

Intro - Ds:  
 $\langle P, \mathcal{F}, \Omega \rangle$

$$P: \mathcal{F} \rightarrow [0, 1]$$

$$\textcircled{1} \quad 0 \leq P \leq 1$$

$$\textcircled{2} \quad P(\Omega) = 1$$

$$\textcircled{3} \quad P(\bigcup A_i) = \sum P(A_i)$$

\* LTRF  $N \rightarrow P$

Lemma:  $P(A), P(A \cup B),$

RV:  $\langle \Omega, \mathcal{F}, P \rangle \quad X: \Omega \rightarrow \mathbb{R}$

(c) DF:  $F_X = P(X \leq x), \forall x \in \mathbb{R}$

$X,$

$\text{card}(X) = ?$

CDF:  $\exists f: \mathbb{R} \rightarrow [0, 1] \text{ PDF}$

$$f_X = P(X=x) = \int_{\Omega} \delta_{x(\omega)}$$

$$\text{Th: } P(X=x) = 0$$

$$f(x) = F'(x)$$

$$\int_0^1 f(x) dx = 1$$

① transformation

② Expectation:  $E(X) = \int x dF_X / (\lambda \in L^1(P))$

Variance:  $X \in L^2(P) \quad \text{lemma: } E(E(X|Y)) = E(X)$

③  $\mathbb{R}^n: JDF: F_{(X)}: \mathbb{R}^n \rightarrow [0, 1]$

$$F_{(X)}(x_1, x_2) = \frac{P(X_1 \leq x_1, X_2 \leq x_2)}{P(\Omega)}$$

$$F_{(X)}(x_1, x_2) = \frac{P(X_1 \leq x_1, X_2 \leq x_2)}{P(\Omega)}$$

Convergence:

- $\hat{X}_n \rightarrow X: \lim_{n \rightarrow \infty} F_{(X_n)}(x) = F(x)$
- $\hat{X}_n \xrightarrow{P} X: \lim_{n \rightarrow \infty} P\{\hat{X}_n > x\} = 0$
- $\hat{X}_n \xrightarrow{D} X: P\{\lim_{n \rightarrow \infty} X_n = X\} = 1$
- $E[(X_n - X)^2] \rightarrow 0$

# Probability Theory

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Bayes

①  
②  
③ ✓

independence  
Total probability

Th: Markov:  $X \in L^1(P), X \geq 0$

$$P(X \geq \epsilon) \leq \frac{E_X}{\epsilon}, \forall \epsilon > 0$$

Chebychev:  $P(|X - E(X)| \geq \epsilon) \leq \frac{V_X}{\epsilon^2}$

Lemma:  $P(X \in (a, b)) = 1 \quad \& \quad E(X_I) = 0$

$$\forall \lambda \in \mathbb{R}, \quad E(e^{\lambda X}) \leq e^{\frac{\lambda^2 V_X}{2}}$$

\* TH:  $X_1, X_2 \in L^1(P), \quad P(X_1 \in (a, b)) = 1$

$$P(|X_1 - E(X_1)| \geq \epsilon) \leq 2e^{-\frac{\epsilon^2 V_{X_1}}{2}}$$

→ AP:  $CI = (\bar{X}_n - \epsilon, \bar{X}_n + \epsilon), \alpha$

? sub-Gaussian (a)  $E(e^{s(X - E(X))}) \leq e^{\frac{s^2 V_X}{2}}$

sub-exponential (b)  $E(e^{s(X - E(X))}) \leq e^{\frac{s^2 V_X}{2}}, |s| \leq \frac{1}{\sqrt{V_X}}$

Th:  $s \cdot G(\sigma) \Rightarrow P(X_n - E(X_n) \geq \epsilon) \leq e^{-\frac{\epsilon s}{2\sigma^2}}$

$$s \cdot e(\alpha) \Rightarrow \forall \epsilon \leq C \cdot \frac{\epsilon}{\alpha^2} \vee e^{-\frac{C \epsilon^2}{\alpha}}$$

Lemma:  $\begin{cases} X \sim s \cdot G(\sigma) \Rightarrow \alpha X \sim (s/\alpha)\sigma \\ X \sim s \cdot e(\alpha) \Rightarrow \alpha X \sim (s/\alpha) \\ X \sim s \cdot G(\sigma) \Rightarrow X \sim s \cdot G(\sigma) \end{cases}$

$$P(X \in (a, b)) = 1 \Rightarrow X \sim s \cdot G(\frac{b-a}{2})$$

$$X \sim s \cdot G(\sigma) \Rightarrow X^2 \sim s^2 \sigma^2$$

$$X, Y \sim s \cdot G(\sigma) \Rightarrow X+Y \sim s \cdot G(\sqrt{2s^2 \sigma^2})$$

\* form

LLN:  $X_i \in L^2(P) \text{ iid. } E(X_i) = \mu$

$$\begin{array}{l} \hat{X}_n \xrightarrow{P} \mu \\ \hat{X}_n \xrightarrow{D} \mu \end{array}$$

CLT:  $X_i \in L^2(P) \text{ iid. } E(X_i) = \mu, V(X_i) = \sigma^2$

$$Z_n = \frac{\hat{X}_n - E(\hat{X}_n)}{\sqrt{V(\hat{X}_n)}} \xrightarrow{D} Z \sim N(0, 1)$$

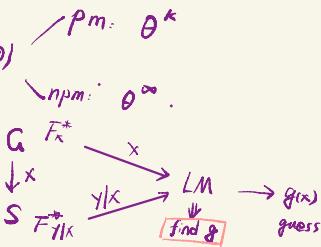
\* Bennett's Inequality

# Estimation

Fun. (data)

Stat model:  $S = \{f_{\theta(x)}, \theta \in \Theta\}$

SLP:  $\langle G, S, LM \rangle$ :



how produce an approx?

a measure of quality in math  
loss function  $L(z_{true}, z_{pred})$

$$R(\beta) = E[L(z_{true}, z_{pred})]$$

$$\text{of regression: } (x_i, y_i) \sim F_x^*, L(a, b) = (a - b)^2$$

$$\Rightarrow R(\beta) = E[(Y - f(x))^2]$$

$$= E[(Y - r + r - f(x))^2] = I + II + III$$

$$= E[(Y - r)^2] + E[(r - f(x))^2]$$

$$\text{argmin}_{\beta} R(\beta) = \frac{\partial}{\partial \beta} R(\beta) = 0 \Leftrightarrow \text{MSE}_F = \frac{\text{bias}}{\text{var}}$$

Q2: Classification  $\Rightarrow Y_i$  is discrete  $\Rightarrow 0-1$  loss  $\rightarrow R$   $\rightarrow R(\hat{f}, f)$

$$Q3: S = \{P_{\theta}(x), \theta \in \Theta\}$$

JDF

$$\text{Empirical: } \hat{P}_{(n)} = \frac{1}{n} \sum (-\ln P_{\theta}(x_i))$$

$$Q_f \Leftrightarrow Q_{\hat{f}} \text{ if } f \text{ is BM}$$

$$f_{\hat{y}x} = f_{y|x} - f_x \Rightarrow \text{Info}_{\hat{f}} = \text{Info}_{y|x} + \text{Info}_x$$

$$f_{\hat{y}x} = P_{\theta}^*(\hat{f}_x)$$

$\Rightarrow$  linear Reg.

$$\text{logistic Reg.: } P_{\theta(x)}(y|x) = \text{Ber}(p) = \text{Ber}(G(x; \alpha, \beta)), G = \frac{1}{1 + e^{-x}}$$

$$P_{\theta}(y|x) = p^y (1-p)^{1-y} \quad \hat{P} \rightarrow \text{cross entropy}$$

Point:

parametric:  
Point est: guess (point)



data:  $X_{\text{obs}}$ :  $S \rightarrow X$   $\text{CR}^{\text{obs}}$   
Matrix

Statistic:  $T(x) : X \rightarrow T$   
point esti:  $\hat{\theta} : X \rightarrow \Theta$   
Up sampling distri

unbiased (asym)

$\hat{\theta} \xrightarrow{D} \theta$  (asympt.)

$\Rightarrow$  measure quality:

$$\text{Bias}(\hat{\theta}(z)) = E(\hat{\theta}) - \theta$$

$$\text{std}(\hat{\theta}(z)) = \sqrt{V(\hat{\theta}(z))}$$

$$\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$$

$$\text{prop: } \text{MSE} = \text{bias}^2 + \text{Var}$$

Non-parametric

$$\text{empirical DF: } \hat{F}_n(x, \bar{x}) = \frac{1}{n} \sum \mathbb{1}_{x_i \leq x}$$

$$\text{lemma: } E(\hat{F}_n) = F_{\bar{x}}$$

$$V(\hat{F}_n) = \frac{F(1-F)}{n}$$

DKW inequality:  $\mathbb{P}(\hat{F}_n \geq F + \epsilon) \geq 1 - \delta$

$$\mathbb{P}(\sup_x |\hat{F}_n(x) - F(x)| > \epsilon) \leq 2e^{-2n\epsilon^2}$$

$$\Rightarrow \text{CI: } (\hat{F} - \epsilon, \hat{F} + \epsilon)$$

# RV - Generation

UPRNG:  $\frac{N(x)}{n} \rightarrow f_1 dx$

pseudo-random on  $M$ :  $\frac{N(x)}{n} \rightarrow \frac{1}{M}$

→ how to make a UPRNG

① Congruential generator:  $(a, b, M)$ ,  $D(x) = (ax + b) \bmod M$

lemma:  $a \in D$  or  $\exists M = 50 \dots M-1$  with period  $M$

? pseudo-random

? how to have  $M$  as period

Th: Full-Dobell  $D$ :  $(a, b, M)$  has period  $M$ .

2nd:  $\gcd(a, M) = 1$

$a^{-1} \bmod p = 0$  for every prime  $p$  that divides  $M$

? how to generate any distribution

Alg: Accept-Reject sampler.

In:  $f(x)$ : target

$g(x)$ :  $f \leq M g(x)$

Out:  $(x_1, \dots, x_n) \sim f$

do

iter:  $t$

$$x = g(t) \Rightarrow r_{t+1} = \frac{f(x)}{M g(x)}$$

while: length  $= l$   $x \leftarrow U_{[0, 1]} \cdot r_{t+1}$

length  $= l$   $x \leftarrow U_{[0, 1]} \cdot r_{t+1}$

# Markov Chains (Stochastic process)

Stochastic process:  $(X_t)_{t \in \mathbb{N}}$  ( $A = N$ )

⇒ Markov chain:  $P(X_{t+1} = x | F_t) = P(X_{t+1} = x | X_t)$

→ homogeneous:  $P(X_{t+1} | F_t) = P(X_{t+1} | X_t)$

Irreducible:  $S_i \leftrightarrow S_j$  ( $S_i \rightarrow S_j \Rightarrow P(X_t = j | X_0 = i) > 0$ )

Stationary distri.:  $\pi P = \pi$ ,  $\sum_{i \in S} \pi_i = 1$

Reversible:  $\pi_{ik} P_{kj} = \pi_{kj} P_{ik}$

i-th state → other state = stay in i-th

Question: ① Now in i-state, n step → ?-state:  
 $* P^n$

② the expectation step from a first to b

$$\begin{cases} E_a = P_{ab} E_a + P_{ba} E_b + P_{cb} E_c \\ \vdots \\ E_c = \dots \end{cases}$$

$E_a$ : a → b step  
 $E_b$ :

③ period:  $i \rightarrow i$ :  $\gcd(g_i(d))$

④ steady state vector:  $s = \lim_{n \rightarrow \infty} p^n$

$p^{(n)} = p^{(n-1)} p^n$  is  $p^n$ 's eigen vector

if  $\lim_{n \rightarrow \infty} p^n = s$ ,  $p$  is normal  $\Rightarrow s_i = s_j$   
 $p_{ij} > 0$

# Pattern Recognition

\* when  $\hat{Y} \in S_0 \dots k_1$ : class  
 $U_k = \{g_k(x) : g_k(x) \in S_0 \dots k_1\}$   
 $0-1$  loss  $L_{0-1}(y_i, u_j) = \int_{u_j}^1 I(u_i \neq u_j)$   
 $R(\alpha) = E[L_{0-1}(g_\alpha(x), y)]$   
 $\min R \Leftrightarrow \hat{\alpha}^* = \arg \min R(\alpha)$

# linear classification:  $\text{LRD}, \text{PERF}_1, \dots$   
 $g_\alpha(x) = \text{sign}(w^\top x + b)$  (WEIRD, CCR)  
 $L = -\text{sign}(w^\top x + b) \cdot y + 1 \rightarrow \text{cross-entropy}$   
 ref:  $\tilde{w} = (w, c), \tilde{x} = (x, 1)$   
 $\langle w, x \rangle + c = \langle \tilde{w}, \tilde{x} \rangle$   
 $\Rightarrow \langle \tilde{w}, \tilde{x} \rangle = 0 \Rightarrow \text{super-plane using co.}$

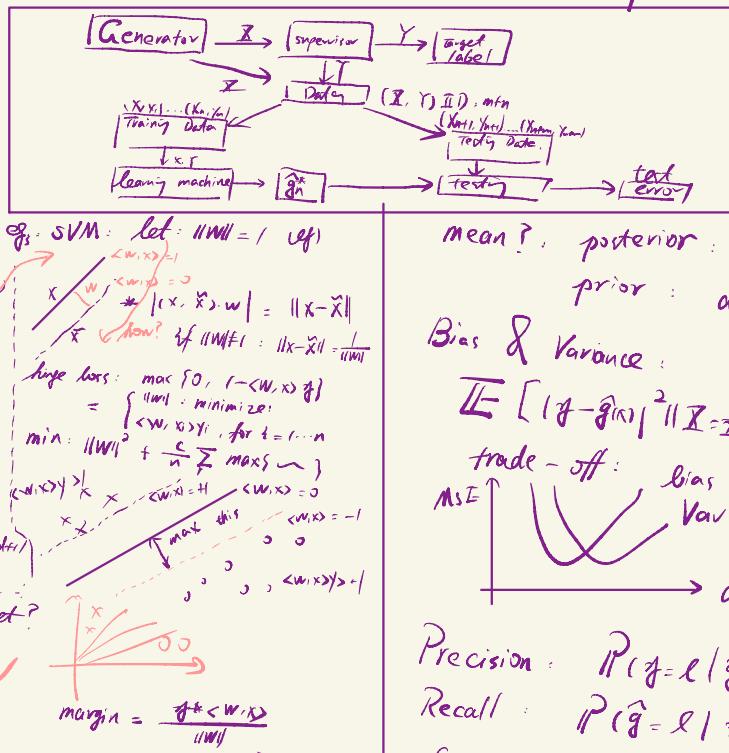
Qf2: Perception AI.  
 $0 \cdot W = (0 \dots 0)_d \cdot 1$   
 $\exists i \in \{1, \dots, d\} : 0_i < 0$   
 $W \leftarrow W + x_i \cdot y_i$

Th. most  $\ell^2$  norm  $\ell^2$  updates,  $R = \max_i |x_i|$   
 if  $\exists i \in \{1, \dots, d\} : y_i \geq 1$ , to find a  $w$ :  
 $w \cdot x_i y_i \geq 0 \quad \forall i$

Why:  $\nabla_w (\langle w, x_i \rangle y_i) \quad \text{dim}(x_i) = d+1$   
 $= x_i y_i$   
 But for non-linear divisible set?  
 → The kernel trick.

If:  $w \leftarrow w + x_i y_i$  ( $w = 0$  initial)  
 $\Rightarrow w = \sum_i x_i y_i$  ( $y_i \in \{-1, 1\}$ )  
 $\langle w, x_k \rangle = \sum_i c_i \langle x_i, x_k \rangle$   
 We have:  $\phi \rightarrow z_i = \phi(x_i)$   
 $\langle \tilde{w}, z_k \rangle = \sum_i c_i \langle \phi(x_i), \phi(x_k) \rangle$   
 $\Rightarrow \text{function } K(a, b), a, b \in \mathbb{R}^d$   
 $K(a, b) = K(x_i, x_k) \quad \phi: \mathbb{R}^d \rightarrow \mathbb{R}^d$   
 $\Leftrightarrow K(a, b) = \langle \phi(a), \phi(b) \rangle$   
 $K(a, b) = \langle r(a, b) + m \rangle \dots$

⇒ lemma:  $K$



## Training and Testing

LM: minimizes

$$g^* = \arg \min \frac{1}{n} \sum L(g, y_i)$$

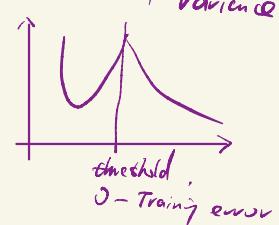
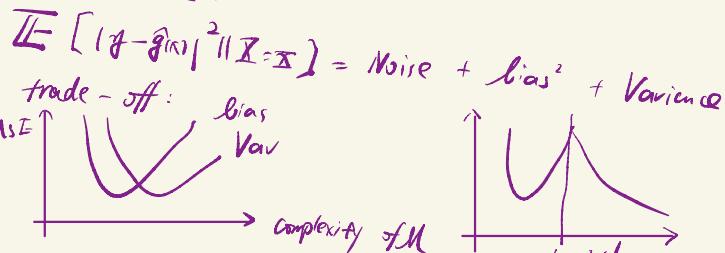
test error:  $\frac{1}{m} \sum L(\hat{g}_n, y_i)$

test metric: make decision correctly

prior: once trained we fix

mean? posterior: fed error fixed

Bias & Variance:



Precision:  $P(f = l | \hat{g} = l)$

Recall:  $P(\hat{g} = l | f = l)$

Covariates:  $L: 0-1$  loss

$$L(\hat{g}_n, y_n) \sim \text{Ber}(p)$$

Theoretical guarantee:

$$P(R(\hat{g}) - R_g / \epsilon) < 2e^{-\epsilon^2}$$

If loss is bounded even if they're off

Linear classification:

$$P(R(\hat{g}) - \inf_{\text{cls}} R(g) + \epsilon) < e^{-\epsilon^2 / (Cn)} e^{-\epsilon^2} \quad \text{how to do it}$$

$$R(\hat{g}) = \frac{1}{n} \sum R_i$$

$$\text{VC theory: } M(A) = \frac{1}{n} \sum \mathbb{1}_{x_i \in A}$$

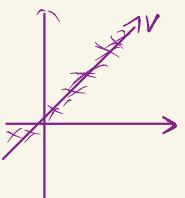
$$D(A) = P(\alpha_i \in A)$$

$$= P(L_i(Y, f(x)) = 1)$$

$$P(R(\hat{g}) - \inf_{\text{cls}} R(g) + \epsilon) < 8 \left( \frac{n}{K} \right)^{\frac{2}{\epsilon^2}}$$

$K$ : VC-dimension, complexity

## Dimensionality Reduction



$$y_i = x_i \cdot V, \|V\| = 1$$

$$\{y_i\} \text{ assume } \bar{y}_n = 0 \quad \text{or} \quad \tilde{y}_i = y_i - \bar{y}_n$$

$$\frac{1}{n} \sum (y_i - \bar{y}_n)^2 = \frac{1}{n} \sum (x_i \cdot V)^2$$

Goal: select  $V$  such that  $\max_{V \in \mathbb{R}^m} \sum (x_i \cdot V)^2$

$V_1 = \arg\max_{\|V\|=1} \frac{1}{n} \sum (x_i \cdot V)^2$  \* not unique?  
→ first singular vector

else: Greedy algorithm: first singular value  
 $\sigma_1 = \sqrt{\sum (x_i \cdot V_1)^2}$

$$V_2 = \arg\max_{\substack{\|V\|=1 \\ V \perp V_1}} \frac{1}{n} \sum (x_i \cdot V)^2$$

$$V \perp V_1, V_2, \dots, V_r \quad (r \leq \min(m, n))$$

\* Linear Algebra:  $A = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}_{n \times m}$   $AV = \begin{pmatrix} x_1 \cdot V \\ x_2 \cdot V \\ \vdots \\ x_n \cdot V \end{pmatrix}$   $\|AV\|^2 = \sum (x_i \cdot V)^2, V_i = \arg\max_{\|V\|=1} \|AV\|^2$   
Once:  $AV = 0$  for all  $V \perp V_1, \dots, V_r$ . choose a base for the null space of  $A$

Th:

## SVD

$$\text{proj}_{\text{im } A} A$$

$V_1, \dots, V_m$  is a base in  $\mathbb{R}^m$

$$x_i = \sum_{j=1}^m (x_i \cdot V_j) V_j$$

$$\Rightarrow A = \sum_{j=1}^m A V_j V_j^T, \quad U_j = \frac{A V_j}{\sigma_j} \text{ if } \sigma_j > 0.$$

$$\Rightarrow A = \sum_{j=1}^m A V_j V_j^T = \sum_{j=1}^m \sigma_j U_j V_j^T \quad U_j \text{ is the left singular vector}$$

$$U = (u_1 | u_2 | \dots | u_m)_{n \times m}, D = (\sigma_1, \sigma_2, \dots, \sigma_m)_{m \times m}, V = (V_1 | V_2 | \dots | V_m)_{m \times m}$$

## Low rank approx.

$$Ak = \sum_{j=1}^k \sigma_j U_j V_j^T, \quad \sigma_j > 0$$

Rank( $A_k$ ) =  $k$

→ error:  $\tilde{x}_i = (A - Ak)x_i$

$$\|A - Ak\|_{Fro}^2 = \sum_{i=1}^n (a_{ij} - a_{ik})^2$$

$$= \sum_{i=1}^n \|x_i - \tilde{x}_i\|_2^2 = \sum_{i=1}^n \sum_{j=k+1}^m \|\tilde{x}_i\|_2^2 = \sum_{j=k+1}^m \sigma_j^2$$

$(V_1 | \dots | V_m)$  all orthogonal eigenvectors  
 $\sigma_1, \dots, \sigma_m$  eigenvalues of  $A^T A$

$$\|AV\|^2 = \langle AV, AV \rangle = \langle A^T AV, V \rangle$$

$$\tilde{A} = A^T A \quad \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m \geq 0$$

$$\Rightarrow V = \sum_{i=1}^m \sigma_i V_i = \sum_{i=1}^m \sigma_i^2 V_i$$

$$\|V\| = \sqrt{\sigma_1^2 + \dots + \sigma_m^2}$$

⇒ Lemma:  $A_{n \times m}, V_i$  first singular vec. of  $A$ .  
 $\sigma_i$ : first singular value;  $\Rightarrow V_i$  is eigenvector of  $A^T A$  &  $\max_{i \in \{1, \dots, m\}} \|AV_i\| = \|AV_i\| = \sqrt{\sigma_i^2} = \sigma_i$

④ power method:

$$\tilde{A}^T \tilde{A} = V D^2 V^T \Rightarrow A^T A V_i = \sigma_i^2 V_i$$

$$\text{Ex: } U V^T = I_m$$

⑤ PCA:

$$AV = UD$$