

1. PROBABILITY: Lecture 1

• DEFINITIONS:

• Experiment is an activity that produces distinct well defined possibilities called outcomes

• $\Omega = \{\text{all outcomes}\}$:= sample space

• Trial: doing an experiment once and getting an outcome

• Events: subsets of Ω

• $w \in E \Rightarrow E$ occurred

• "THE LONG TERM RELATIVE FREQUENCY": $N(A, n) = \frac{1}{n} \cdot (\# \text{ of times } A \text{ happens})$

$$\boxed{1} \quad N(\Omega, n) = 1 \quad \wedge \quad N(A, n) \in [0, 1]$$

$$\boxed{2} \quad A \cap B = \emptyset \rightarrow N(A \cup B, n) = N(A, n) + N(B, n)$$

$$\boxed{2.1} \quad A_1, \dots, A_j, \dots, A_i \cap A_k = \emptyset \implies N(\bigcup_{k=1}^n A_k, n) = \sum_{k=1}^n N(A_k, n)$$

$\boxed{3}$ Each experiment is independent of each other

• DEF: \mathcal{F} a collection of subsets of Ω is σ -algebraic if:

$$(1) \quad \Omega \in \mathcal{F}$$

$$(2) \quad A \in \mathcal{F} \implies A^c \in \mathcal{F}$$

$$(3) \quad A_1, A_2, \dots \in \mathcal{F} \implies \bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$$

• DEF: \mathcal{F} σ -algebra. $P: \mathcal{F} \mapsto [0, 1]$ is a probability measure if:

$$(1) \quad P(\Omega) = 1$$

$$(2) \quad A \cap B = \emptyset \implies P(A \cup B) = P(A) + P(B)$$

• DEF: $P(B|A) = \frac{P(A \cap B)}{P(A)}$:= conditional prob. (B given A)

• EX: $\Omega = \{\text{"free,spam"}, \text{"no free,spam"}, \text{"free,notspam"}, \text{"no free,notspam"}\}$
Probability that the text is spam given that it contains "Free"?

$$B = \{\text{"free,spam"}, \text{"free,notspam"}\}$$

$$A = \{\text{is spam}\} = \{\text{"free,spam"}, \text{"no free,spam"}\}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(\{\text{is free,spam}\})}{P(B)}$$

• FORMULAS:

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

$$\rightarrow P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$\rightarrow P(A \cup B) \leq P(A) + P(B)$$

$$\rightarrow A \wedge B \text{ INDEPENDENT} \Rightarrow P(A \cap B) = P(A)P(B)$$

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

2. RANDOM VARIABLES

• DEF: (Ω, \mathcal{F}, P) prob. triple. A random variable $X: \Omega \rightarrow \mathbb{R}$ is a function such that $\forall x \in \mathbb{R}, X^{-1}((-\infty, x]) := \{w: X(w) \leq x\} \in \mathcal{F}$.

We assign probability to the RV X as follows:

$$P(X \leq x) = P(X^{-1}((-\infty, x])) = P(\{w: X(w) \leq x\})$$

• EX: $\Omega = \{H, T\}, X(w) := \begin{cases} 1, & \text{if } w = H \\ 0, & \text{if } w = T \end{cases}$

$$X^{-1}((-\infty, 0]) = \{T\}$$

$$X^{-1}((-\infty, 1]) = \{H, T\} = X^{-1}((-\infty, 2])$$

• DEF: X is a RV DISCRETE, if it takes discrete values: $0, 1, 2, 3, \dots$

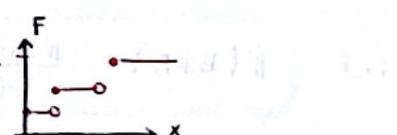
• DEF: Let X be a \mathbb{R} -valued discrete RV. $f: \mathbb{R} \rightarrow [0, 1]$ where

$$f(x) := P(X=x) = P(\{w: X(w)=x\}) = \begin{cases} \theta_i, & x = X_i \in X \\ 0, & \text{otherwise} \end{cases}$$

• PROBABILITY MASS FUNCTION "PMF"

• DEF: $F(x) = P(X \leq x) = P(\{w: X(w) \leq x\})$ = "DISTRIBUTION FUNCTION" (or cumulative DF)

- non decreasing (order): $x_1 < x_2 \rightarrow F(x_1) \leq F(x_2)$
- continuous (right)
- X has DF $F(x) \Rightarrow X \sim F$



• FORMULAS: Definitions and properties

$$\rightarrow F_X(x) = P(X \leq x) = P(\{X = X_1\} \cup \{X = X_2\} \cup \dots \cup \{X = X_k\} \mid X_{k+1} > x)$$

$$= \sum_{X_i \leq x} P(X = X_i) = \sum_{X_i \leq x} f_X(X_i)$$

$$\rightarrow \lim_{x \rightarrow -\infty} F(x) = 0$$

$$\rightarrow F_X(b) - F_X(a) = \sum_{a < X_i \leq b} f_X(X_i)$$

$$\rightarrow \lim_{x \rightarrow +\infty} F(x) = 1$$

$$\rightarrow \sum_{X_i} f_X(X_i) = 1$$

• DEF: $E[X] = \sum_x x f(x)$:= expectation / mean

• COMMON EXPECTATIONS:

- $V[X] = E[(X - E(X))^2] = \sum_{x_i} (x_i - E(X))^2 f_X(x_i)$

- $E[(X - E(X))^k]$:= k -th central moment

- $\sigma(X) = \sqrt{V(X)}$:= Standard Deviation

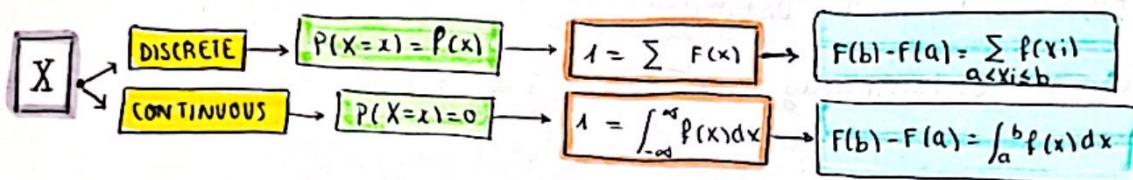
- $E\left[\left(\frac{X - E(X)}{\sigma(X)}\right)^k\right]$:= k -th standardized moment

$k=3$: Skewness
 $k=4$: Kurtosis

• DEF: $\exists f_X: \mathbb{R} \mapsto [0, \infty)$ s.t $F_X(x) = \int_{-\infty}^x f_X(s) ds = P(X \leq x)$ $\rightarrow X$ RV CONTINUOUS

⚠ f_X := "probability density"

- $P(X=x)=0$



• DEF: $F_X(x) = P\left(\bigcap_{i=1}^m (X_i \leq x_i)\right) = P(X_1 \leq x_1, \dots, X_m \leq x_m) = P(\{w: X_i(w) \leq x_1, \dots, X_m(w) \leq x_m\})$:= Joint Distribution Function

where $X: \Omega \mapsto \mathbb{R}^m : X = (X_1, \dots, X_m)$

• DEF: $F_{X,Y}(x,y) :=$ JDF of $Z = (X, Y)$. We call "Marginal Distribution" to:

$$F_X(x) = F_{X,Y}((x, \infty)) = P(X \leq x, Y \leq \infty) = P(X \leq x)$$

• DEF: $Z = (X, Y)$
 $\forall (x, y) \in \mathbb{R}^2 \quad F_Z(x, y) = F_X(x) F_Y(y)$ \Rightarrow X and Y INDEPENDENT

• DEF: $Z = (X, Y)$. $f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}$:= Conditional PMF / Density

- IF Y is discrete: $f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)} = \frac{P(X \leq x, Y=y)}{P(Y=y)}$

$(X, Y) \in \mathbb{R}^2$ RV
AC IR : $P(Y \in A) > 0$

$$F_{X|Y}(x|A) := \frac{P(X \leq x, Y \in A)}{P(Y \in A)}$$

"KASU OROKORRA"

DEF: $g(X)$ function of a RV X .

$$\mathbb{E}(g(X)) = \begin{cases} \sum g(x) f(x), & \text{DISCRETE} \\ \int_{-\infty}^{\infty} g(x) f(x) dx, & \text{CONTINUOUS} \end{cases}$$

«SEQUENCES»

• It's a random vector: $X = (X_1, \dots, X_n)$ for a fixed n , where $\bar{X} = (X_1, \dots, X_n)$.

• $F(X \leq x) = F_1(X_1 \leq x_1) \dots F_n(X_n \leq x_n) \rightarrow$ Independent Sequence

• $F_{X_i} = F_X \rightarrow$ Identically Distributed

• Independent \oplus Identically Distributed : = "I.I.D"

PROPERTIES OF EXPECTATION:

(1) $\alpha \in \mathbb{R} : \mathbb{E}(\alpha X) = \alpha \mathbb{E}(X)$

(2) $\mathbb{E}(X+Y) = \mathbb{E}(X) + \mathbb{E}(Y)$

(3) $X \wedge Y$ INDEP $\Rightarrow \mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$

(4) $\mathbb{E}[X | Y=y] = \int_{-\infty}^{\infty} x f_{X|Y}(x|y) dx$

(5) TOWER PROPERTY: $\mathbb{E}[\mathbb{E}(X|Y)] = \mathbb{E}[X]$

⚠ $g(y) = \mathbb{E}(X|Y=y)$. Define: $g(Y) = \mathbb{E}(X|Y)$

so, $\mathbb{E}[X|Y]$ is another random variable

3. CONCENTRATION

(1) LEARNING FROM DATA:

- Experiment: Randomly pick a swedish person and weighing them
- RV: X represents the weight of the random individual. Assume that the weight is from 0 to 300: $P(0 \leq X \leq 300) = 1$.
- We want to learn: $E[X]$
- n-product experiment: check n people $X = \{X_1, \dots, X_n\}$ where X_i is the weight of person i.
- Estimator: $\frac{1}{n} \sum_{i=1}^n X_i \approx E[X]$ empirical mean

"How much is the empirical mean concentrated around $E(X)$?"

- THEOREM: "Chebychev's inequality"

$$\forall \text{ RV } X, \forall \epsilon > 0, P(|X - E(X)| > \epsilon) \leq \frac{V(X)}{\epsilon^2}$$

$\bar{X}_n :=$ empirical mean

choose $\epsilon = 10$, $X_i \leq 300 \rightarrow E(\bar{X}_n) \leq 300$

$$P(|\bar{X}_n - E(\bar{X}_n)| > 10) \leq V(\bar{X}_n) / 100$$

} ... Lecture 4

(2) CONFIDENCE INTERVAL:

- $P(|\bar{X}_n - E(X)| > \epsilon) \leq \delta$
- $P(|\bar{X}_n - E(X)| < \epsilon) \geq 1 - \delta$
- $P(\bar{X}_n - \epsilon < E(X) < \bar{X}_n + \epsilon) \geq 1 - \delta$

$$I = (\bar{X}_n - \epsilon, \bar{X}_n + \epsilon) \Rightarrow P(E(X) \in I) \geq 1 - \delta$$

"Can We Do Better?":

- THEOREM: "Hoeffding's inequality"

$$\left. \begin{array}{l} X_1, \dots, X_n \stackrel{\text{IID}}{\sim} F \\ P[X_i \in [a, b]] = 1 \\ \forall \epsilon > 0 \end{array} \right\} \Rightarrow P(|\bar{X}_n - E(\bar{X}_n)| > \epsilon) \leq 2e^{-\frac{2n\epsilon^2}{(b-a)^2}}$$

$$E = \sqrt{\frac{(b-a)^2}{2n} \ln \frac{4}{\alpha}} \rightarrow \text{Confidence Interval: } I = (\bar{x}_n - E, \bar{x}_n + E)$$

$P(\text{IE}[\bar{x}_n] \text{ is in } I) \geq 1 - \alpha$

Ex: X_1, \dots, X_n i.i.d Bernoulli(p) $\rightarrow P(\bar{x}_n - p \geq E) \leq e^{-2nE^2}$

- $P(\bar{x}_n - E < p) \geq 1 - e^{-2nE^2} \rightarrow P(\bar{x}_n - E < p < \bar{x}_n + E) \geq 1 - 2e^{-2nE^2}$

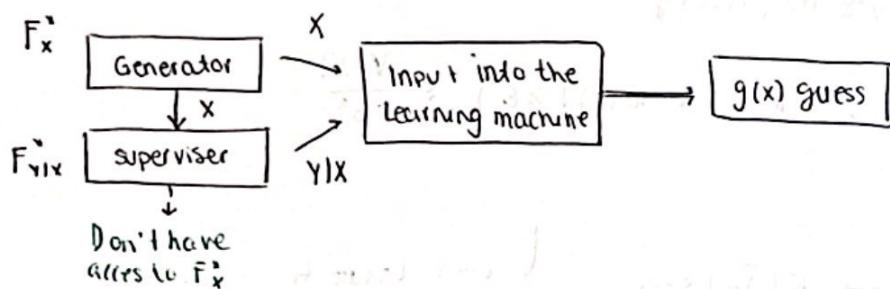
$E = \sqrt{\frac{-1}{2n} \ln \left(\frac{\alpha}{2} \right)}$

BENNETS INEQUALITY:

$$P(|\bar{x}_n - E(\bar{x}_n)| \geq E) \leq 2e^{-\frac{n\alpha^2}{b^2} h\left(\frac{bE}{\alpha^2}\right)}$$

where $h(u) = (1+u) \ln(1+u) - u$

«SUPERVISED LEARNING»



DEF. $L: \mathbb{R}^2 \rightarrow \mathbb{R}$ loss function. Let $(x, y) \sim F_{xy}$. The RISK of a function

$$g: \mathbb{X} \rightarrow \mathbb{Y} \text{ is } \text{IR}[g] := \mathbb{E}[L(g(x), y)]$$

EXAMPLES:

(1) Quadratic loss: $L(a, b) = (a-b)^2 \rightarrow \text{R}(g) = \mathbb{E}[(g(x) - y)^2]$

↳ care about LARGE errors

(2) Absolute loss: $L(a, b) = |a-b| \rightarrow \text{R}(g) = \mathbb{E}[|g(x) - y|]$

↳ care about SMALL errors

(3) 0-1 LOSS: $Y = \{0, 1\} \quad L(a, b) = \begin{cases} 1, & a \neq b \text{ (guess+expected)} \\ 0, & \text{otherwise} \end{cases} \rightarrow \text{R}(g) = P(g(x) \neq y)$

↳ $1 - \text{R}(g)$:= "accuracy"

DEF: "Training Error": $\hat{R}_n(g) = \frac{1}{n} \sum_{i=1}^n L(g(x_i), y_i)$ empirical: $(x_1, y_1), \dots, (x_n, y_n)$

"Testing Error": $\hat{R}_m(g) = \frac{1}{m} \sum_{j=1}^m L(g(x_{n+j}), y_{n+j})$: $(x_{n+1}, y_{n+1}), \dots, (x_{n+m}, y_{n+m})$

«DATA»

- $D_{\text{Tr}}: \{(x_1, y_1), \dots, (x_n, y_n)\} \rightarrow \text{Training Data}$
- $D_{\text{Te}}: \{(x_{n+1}, y_{n+1}), \dots, (x_{n+m}, y_{n+m})\} \rightarrow \text{Testing Data}$

► The LH minimizes the empirical risk:

$$\hat{g} = \arg \min_{g \in \mathcal{G}} \left(\frac{1}{n} \sum_{i=1}^n L(g(x_i), y_i) \right)$$

We test it
on D_{Te}

$$\frac{1}{n} \sum_{j=1}^m L(\hat{g}(x_{n+j}), y_{n+j}) \approx \mathbb{E}_{\text{R}(\hat{g})} L(\hat{g})$$

$\mathbb{E}_{\text{R}(\hat{g})} | D_{\text{Tr}}$

► TRUE RISK: $R(g) = \mathbb{E}[L(g(x), y)]$

«ESTIMATION»

- STATISTIC T : something we can calculate based on data $T: X_n \rightarrow \mathbb{T}$

► EXAMPLE: \hat{g} is a statistic,

► EXAMPLE: Testing Error: $\frac{1}{m} \sum_{i=1}^m L(\hat{g}(x_{n+i}), y_{n+i}) = T(D_{\text{Te}})$

- DEF: $\text{bias}(T) = \mathbb{E}[T] - \theta$ where T is the estimator of θ

- $\text{guess} > \theta$: positive bias
- $\text{guess} < \theta$: negative bias
- $\text{bias}(T) = 0 \rightarrow \text{unbiased}$

possible values
of T

↑ samples ↓ bias

- DEF: $S_e(T) = \sqrt{V(T)}$ "standard error"

- EXAMPLES:
 - $T = \frac{1}{n} \sum x_i \rightarrow S_e(T) = \frac{\sqrt{V(x_i)}}{\sqrt{n}}$
 - $T = x_i \rightarrow S_e(T) = \sqrt{V(x_i)}$

- DEF: "Mean Squared Error":

$$\text{MSE}(T) = \mathbb{E}[(T - \theta)^2] = \mathbb{E}[(T - \mathbb{E}[T])^2] + [\mathbb{E}[T] - \theta]^2 = [S_e(T)]^2 + [\text{bias}(T)]^2$$

► The risk of the estimator w.r.t the quadratic loss can be decomposed as above.

«REGRESSION»

- 1 Find the regression function:

- $(X, Y) \sim F_{X,Y}$ (unknown)
- $L(a, b) = (a - b)^2$ LOSS FUNCTION
- GUESS: f
- $R(f) = \mathbb{E}[(Y - f(X))^2]$ RISK
- $r(x) := \mathbb{E}[Y | X=x]$ (x fixed)

$$\begin{aligned} R(f) &= \mathbb{E}[(Y - r(x)) + (r(x) - f(x))^2] = \dots = \\ &= \underbrace{\mathbb{E}[(Y - r(x))^2]}_{\text{Measures the Noise (CAN'T REDUCE IT)}} + \underbrace{\mathbb{E}[(r(x) - f(x))^2]}_{\text{BIGS (CAN REDUCE IT)}} \end{aligned}$$

Best Guess $\rightarrow R(r)$ smallest

2 FIND h :

- Assume $E[(Y - r(x))^2] = 0$

"PATTERN RECOGNITION":

- $Y \in \{0, 1\}$ class

- $L(a|b) = \begin{cases} 1, & a \neq b \\ 0, & a = b \end{cases}$

- $r(x) = E(Y|X=x) = P(Y=1|X=x)$ Regression Function

- $\hat{h}^*(x) = \begin{cases} 1, & r(x) > 1/2 \\ 0, & r(x) \leq 1/2 \end{cases} := \text{Bayes Classifying Rule}$

$\hookrightarrow R(\hat{h}^*) \leq R(h)$ for any guessing function h

"MAXIMUM LIKELIHOOD":

- $\mathcal{P} = \{P_\alpha | X\}: \alpha \in \mathbb{R}^d$ parametric model

- $L(x, \alpha) = -\ln P_\alpha(x)$ log-loss
↑ parameter

- $X \sim P_{\alpha^*} \in \mathcal{P}: R(\alpha) = E[L(X, \alpha)] = E[-\ln P_\alpha(x)] = -\int \ln P_\alpha(x) P_{\alpha^*}(dx)$

- $x_1, \dots, x_n \stackrel{\text{iid}}{\sim} P_{\alpha^*}: \hat{R}(\alpha) = \frac{1}{n} \sum_{i=1}^n -\ln P_\alpha(x_i)$ empirical Risk

EXAMPLES:

(1) Linear Regression:

$$P_\alpha(y|x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(y-(\alpha_0 + \alpha_1 x))^2}{2\sigma^2}} \quad \Rightarrow \quad \hat{R}(\alpha) = \ln(1/\sigma) + \frac{1}{n} \sum_{i=1}^n \frac{(y_i - (\alpha_0 + \alpha_1 x_i))^2}{2\sigma^2}$$

- $\alpha = (\alpha_0, \alpha_1, \sigma)$

(2) Logistic Regression

- $P_\alpha(y|x) = \text{Bernoulli}(p) = p^y (1-p)^{1-y}$

- $G(x) = \frac{1}{1+e^{-x}}$ Logistic Function

$$L(z, \alpha) := -\ln P_\alpha(z) \quad \text{where } P_\alpha \text{ proposal density for our data}$$

- $\hat{R}(\alpha) = \frac{1}{n} \sum_{i=1}^n -\ln (P_\alpha(x_i))$

- $R(\alpha) = E[\ln(P_\alpha(x))]$

« GENERATING RANDOM VARIABLES » (Markov Chain)

1 PRNG: PSEUDO RANDOM NUMBER GENERATOR

- DEF: A sequence z_0, u_1, \dots is pseudorandom if for any $a \in M = \{0, 1, \dots, M-1\}$

$$\frac{N_n(a)}{n} \rightarrow \frac{1}{M}.$$

- EX: $0, 1, 0, 1, \dots$ is pseudorandom on $\{0, 1\}$

- DEF: A congruential generator with parameters (a, b, M) :

- $D(x) = (ax + b) \bmod M$
- start at some u_0 - random seed

- $u_i = D(u_{i-1})$

- Period: of D starting at u_0 is the smallest T_0 such that $u_i + T_0 = u_i \forall i$
- Full period: on $\{0, 1, \dots, M-1\}$ is M (you get all the numbers)

$$\rightarrow u_0 = 0 \xrightarrow{\text{steps}} \text{get } 0 : \frac{N_{K \cdot M}(0)}{K \cdot M} = \frac{1}{K \cdot M} = \frac{1}{N}$$

\hookrightarrow so: u_0, u_1, \dots is pseudorandom

IMPROVE IT:

- Large M
- Assume full period: u_0, \dots, u_M
- Fix K s.t $K \mid M$
- Construct: $v_i = \lfloor u_i \frac{K}{M} \rfloor$ (floor in python)
 - v_i pseudorandom on $\{0, 1, \dots, K-1\}$
 - Period: M

⚠ To create $\text{unif}([0, 1])$ we take:

$$v_i = \frac{u_i}{M}$$

« SAMPLING »

1 INVERSION SAMPLING

- We want to sample from F

① calculate F^{-1} (if we can)

② Generate $\text{Unif}([0, 1])$

③ consider: $X = F^{-1}(u) \rightarrow X \sim F$

$$F_X(x) = P(X \leq x) = P(F^{-1}(u) \leq x) = P(u \leq F(x)) = F(x)$$

$u \sim \text{uniform}$

$$F_u(u) = u$$

• EX: $F(x) = 1 - e^{-\lambda x}$ exponential
 $y = 1 - e^{-\lambda x} \rightarrow e^{-\lambda x} = 1 - y \rightarrow -\lambda x = \ln(1-y) \rightarrow x = -\frac{1}{\lambda} \ln(1-y) = F^{-1}(y)$

2 ACCEPT - REJECT METHOD

- Input target density: f
 - Sampling density: g
- We need a sampling density that dominates f : $f \leq Ng(x)$
for all $x \in \mathbb{N} < 1$

[1] Draw $X \sim g$ and calculate

$$r(X) = \frac{f(X)}{Ng(X)}$$

[2] Draw a $U \sim \text{Unif}[0,1]$

[3] $\begin{cases} \text{accept } X, & U \leq r(X) \\ \text{reject } X, & \text{otherwise + try again} \end{cases}$

↳ $I = \begin{cases} 1, & U \leq r(X) \\ 0, & \text{otherwise} \end{cases}$

↳ If $I=1$ (ACCEPT):

$$P(X \leq x | I=1) = F(x)$$

↑
• Bayes theorem: $f_{I|X}(x|1) = \frac{f_{I|X}(1|x) f_X(x)}{f_I(1)}$

• $f_{I|X}(1|x) = P(U \leq r(x) | X=x) = r(x)$

• $f_I(1) = \int f_{I|X}(1|x) f_X(x) dx = \int r(x) g(x) dx = \int \frac{f(x)}{Ng(x)} g(x) dx = \frac{1}{N} \int f(x) dx = \frac{1}{N}$

• $f_X(x) := g(x)$

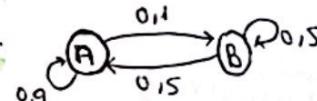
• DEF: A stochastic process is indexed set of random variables $\{X_i | i \in I\}$.

IF $I = \{1, 2, 3, \dots, N\}$, we say it is discrete.

• DEF: A Markov process $\{X_i | i \in \mathbb{N}\}, X_i \in \mathcal{X} = \{S_1, \dots, S_N\}$:

$$P(X_t=y | x_1, \dots, x_{t-1}) = P(X_t=y | x_{t-1}) := \text{Markov Property}$$

↳ Only cares about the last one (previous don't affect)

• EX:  $P = \begin{pmatrix} 0.9 & 0.1 \\ 0.15 & 0.85 \end{pmatrix}$ A = "dry", B = "wet"

► $P_{t-1} = [P(X_{t-1}=\text{dry}), P(X_{t-1}=\text{wet})]$ we know

$$\bullet P(X_t=\text{dry}) = P(X_t=\text{dry} | X_{t-1}=\text{dry}) P(X_{t-1}=\text{dry}) + P(X_t=\text{dry} | X_{t-1}=\text{wet}) P(X_{t-1}=\text{wet})$$

► $P_t = P_{t-1} \cdot P = P_0 \cdot P^t$

• DEF: $X = \{S_1, \dots, S_n\}$. We say that $S_i \rightarrow S_j$ "state S_i communicates with S_j "

if $P(X_t = S_j | X_0) > 0$ for some t

• DEF: $S_i \rightarrow S_j \wedge S_j \rightarrow S_i \rightarrow$ They intercommunicate

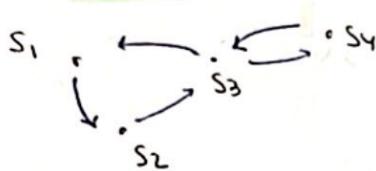
• DEF: • Everything communicates \rightarrow Irreducible

• If not \rightarrow Reducible

• DEF: Return times: How many steps to go back (get stuck):

$$\pi(x) = \inf_{t \in \mathbb{N}} \{ P^t(x, x) > 0 \}$$

↳ EXAMPLE:



$$\pi(S_1) = \{3, 5, 7, \dots\}$$

$$\pi(S_3) = \{2, 3, 5, \dots\}$$

• DEF: $\text{gcd}(\pi(x)) :=$ "Period of the state x "

• DEF: $\text{gcd}(\pi(S_i)) = 1$
for all states \rightarrow The chain is APERIODIC

• THEOREM:

- Irreducible
- Aperiodic

$$P_0 P^t \xrightarrow{\text{(UNIQUE)}} \pi \text{ mixing}$$

$$\pi = \pi P \xrightarrow{\text{STATIONARY DISTRIBUTION}}$$

$$\pi^T = (\pi P)^T = P^T \pi^T \xrightarrow{\text{Eigen vector of } P^T \text{ with eigenvalue 1}}$$

«LINEAR CLASSIFIER»

• $x \in \mathbb{R}^d$, $y \in \{-1, 1\}$, $w \in \mathbb{R}^d$, $c \in \mathbb{R}$ (parameters)

$$g_w(x) = \text{sgn}(\langle w, x \rangle + c)$$

$$L(y, g_w(x)) = -\frac{\text{sgn}(\langle w, x \rangle + c) y + 1}{2} = \begin{cases} 1, & y \neq g_w(x) \\ 0, & y = g_w(x) \end{cases}$$

$$\langle \tilde{w}, \tilde{x} \rangle = \langle w, x \rangle + c$$

$$\cdot \tilde{w} = (w_1, \dots, w_d, c)$$

$$\cdot \tilde{x} = (x_1, \dots, x_d, 1)$$

$\therefore 0-1 LOSS$

Two classes

$\langle \tilde{w}, \tilde{x} \rangle = 0$: Plane from the origin

How do we find a linear classifier?

Margin

① PERCEPTION ALGORITHM:

(1) set $w = (0, \dots, 0)$

(2) Points (x_i, y_i) $i=1, \dots, n$

► If $\exists (x_i, y_i)$ s.t $\langle w, x_i \rangle y_i < 0$, \Rightarrow misclassified

$$w \leftarrow w + \sum_{\epsilon \in \{1, -1\}} y_i x_i$$

(3) Repeat until all satisfy: $\langle w, x_i \rangle y_i > 0$

② SUPPORT VECTOR MACHINE:

We want to solve the problem of uniqueness

(1) Take $\|w\| = 1$: $\langle w, x \rangle = 1$

$$x - \tilde{x} \parallel w: \| (x - \tilde{x}) \cdot w \| \rightarrow \|x - \tilde{x}\| \|w\| = 1$$

$$\Delta \text{ If } \|w\| \neq 1 \rightarrow \|x - \tilde{x}\| = 1/\|w\|$$

(2) Consider the function:

$$\text{Hinge loss} := \max \{0, 1 - \langle w, x \rangle y\} = \begin{cases} 1 - \langle w, x \rangle y, & \langle w, x \rangle y < 1 \\ 0, & \text{otherwise} \end{cases}$$

(3) Minimize: $\|w\|$ under the constraint: " $\langle w, x_i \rangle y_i > 1$ "

"SOFT MARGIN SVM":

$$\min_w \|w\|^2 + C \sum_{j=1}^n \max \{0, 1 - \langle w, x_j \rangle y_j\}$$

③ THE KERNEL TRICK:

► Perception update rule: $w \leftarrow w + x_i y_i$

we will find w of the form: $w = \sum c_i x_i$

► $\langle w, x_k \rangle = \sum c_i \langle x_i, x_k \rangle = \sum c_i k_{ik} \xrightarrow{z_i = \phi(x_i)} \langle \tilde{w}, z_k \rangle = \sum \tilde{c}_i \tilde{k}_{ik}$

"TRICK": If we have a function $K(a, b)$ such that $\tilde{k}_{ik} = K(x_i, x_k)$ (we forget about ϕ).

K := Kernel Function

EXAMPLES:

• $K(a, b) = \langle a, b \rangle \rightarrow$ linear

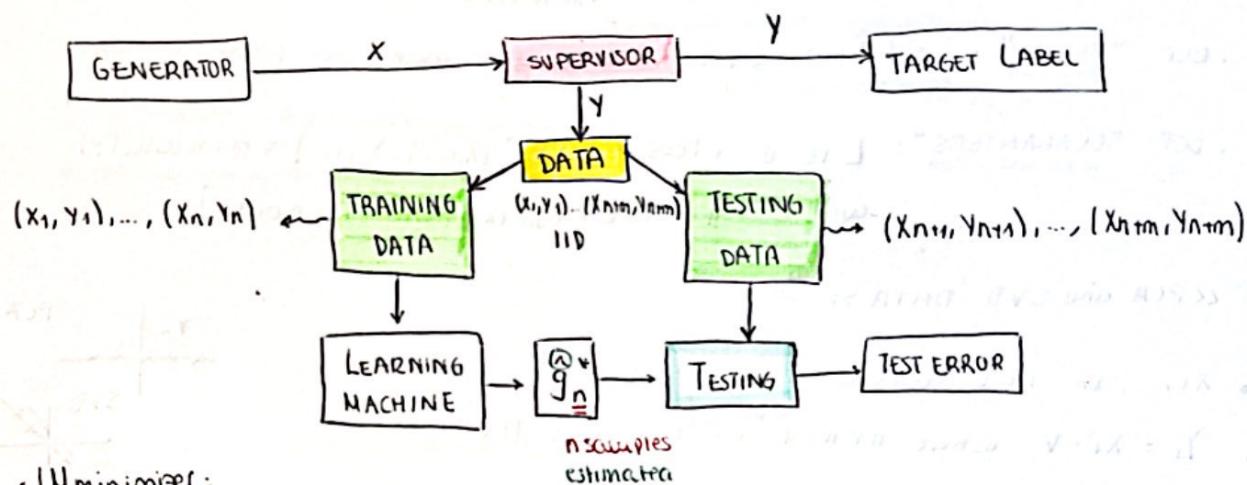
• $K(a, b) = (\gamma \langle a, b \rangle + m)^r \rightarrow$ polynomial

• $K(a, b) = e^{-\|a-b\|^2} \rightarrow$ radial basis function

"RADICAL BASIS FUNCTION":

$$\langle w, x \rangle = \lambda = \sum c_i K(x_i, x)$$

« TRAIN - TEST PROCEDURE »



LH minimizer:

$$\hat{g}_n^* = \arg \min_{g \in \mathcal{U}} \frac{1}{n} \sum_{i=1}^n L(g(x_i), y_i) \quad \text{only training data}$$

Test error:

$$\frac{1}{m} \sum_{i=1}^m L(\hat{g}_n^*(x_{m+i}), y_{m+i})$$

PYTHON: `lr.fit(X-train, Y-train)` := LH

$$\hat{g}_n^*(x) = lr.predict(x)$$

`lr.score(X-test, Y-test)` := Testing Error

« TEST ERROR »

- Once trained, we consider \hat{g}_n^* as fixed: $E[L(\hat{g}_n^*(x), Y) | \text{Training data}]$ is the expected loss.

BIAS VARIANCE:

$$\begin{aligned} E[|Y - \hat{g}(x)|^2 | x=x] &= E[|Y - r(x) + r(x) - E[\hat{g}(x) | x=x]|^2] + E[(\hat{g}(x) | x=x) - E[\hat{g}(x) | x=x]]^2 \\ &= E[|Y - r(x)|^2 | x=x] + E[|r(x) - E[\hat{g}(x) | x=x]|^2 | x=x] + \\ &\quad + E[|E[\hat{g}(x) | x=x] - \hat{g}(x)|^2 | x=x] \end{aligned}$$

◻ Noise

◻ Bias squared: Bias^2

◻ Variance

• DEF: "PRECISION": Given $Y \in \{0, 1, \dots, k-1\}$ and \hat{g} , then the precision for class $l \in \{0, 1, \dots, k-1\}$: $P(Y=l | \hat{g}(x)=l)$

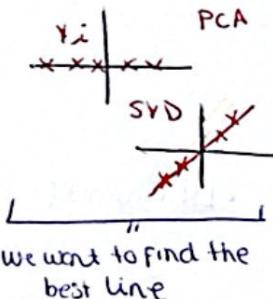
I have covid test tells I have covid

• DEF: "RECALL": $P(\hat{g}(x)=l | Y=l)$ → How many sick people can detect

• DEF: "GUARANTEES": L is 0-1 loss. $L(\hat{g}_n^*(x_{n+i}), y_{n+i}) \sim \text{Bernoulli}(p)$
where $p = \mathbb{E}[L(\hat{g}_n^*(x_n), y_n) | \text{Train data}]$

«PCA and SVD DATA»

- x_1, \dots, x_n 1.D samples
- $y_i = x_i \cdot v$ where $\|v\|=1$ "coordinate on the line"
- $y_i v$: "real coordinate"
- $\sum_{i=1}^n y_i$ assume $\bar{y}_n = 0$ empirical mean
- $\frac{1}{n} \sum (y_i - \bar{y}_n)^2 = \frac{1}{n} \sum (x_i \cdot v)^2$ empirical variance



GOAL: select v s.t. variance is maximize

$$v_1 = \underset{\|v\|=1}{\operatorname{argmax}} \left(\frac{1}{n} \sum (x_i \cdot v)^2 \right) \rightarrow 1^{\text{st}} \text{ SINGULAR VECTOR}$$

$$\hookrightarrow \sigma_1 := \sqrt{\sum (x_i \cdot v_1)^2} \rightarrow 1^{\text{st}} \text{ SINGULAR VALUE}$$

$$v_2 = \underset{\|v\|=1, v \perp v_1}{\operatorname{argmax}} \left(\frac{1}{n} \sum (x_i \cdot v)^2 \right) \rightarrow 2^{\text{nd}} \text{ SINGULAR VECTOR}$$

$$\hookrightarrow \sigma_2 := \sqrt{\sum (x_i \cdot v_2)^2} \rightarrow 2^{\text{nd}} \text{ SINGULAR VALUE}$$

$$v_3 = \underset{\|v\|=1, v \perp v_1, v \perp v_2}{\operatorname{argmax}} \left(\frac{1}{n} \sum (x_i \cdot v)^2 \right) \rightarrow 3^{\text{rd}} \text{ SINGULAR VECTOR}$$

$$\hookrightarrow \sigma_3 := \sqrt{\sum (x_i \cdot v_3)^2} \rightarrow 3^{\text{rd}} \text{ SINGULAR VALUE}$$

(...)

TWO STOPPING CONDITIONS

→ [1] Reach the size of the space

→ [2] There is no vector $v \perp v_1, \dots, v_r$ or $\sum (x_i \cdot v)^2 = 0$ for all vectors, $r \leq \min(m, n)$

- $A = \begin{pmatrix} -x_1 - \\ \vdots \\ -x_n - \end{pmatrix}_{n \times m} := \text{Gram Matrix}$

- $AU = \begin{pmatrix} x_1 \cdot v \\ \vdots \\ x_n \cdot v \end{pmatrix} \implies \|AU\|^2 = \sum_{i=1}^n (x_i \cdot v)^2$

- $v_1 = \arg \max_{\|v\|=1} (\|Av\|^2) \rightarrow \alpha_1 = \|Av_1\|$

- Once $AU = 0$ for all $v \perp v_1, \dots, v \perp v_k$ choose a basis for the null space of $A|_{U_{k+1}, \dots, U_r}$ where $r = \min(m, n)$ so $\alpha_{k+1} = \dots = \alpha_r = 0$.

- THEOREM:** $m \leq n$ (more data than dimension). We have v_1, \dots, v_m . Choose $k \leq m$:

- $W_k = \sum_{j=1}^k (x_i \cdot v_j) v_j$ k -dimensional subspace (passes from 0) "The best fit subspace"

- $\sum_{i=1}^n \|x_i - W_k\|^2 \leq \sum_{i=1}^n \|x_i - \text{Proj}_{H_k}(x_i)\|^2$ where H_k is any k -dim subspace (with any other the error will be bigger)

- v_1, \dots, v_m basis in $\mathbb{R}^m \rightarrow x_i = \sum_{j=1}^m (x_i \cdot v_j) v_j$

$$AU_j v_j^\top = \begin{pmatrix} x_1 \cdot v_j \\ \vdots \\ x_n \cdot v_j \end{pmatrix} v_j^\top = \begin{pmatrix} -(x_1 \cdot v_j) v_j \\ \vdots \\ -(x_n \cdot v_j) v_j \end{pmatrix}$$

$$\sum_{j=1}^m AU_j v_j^\top = \begin{pmatrix} \sum_{j=1}^m (x_1 \cdot v_j) v_j \\ \vdots \\ \sum_{j=1}^m (x_n \cdot v_j) v_j \end{pmatrix} = A$$

$$U_j = \begin{cases} \frac{AU_j}{\alpha_j}, & \alpha_j > 0 \\ 0, & \text{otherwise} \end{cases} \rightarrow \text{"Left Singular Vector"}$$

$$\alpha_j U_j = AU_j$$

- $A = \sum_{j=1}^m \alpha_j U_j V_j^\top = \sum_{j=1}^m AU_j V_j^\top$

- $U = (U_1 | U_2 | \dots | U_m)_{n \times m}$

- $D = \begin{pmatrix} \alpha_1 & 0 \\ 0 & \alpha_m \end{pmatrix}$

- $V = (V_1 | \dots | V_m)_{m \times m}$

- "Low Rank Approximation": $A_K = \sum_{j=1}^K \alpha_j U_j V_j^\top$ if $\alpha_j > 0$ $\rightarrow \text{Rank}(A_K) = K$

- "Error": $\|A - A_K\|_F^2 = \sum (a_{ij} - \hat{a}_{ij})^2$

- $A = \sum_{i=1}^m \alpha_i U_i V_i^\top$ $\rightarrow A - A_K = \sum_{i=K+1}^m \alpha_i U_i V_i^\top$

$$\begin{aligned} \|\mathbf{A} - \mathbf{A}_k\|_2^2 &= \sum_{j=1}^m \|\tilde{\mathbf{x}}_j\|_2^2 = \sum_{j=1}^m \sum_{i=1}^n |\tilde{x}_i \cdot v_j|^2 = \sum_{i=1}^n \underbrace{\sum_{j=1}^m (\mathbf{A} - \mathbf{A}_k) v_i v_j^T}_{\sum_{j=k+1}^m \mathbf{A}_j v_i v_j^T} = \\ &= \sum_{j=k+1}^m \|\mathbf{A}_j v_i\|_2^2 = \sum_{j=k+1}^m \sigma_j^2 \end{aligned}$$

• **SELECT K LARGE** $\rightarrow \sum_{j=k+1}^m \sigma_j^2 \text{ small}$

- $\|\mathbf{AV}\|_2^2 = \langle \mathbf{AV}, \mathbf{AV} \rangle = \langle \mathbf{A}^T \mathbf{AV}, \mathbf{V} \rangle$
- $\tilde{\mathbf{A}} = \mathbf{A}^T \mathbf{A} : \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m > 0$
- $\mathbf{w}_1, \dots, \mathbf{w}_m$ Eigen vectors
- $\mathbf{V} = \sum (\mathbf{v} \cdot \mathbf{w}_i) \mathbf{w}_i$
- $\tilde{\mathbf{A}} \mathbf{V} = \sum_{i=1}^m \lambda_i (\mathbf{v} \cdot \mathbf{w}_i) \mathbf{w}_i$
- $\langle \tilde{\mathbf{A}} \mathbf{V}, \mathbf{V} \rangle = \sum_{i=1}^m \lambda_i (\mathbf{v} \cdot \mathbf{w}_i)^2$

$\max \|\mathbf{AV}\|_2^2 = \lambda_1$

⚠ CENTER THE DATA! ⚡ STANDARDIZE DATA

⚠ Before logistic
Regression:
STANDARDIZE

{ • CLEAN (quit outliers)
• TRANSFORM
• STANDARDIZE } \Rightarrow TRAINING DATA

⚠ PIPELINE (recomendado)

► CALIBRATOR ERROR: $c(f) = \sqrt{\mathbb{E}[(Y \mid f(x)) - f(x)]^2}$

$Y \in \{0, 1\} \rightarrow c(f) = \sqrt{\mathbb{E}[(I(Y=1) \mid f(x)) - f(x)]^2}$

If it is perfectly calibrated $\rightarrow c(f) = 0$

Construct features from data to make a model more powerful

► FEATURE ENGINEERING $\xrightarrow{\quad}$ Transformation of features: replace x with $\log(x)$
 $\xrightarrow{\quad}$ Transform the target Y price $\rightarrow Y' = \log(Y)$

« TEST DATA »

. Dictionary: enumeration of all possible words in your data.

. TF-IDF := Term Frequency - Inverse Document Frequency

. $f_{t,d}$: count how many times t appears in d , document

. $TF(t,d) = f_{t,d} / \sum_{s \in d} f_{s,d} \rightarrow$ How many words in d

. $IDF(t) = \log \left(\frac{N}{1 + \text{tf}(t, \text{every})} \right) \rightarrow$ # documents in how many documents this word appears

$TF \cdot IDF(t,d) = TF(t,d) \cdot IDF(t)$

► $TF \uparrow$: very often
 ► $IDF \uparrow$: not in every doc

PRODUCT ↑:
 "interesting word"