lim_{N \to \infty} \hat{\theta}_N = \theta}$  →  (Lec 6.3)

$\hat{\sigma}^2_{ML}$  bias, but consistent w/ $\sigma^2$

**Estimation & Regression** →

True Model $\quad y = F(x) + V \qquad ; V:$ noise with mean $= 0$

$\qquad\qquad\qquad$ function $\qquad$ function space

approximate $F'$ by $\ f(x;\hat{w}) \in \tilde{F} \qquad :$ estimate $\hat{w}$ from $\underline{x}$

$\qquad\left[\begin{array}{l}
\hat{f}(\vec{x}) = f(\vec{x};\hat{w}) \qquad \text{estimation based on this } \vec{x} \qquad\qquad \rightarrow \text{(Lec 6.4)} \\[8pt]
\bar{f}(\vec{x}) = \mathbb{E}_x[f(\vec{x};\hat{w})] \quad \text{avg estimate over training sets } \bar{x} \\[8pt]
f^*(\vec{x}) = f(\vec{x};\underset{x}{\arg\min}\ \underset{P(x,y)}{\mathbb{E}}[(y-f(x;w))^2]) \quad \text{the best estimate } f \in \tilde{F}
\end{array}\right.$

**Bias + Variance** →

▷ $\mathbb{E}_x[\text{square loss}] = \mathbb{E}_x[(y_0-\hat{f}(x_0))^2] = (y_0-\bar{f}(x_0))^2 + \mathbb{E}_x[(\hat{f}(x_0)-\bar{f}(x_0))^2]$

$\qquad\qquad\qquad\qquad\qquad$ true model $\qquad\qquad\qquad\qquad \underleftrightarrow{\quad\square\quad} \quad \leftarrow$ Variance $\longrightarrow$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \searrow\text{(Lec 6.5)}$

☐ $(y_0-\bar{f}(x_0))^2 = (y_0-F(x_0))^2 + (F(x)-\bar{f}(x_0))^2 \quad\nearrow$

$\qquad\qquad\qquad\qquad \underleftrightarrow{\text{noise}} \qquad\qquad \underleftrightarrow{\text{bias}^2} \quad$ different btw

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ avg estimation and true model

$\qquad\qquad\uparrow$ difference btw

$\qquad\qquad$ observe value and true model

$(y_0-F(x_0))^2$ **noise** = irreducible (independent to data)

$(F(x_0)-\bar{f}(x_0))^2$ **bias²** = different $f \& F$

$\mathbb{E}_x[(\hat{f}(x_0)-\bar{f}(x_0))^2]$ **Variance** = $\qquad$ $\left.\begin{array}{c}\\\\\\\\\end{array}\right\}$ try to minimize

**Fisher Information ↓**
(just information)
→

measuring information that observe $X$ carry about unknown parameter $\theta$

$\qquad\qquad\qquad\qquad\qquad\qquad$ given $\theta \to$ **no** $\int d\theta$

$I[\theta] = \mathbb{E}\left[\left(\frac{\partial}{\partial\theta}\log f(x;\theta)\right)^2\Big|\theta\right] = \int\left(\frac{\partial}{\partial\theta}\log f(x;\theta)\right)^2 f(x;\theta)dx$

$\qquad\quad = \int\left(\frac{1}{f^2}\left(\frac{\partial}{\partial\theta}f\right)^2\right)f\,dx \qquad\qquad\longrightarrow \text{(Lec 6.6)}$

$I[\theta] = -\mathbb{E}\left[\frac{\partial^2}{\partial\theta^2}\log f(x;\theta)\Big|\theta\right] = \int\left(-\frac{\partial^2}{\partial\theta^2}\{\log f(x;\theta)\}\right)f(x;\theta)dx$

$\qquad\quad = \int\left(-\frac{\partial}{\partial\theta}\left\{\frac{1}{f}\frac{\partial}{\partial\theta}f\right\}\right)f\,dx = \int\left\{\frac{1}{f^2}\left(\frac{\partial}{\partial\theta}f\right)^2 - \frac{1}{f}\frac{\partial^2}{\partial\theta^2}f\right\}f\,dx$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \underleftrightarrow{\quad①\quad} \quad \underleftrightarrow{\quad②\quad}$

$I[\theta] = \mathbb{E}\left[\frac{\partial}{\partial\theta}\log f(x;\theta)\Big|\theta\right] - \frac{\partial^2}{\partial\theta^2}\int f(x;\theta)dx$

$\qquad\qquad \underleftrightarrow{\quad①\quad} \qquad\qquad\qquad\quad 0$

$\boxed{I[\theta] = -\mathbb{E}\left[\frac{\partial^2}{\partial\theta^2}\log f(x;\theta)\Big|\theta\right]} \qquad \overset{\text{(Lec 6.6)}}{\longleftarrow}$

**Fisher Info**

$\mathbb{E}\left[\frac{\partial}{\partial\theta}\log f(x;\theta)\Big|\theta\right] = \int\left(\frac{1}{f}\frac{\partial f}{\partial\theta}\right)f\,dx = \frac{\partial}{\partial\theta}\int f\,dx = \frac{\partial}{\partial\theta}1 = 0$

**Fisher** 3. Prove $Var[\hat{\theta}] \gtrsim \frac{1}{I(\theta)}$ ———→ (lec 6.7)

**Cramér-Rao Bound**

*Proof*

unbiased estimator $\hat{\theta}(X) \Rightarrow E[\hat{\theta}(X)-\theta|\theta] = \int (\hat{\theta}(x)-\theta) f(x;\theta) dx = 0$

$$0 = \int \frac{\partial}{\partial\theta} (\hat{\theta}(x)-\theta) f(x;\theta) d\theta = \int (\hat{\theta}(x)-\theta)\frac{\partial f}{\partial\theta} dx - \underbrace{\int f dx}_{1}$$

$$f \frac{\partial \log f}{\partial \theta}$$

$$\int (\hat{\theta}-\theta) f \frac{\partial \log f}{\partial \theta} dx = 1$$

$$\boxed{1 = \left( \int \left[ (\hat{\theta}-\theta)\sqrt{f} \right] \cdot \left[ \sqrt{f} \frac{\partial \log f}{\partial \theta} \right] dx \right)^2 \leq \underbrace{\left[ \int (\hat{\theta}-\theta)^2 f dx \right]}_{Var[\theta]} \underbrace{\left[ \int \left( \frac{\partial \log f}{\partial \theta} \right)^2 f dx \right]}_{I(\theta)}}$$

similar to $\quad |cov[A,B]|^2 \leq Var[A] Var[B]$

$E[loss]$ ⌣ Meansquare = Noise + bias² + Var

bias²
Variance

---

# Lecture 7 — Logistic Regression

**classification as regression** → $y, \hat{y}$ are classes ex. $\{-1,1\}$

$f(X,\hat{w}) = w_0 + \hat{w}\cdot \underline{x}$ = function ≠ classes

decision rule $\quad \hat{y}=1$ if $f(x;\hat{w})\geq 0$, otherwise $\hat{y}=-1$ $\quad$ (L 7.1)

$$\hat{y} = sign(f(x;\hat{w})) = sign(w_0 + \hat{w}\cdot x) = h(x)$$

in here $0 =$ **decision boundary**

**Loss calculation** → (L7.1)
$$L(h(\underline{x}),y) = L(\hat{y},y) = \begin{cases} 0 & \text{if } h(\underline{x})=y \\ 1 & \text{if } h(\underline{x}) \neq y \end{cases}$$

**Risk = Expected Loss** →
$$R(h) = E_{x,y}[L(h(\underline{x}),y)] = \int_x \sum_{c=1}^{C(classes)} L(h(\underline{x}),c) \underbrace{p(y=c|x)p(x)}_{p(x,y=c)} d\underline{x} \quad \to (L7.2)$$

minimize = $R(h|x)$

$$\int_x R(h|\underline{x}) p(x) dx$$

only count when $L(\hat{y},y)=1$

$$R(h|x) = \sum_{c\neq h(\underline{x})}^{C} (1) \, p(y=c|\underline{x}) = ? \quad 1-p(y=h(\underline{x})|\underline{x})$$

timum
Rule

e logistic
Model

elationship
$\sigma(x)$ and $p(y|x)$

Decision
Boundary

$$h(x) : \arg\max_c \; p(y=c|x)$$

$$h(x) : c^* \;\longleftrightarrow\; \frac{p(y=c^*|x)}{p(y=c|x)} \geq 1 \quad \text{or} \quad \ln\left\{\frac{p(y=c^*|x)}{p(y=c|x)}\right\} \geq 0 \quad \forall c$$

$\downarrow$ (L7.3)

 $\rightarrow$ logistic

$$\frac{dN}{dt} = rN\left(1-\frac{N}{k}\right)$$

$$r = \left(\frac{1}{N} - \frac{-\frac{1}{k}}{1-\frac{N}{k}}\right)\frac{dN}{dt} = \; dt$$

$$r(t)+C = \ln(N) - \ln\left(1-\frac{N}{k}\right) = \ln\left(\frac{N}{1-\frac{N}{k}}\right) \;\Rightarrow\; \left(\frac{N}{1-\frac{N}{k}}\right) =$$

$$N(t) = \frac{1}{c_3 e^{-rt} + \frac{1}{k}} \;=\; \frac{N_0 k}{\underbrace{(k-N_0)e^{-rt}+N_0}_{\text{logistic}}} \quad\rightarrow (L7.4)$$

## simple logistic (growth) model

simple
$$\boxed{\sigma(x) = \left(\frac{1}{1+e^{-x}}\right)} \quad ; \quad \begin{cases} \sigma(-\infty)=0 \\ \sigma(0)=\frac{1}{2} \\ \sigma(\infty)=1 \end{cases}$$



complex
$$\sigma(\;) = \frac{(y_{high}-y_{low})}{1+e^{-a(x-x_0)}} + y_{low} \quad \rightarrow (L7.5)$$



$y$ = vector    $y_i$ = value
$X$ = Matrix    $x_i$ = vector

(L23)

Binary     **Boundary**

$\triangleright$ $\ln\left\{\frac{p(y=1|x)}{p(y=0|x)}\right\} \geq 0 = w_0 + w\cdot x$

$\triangleright$ $\dfrac{p(y=1|x)}{1-p(y=1|x)} = e^{(w_0+w\cdot x)} = 1$

(L7.5)

$\Rightarrow p(y=1|x) = \left(\dfrac{1}{1+e^{-(w_0+w\cdot x)}}\right) = \sigma(w_0+\bar{w}\cdot\bar{x}) = \frac{1}{2} = \sigma(0)$

$p(y=0|x) = \left\{\dfrac{e^{-(w_0+\bar{w}\cdot\bar{x})}}{1+e^{-(w_0+\bar{w}\cdot\bar{x})}}\right\} = 1-\sigma(w_0+\bar{w}\cdot\bar{x}) = \frac{1}{2}$

$\Big\}$ (L7.6)

$y_i=1 \qquad y_i=0$
$$p(y_i|\bar{x}_i;w) = \sigma(w_0+\bar{w}\cdot x_i)^{y_i}\,(1-\sigma(w_0+\bar{w}\cdot x_i))^{1-y_i}$$

$$\log p(\bar{y}|x;\bar{w}) = \sum_{i=1}^{N} \log p(y_i|\bar{x}_i;w)$$

**Lecture 8** | Regulation in logistic regression; stochastic gradient descent; Softmax

optimal regressor  $\hat{y} = E[y|x]$

optimal classifier  $\hat{y} = \underset{c}{\text{argmax}} \; p(y=c|\underline{x})$

**Review:**
**Logistic**
**Regression**

> - log-odds as a : $\log \dfrac{p(y=1|x)}{p(y=0|x)} = f(\emptyset(\underline{x}); \bar{w}) = 0$
>   function of X $\quad (1 - p(y=1|x))$

- $p(y=1|\underline{x}) = \dfrac{1}{1+e^{(-f(\emptyset(x);\underline{w}))}} = \dfrac{1}{1+e^{(-w_0 - \underline{w}\cdot\underline{x})}}$

  can be non linear  $\emptyset[x] = \left(1, x_1, x_2, x_1x_2\right)$

**Gradient**
**Descent**

> - $1^{st}$ order iterative optimization algorithm
>   $(L2.9) \longrightarrow$  $\log p(\bar{Y}|\underline{x}, \bar{w}) = \sum\limits_{i=1}^{N}\left[ y_i \log\sigma(w_0 + \bar{w}\cdot\bar{x_i}) \; (1-y_i)\log(1-\sigma(w_0 + \bar{w}\cdot\bar{x_i}) \right]$
>   $($

$(L7.8)$

$\dfrac{\partial}{\partial w_0}\log p(\bar{Y}|\underline{x}; \bar{w}) = \sum\limits_{i=1}^{N}\left[ \dfrac{y_i \; \sigma(1-\sigma)}{\sigma}\left(\dfrac{\partial w_0}{\partial w_0}\right) + \dfrac{(1-y_i)(-\sigma)(1-\sigma)}{(1-\sigma)}\dfrac{\partial w_0}{\partial w_0} \right]$

$\qquad = \sum\limits_{i=1}^{N}\left[ y_i - y_i\sigma - \sigma + y\sigma \right] = \sum\limits_{i=1}^{N}\left[ y_i - \sigma(w_0 + \bar{w}\cdot\bar{x_i}) \right] = 0 \quad (L7.9A)$

$\dfrac{\partial}{\partial w_j}\log p(\bar{Y}|\bar{x}; \bar{w}) = \sum\limits_{i=1}^{N}\left[ \dfrac{y_i \; \sigma(1-\sigma)}{\sigma} \dfrac{\partial \bar{w}\cdot\bar{x_i}}{\partial w_j} + \dfrac{(1-y_i)(-\sigma)(1-\sigma)}{(1-\sigma)}\left(\dfrac{\partial \bar{w}\cdot\bar{x_i}}{\partial w_j}\right) \right]$

$\qquad = \sum\limits_{i=1}^{N}\left[ \left(y_i - \sigma(w_0 + \bar{w}\cdot\bar{x_i})\right) x_{ij} \right] \quad (L7.9B)$

**Updated**
**W**

$w_{new}^{(t+1)} = \bar{w}^{(t)} + \eta\dfrac{\partial}{\partial \bar{w}}\log p(\underline{X}; \bar{w})$

$\qquad = \bar{w}^{(t)} + \eta\sum\limits_{i=1}^{N}(y_i - \sigma(w_0 + \bar{w}\cdot\bar{x_i}))\begin{bmatrix}1\\x_i\end{bmatrix}$
$\qquad\qquad (L7.9A) \quad \longleftarrow (L7.9B) \qquad (L7.9)$

$$\frac{1}{N} \sum_{i=1}^{N} \frac{\partial}{\partial w} L(y_i, \bar{x}_i; w) \approx \frac{\partial}{\partial w} L(y_t, \bar{x}_t; \bar{w})$$

$$\nabla L(w) \approx N \nabla L_i(w)$$

$$w \approx w + \eta \frac{\partial}{\partial w} \log p(y_i | \bar{x}_i; w) \quad \longleftarrow \quad \text{L7.9} \qquad \text{'; } \log p(\bar{y} | \underline{x}; \bar{w}) \text{ for full update}$$

$$w = w + \eta (y_i - \sigma(\bar{w}'\bar{x}_i)) \bar{x}_i$$

$$w = w + \eta (y - \sigma(\bar{w}^T \underline{x})) \underline{x} \qquad \text{for full update}$$

Example 2D, $w_0 = 0$

$$\therefore \hat{p}(y=1|x) = \sigma(w_0 + \bar{w} \cdot \underline{x}) = \sigma(w_1 x_1 + w_2 x_2)$$

Mapping b.c. to parameter



$$\log p(Y|X, \bar{W}; \sigma) = \log p(\bar{Y}|\underline{x}, \bar{w}) + \log p(\bar{w}; \sigma)$$

(lec 5)

similar to penalty $\quad \frac{1}{2\sigma^2} \|w\|_2^2$

$$= \sum_{i=1}^{N} \log p(y_i | \bar{x}_i, \bar{w}) - \left(\frac{1}{2\sigma^2} \sum_{j=1}^{d} w_j^2\right) + \text{const}(w)$$

ML (min Less) $\qquad \sigma^2 = 1$

$$\sigma(\bar{z})_i = \frac{e^{z_i}}{\sum_{j=1}^{k} e^{z_j}} \qquad \text{for } i = 1,\dots,k \quad \text{and } \bar{z} = (z_1,\dots,z_k) \in \mathbb{R}^k$$
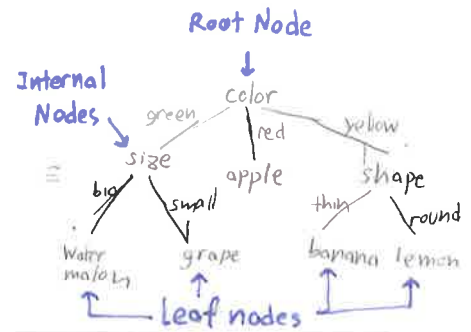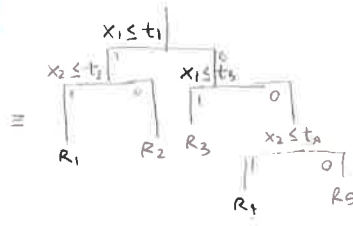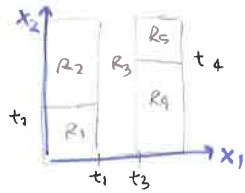
$$p(y=c|x) = \frac{e^{\bar{w}_c \cdot \phi(x) - a}}{\sum_{k=1}^{c} e^{\bar{w}_k \phi(x) - a}} \qquad ; \ a = \max_c \ \bar{w}_c \cdot \phi(x) \qquad \therefore \ \frac{e^{\bar{w}_c \cdot \phi(x) - a}}{\square} = \frac{e^{\leq 0}}{\square} = \frac{\leq 1}{\square}$$

**Lecture 9** | Decision Tree

**Space Partition** → Ex 2D



**Regression Tree** →

Algorithm → CART (Classification and Regression Trees)

Usages

1) To calculate **Prob** that a given data belong to each class

2) To classify the new data to the most likely **class**