



## SIMULATION:

• If  $X$  has CDF  $F(x) = P(X \leq x)$

ex.

$$X \sim \text{Exp}(\lambda) \rightarrow f_X(x) = \lambda e^{-\lambda x}, x > 0$$

$$F_X(x) = \int_0^x \lambda e^{-\lambda t} dt = -e^{-\lambda t} \Big|_0^x = 1 - e^{-\lambda t}$$

$$(Y = 1 - e^{-\lambda X} = 1 - Y = e^{-\lambda X} = \ln(1 - Y) = -\lambda X \rightarrow X = -\frac{\ln(1 - Y)}{\lambda})$$

$$\cdot Q(u) = F_X^{-1}(u)$$

• How to approx.  $F(x)$ :

1. Generate  $X_1, \dots, X_B$  iid  $X$

$$2. \hat{F}(x_0) = \frac{1}{B} \sum_{i=1}^B (X_i \leq x_0)$$

→ QUANTILE FUNCTION:

$$Q(u) = \inf \{x : F(x) \geq u\}$$

$$f(x)$$

$$F(x)$$

$$x$$

$$u_1 \in (0,1)$$

$$u_2 \in (0,1)$$

$$u_3 \in (0,1)$$

$$u_4 \in (0,1)$$

$$u_5 \in (0,1)$$

$$u_6 \in (0,1)$$

$$u_7 \in (0,1)$$

$$u_8 \in (0,1)$$

$$u_9 \in (0,1)$$

$$u_{10} \in (0,1)$$

$$u_{11} \in (0,1)$$

$$u_{12} \in (0,1)$$

$$u_{13} \in (0,1)$$

$$u_{14} \in (0,1)$$

$$u_{15} \in (0,1)$$

$$u_{16} \in (0,1)$$

$$u_{17} \in (0,1)$$

$$u_{18} \in (0,1)$$

$$u_{19} \in (0,1)$$

$$u_{20} \in (0,1)$$

$$u_{21} \in (0,1)$$

$$u_{22} \in (0,1)$$

$$u_{23} \in (0,1)$$

$$u_{24} \in (0,1)$$

$$u_{25} \in (0,1)$$

$$u_{26} \in (0,1)$$

$$u_{27} \in (0,1)$$

$$u_{28} \in (0,1)$$

$$u_{29} \in (0,1)$$

$$u_{30} \in (0,1)$$

$$u_{31} \in (0,1)$$

$$u_{32} \in (0,1)$$

$$u_{33} \in (0,1)$$

$$u_{34} \in (0,1)$$

$$u_{35} \in (0,1)$$

$$u_{36} \in (0,1)$$

$$u_{37} \in (0,1)$$

$$u_{38} \in (0,1)$$

$$u_{39} \in (0,1)$$

$$u_{40} \in (0,1)$$

$$u_{41} \in (0,1)$$

$$u_{42} \in (0,1)$$

$$u_{43} \in (0,1)$$

$$u_{44} \in (0,1)$$

$$u_{45} \in (0,1)$$

$$u_{46} \in (0,1)$$

$$u_{47} \in (0,1)$$

$$u_{48} \in (0,1)$$

$$u_{49} \in (0,1)$$

$$u_{50} \in (0,1)$$

$$u_{51} \in (0,1)$$

$$u_{52} \in (0,1)$$

$$u_{53} \in (0,1)$$

$$u_{54} \in (0,1)$$

$$u_{55} \in (0,1)$$

$$u_{56} \in (0,1)$$

$$u_{57} \in (0,1)$$

$$u_{58} \in (0,1)$$

$$u_{59} \in (0,1)$$

$$u_{60} \in (0,1)$$

$$u_{61} \in (0,1)$$

$$u_{62} \in (0,1)$$

$$u_{63} \in (0,1)$$

$$u_{64} \in (0,1)$$

$$u_{65} \in (0,1)$$

$$u_{66} \in (0,1)$$

$$u_{67} \in (0,1)$$

$$u_{68} \in (0,1)$$

$$u_{69} \in (0,1)$$

$$u_{70} \in (0,1)$$

$$u_{71} \in (0,1)$$

$$u_{72} \in (0,1)$$

$$u_{73} \in (0,1)$$

$$u_{74} \in (0,1)$$

$$u_{75} \in (0,1)$$

$$u_{76} \in (0,1)$$

$$u_{77} \in (0,1)$$

$$u_{78} \in (0,1)$$

$$u_{79} \in (0,1)$$

$$u_{80} \in (0,1)$$

$$u_{81} \in (0,1)$$

$$u_{82} \in (0,1)$$

$$u_{83} \in (0,1)$$

$$u_{84} \in (0,1)$$

$$u_{85} \in (0,1)$$

$$u_{86} \in (0,1)$$

$$u_{87} \in (0,1)$$

$$u_{88} \in (0,1)$$

$$u_{89} \in (0,1)$$

$$u_{90} \in (0,1)$$

$$u_{91} \in (0,1)$$

$$u_{92} \in (0,1)$$

$$u_{93} \in (0,1)$$

$$u_{94} \in (0,1)$$

$$u_{95} \in (0,1)$$

$$u_{96} \in (0,1)$$

$$u_{97} \in (0,1)$$

$$u_{98} \in (0,1)$$

$$u_{99} \in (0,1)$$

$$u_{100} \in (0,1)$$

$$u_{101} \in (0,1)$$

$$u_{102} \in (0,1)$$

$$u_{103} \in (0,1)$$

$$u_{104} \in (0,1)$$

$$u_{105} \in (0,1)$$

$$u_{106} \in (0,1)$$

$$u_{107} \in (0,1)$$

$$u_{108} \in (0,1)$$

$$u_{109} \in (0,1)$$

$$u_{110} \in (0,1)$$

$$u_{111} \in (0,1)$$

$$u_{112} \in (0,1)$$

$$u_{113} \in (0,1)$$

$$u_{114} \in (0,1)$$

$$u_{115} \in (0,1)$$

$$u_{116} \in (0,1)$$

$$u_{117} \in (0,1)$$

$$u_{118} \in (0,1)$$

$$u_{119} \in (0,1)$$

$$u_{120} \in (0,1)$$

$$u_{121} \in (0,1)$$

$$u_{122} \in (0,1)$$

$$u_{123} \in (0,1)$$

$$u_{124} \in (0,1)$$

$$u_{125} \in (0,1)$$

$$u_{126} \in (0,1)$$

$$u_{127} \in (0,1)$$

$$u_{128} \in (0,1)$$

$$u_{129} \in (0,1)$$

$$u_{130} \in (0,1)$$

$$u_{131} \in (0,1)$$

$$u_{132} \in (0,1)$$

$$u_{133} \in (0,1)$$

$$u_{134} \in (0,1)$$

$$u_{135} \in (0,1)$$

$$u_{136} \in (0,1)$$

$$u_{137} \in (0,1)$$

$$u_{138} \in (0,1)$$

$$u_{139} \in (0,1)$$

$$u_{140} \in (0,1)$$

$$u_{141} \in (0,1)$$

$$u_{142} \in (0,1)$$

$$u_{143} \in (0,1)$$

$$u_{144} \in (0,1)$$

$$u_{145} \in (0,1)$$

$$u_{146} \in (0,1)$$

$$u_{147} \in (0,1)$$

$$u_{148} \in (0,1)$$

$$u_{149} \in (0,1)$$

$$u_{150} \in (0,1)$$

$$u_{151} \in (0,1)$$

$$u_{152} \in (0,1)$$

$$u_{153} \in (0,1)$$

$$u_{154} \in (0,1)$$

$$u_{155} \in (0,1)$$

$$u_{156} \in (0,1)$$

$$u_{157} \in (0,1)$$

$$u_{158} \in (0,1)$$

$$u_{159} \in (0,1)$$

$$u_{160} \in (0,1)$$

$$u_{161} \in (0,1)$$

$$u_{162} \in (0,1)$$

$$u_{163} \in (0,1)$$

$$u_{164} \in (0,1)$$

$$u_{165} \in (0,1)$$

$$u_{166} \in (0,1)$$

$$u_{167} \in (0,1)$$

$$u_{168} \in (0,1)$$

$$u_{169} \in (0,1)$$

$$u_{170} \in (0,1)$$

$$u_{171} \in (0,1)$$

$$u_{172} \in (0,1)$$

$$u_{173} \in (0,1)$$

$$u_{174} \in (0,1)$$

$$u_{175} \in (0,1)$$

$$u_{176} \in (0,1)$$

$$u_{177} \in (0,1)$$

$$u_{178} \in (0,1)$$

$$u_{179} \in (0,1)$$

$$u_{180} \in (0,1)$$

$$u_{181} \in (0,1)$$

$$u_{182} \in (0,1)$$

$$u_{183} \in (0,1)$$

$$u_{184} \in (0,1)$$

$$u_{185} \in (0,1)$$

$$u_{186} \in (0,1)$$

$$u_{187} \in (0,1)$$

$$u_{188} \in (0,1)$$

$$u_{189} \in (0,1)$$

$$u_{190} \in (0,1)$$

$$u_{191} \in (0,1)$$

$$u_{192} \in (0,1)$$

$$u_{193} \in (0,1)$$

$$u_{194} \in (0,1)$$

$$u_{195} \in (0,1)$$

$$u_{197} \in (0,1)$$

$$u_{199} \in (0,1)$$

$$u_{200} \in (0,1)$$

$$u_{201} \in (0,1)$$

$$u_{202} \in (0,1)$$

$$u_{203} \in (0,1)$$

$$u_{204} \in (0,1)$$

$$u_{205} \in (0,1)$$

$$u_{206} \in (0,1)$$

$$u_{207} \in (0,1)$$

$$u_{208} \in (0,1)$$

$$u_{209} \in (0,1)$$

$$u_{210} \in (0,1)$$

$$u_{211} \in (0,1)$$

$$u_{212} \in (0,1)$$

$$u_{213} \in (0,1)$$

F TEST;  $\rightarrow$  TEST SMALL VS BIG MODEL (IS SIGNIFICANT TO UPGRADe?)

$$F = \frac{\frac{\|\hat{Y} - \hat{Y}_0\|^2}{p-1}}{\frac{\|\hat{Y} - \hat{Y}\|^2}{n-p}} = \frac{\frac{\sum (\hat{Y}_i - \bar{Y})^2}{p-1}}{\frac{\sum (Y_i - \hat{Y}_i)^2}{n-p}}$$

EXPLAINS VARIANCE  
UNEXPLAINED VARIANCE

$H_0$  REJECTED ( $\Rightarrow F > F_{\alpha, p-1, n-p}$ )  
 $\hookrightarrow$  FROM COMPUTERS: DETECT IF P-VALUE  $< \alpha$

IN R:

$$\cdot \text{DOF}_1 = n - p_1 - 1 = 200 - 2 - 1 = 197 \quad | \quad \cdot \text{DOF}_2 = n - p_2 - 1 = 200 - 3 - 1 = 196$$

$$\cdot \text{RSS} = \sum (Y_i - \hat{Y}_i)^2 = e^T e \quad | \quad \cdot \text{Df} = \text{DOF}_2 - \text{DOF}_1 \quad | \quad \cdot \text{SUM OF Sq} = \text{RSS}_2 - \text{RSS}_1$$

$$\rightarrow F = 0.0034 \rightarrow P[X > 0.034] \approx 0.95 \approx \text{p-value} > \alpha \rightarrow H_0 \text{ NOT REJ}$$

PREDICTIONS:

$$\cdot \text{FITTING VALUE OF } X_f \text{ IS: } \hat{Y}_f = X_f \cdot \hat{\beta} \rightarrow \hat{Y}_f = X_f \cdot \hat{\beta} \sim N(X_f \cdot \hat{\beta}, \sigma^2 \cdot X_f (X^T X)^{-1} X_f^T)$$

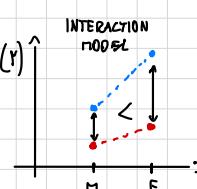
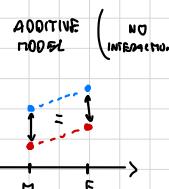
$$\rightarrow \text{C.I.}_{1-\alpha}: X_f \cdot \hat{\beta} \pm t_{\frac{\alpha}{2}, n-p} \sqrt{\frac{e^T e}{n-p} \cdot X_f (X^T X)^{-1} X_f^T}$$

$$\rightarrow \text{PREDICTION INTERVAL: } X_f \cdot \hat{\beta} \pm t_{\frac{\alpha}{2}, n-p} \sqrt{\frac{E^T E}{n-p} (1 + X_f (X^T X)^{-1} X_f^T)}$$

A BIT LARGER THAN C.I.

GENERAL INTERACTION MODEL:

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 1 & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 1 & 0 \\ \vdots & \vdots & \vdots & \ddots \\ 1 & 0 & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 1 & 1 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}$$



{ ANOVA (SMALL, LARGE) GIVES A TEST FOR NULL  $H_0$ ,  $H_0: \text{ALL INTERACTIONS} = 0$

MEASURES OF FIT:

$$\cdot R^2 = \frac{\text{NUM}}{\text{DEN}} = \frac{\|\hat{Y} - \hat{Y}_0\|^2}{\|Y - \hat{Y}_0\|^2} = \frac{\sum (\hat{Y}_i - \bar{Y})^2}{\sum (Y_i - \bar{Y})^2} = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}$$

FREQ.  
BT MEAN OF DATA

$$\begin{aligned} \cdot \text{NUM: } & \|\hat{Y} - \hat{Y}_0\|^2 \\ & \|Y - \hat{Y}_0\|^2 - \|Y - \hat{Y}\|^2 \\ & = TSS - E^T E \end{aligned}$$

$$\cdot \text{DEN: } \|Y - \hat{Y}_0\|^2 = E_0^T E_0$$

TOTAL SUM OF SQUARES

$$\cdot R_{\text{adj}}^2 = 1 - \frac{n-p}{n-p} \cdot \frac{RSS}{TSS}$$

K: # ESTIMATED PARAMETERS IN THE MODEL  
L: MAXIMUM VALUE OF LIKELIHOOD FUNCTION

$$\downarrow \text{AIC, BIC} \rightarrow \text{BETTER MODEL}$$

GLM: `my_fit <- glm(y~predictors, family=binomial)`

$$\cdot \text{DOOS: } \frac{P}{1-P} / [0, 1] \mapsto [0, +\infty]$$

Family	Default link function	Description
binomial	(link = "logit") <span style="color: red;">↳ LOGISTIC REGRESSION</span>	
gaussian	(link = "identity") <span style="color: red;">↳ LINEAR MODELS</span>	
gamma	(link = "inverse")	
inverse.gaussian	(link = "1/mu^2")	
poisson	(link = "log") <span style="color: red;">↳ for count variables, like number of accidents</span>	
quasi	(link = "identity", variance = "constant")	
quasibinomial	(link = "logit")	
quasipoisson	(link = "log")	

$$\cdot \text{LOGIT (LOG-ODDS): } \log\left(\frac{P}{1-P}\right) / [0, 1] \mapsto [-\infty, +\infty]$$

$$\cdot \text{GLM: } \begin{pmatrix} \text{logit}(P_1) \\ \vdots \\ \text{logit}(P_n) \end{pmatrix} = X \cdot \bar{\beta}$$

$E(P) = E(X\bar{\beta}) = \begin{pmatrix} P_1 \\ \vdots \\ P_n \end{pmatrix}$   
 $\hookrightarrow P_i \sim \text{Bernoulli}(P_i)$

$$\begin{cases} \text{logit}(P) \mapsto P: & P = \frac{e^P}{1+e^P} \rightarrow P_i = \frac{e^{X \cdot \bar{\beta}}}{1+e^{X \cdot \bar{\beta}}} \\ [-\infty, +\infty] \mapsto [0, 1] & \end{cases}$$

$$\begin{cases} \text{logit}(P) \mapsto P: & P = \frac{e^P}{1+e^P} \rightarrow P_i = \frac{e^{X \cdot \bar{\beta}}}{1+e^{X \cdot \bar{\beta}}} \\ [-\infty, +\infty] \mapsto [0, 1] & \end{cases}$$

```
# Perform the F-test
f_test <- anova(reduced_model, full_model)
```

Analysis of Variance Table				
	Model 1: Sales ~ TV + Radio	Model 2: Sales ~ TV + Radio + Newspaper	Res.Df	RSS Df Sum of Sq F Pr(>F)
1	197	541.21	1	0.0092861
2	196	541.20	1	0.0034 0.9538

## CONFUSION MATRIX:

PREDICTED		1	0
TRUE	1	TP	FN
	0	FP	TN

- SENSITIVITY = RECALL( $c=1$ ) =  $P(\hat{Y}=1 | Y=1)$   
 $\hookrightarrow \text{SENSITIVITY} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \text{TPR}$

- SPECIFICITY = RECALL( $c=0$ ) =  $P(\hat{Y}=0 | Y=0)$   
 $\hookrightarrow \text{SPECIFICITY} = \frac{\text{TN}}{\text{TN} + \text{FP}} = 1 - \text{FPR}$

$$\int_{X|Y} f(x|y) (x_i | Y=y) = \frac{\int f(x_i, y)}{\int_Y f(y)} \quad \begin{cases} f(x_i, y) dx \\ \uparrow \text{MARGINAL DENSITY OF } Y \end{cases}$$

- DA(GS):

$$f(x_1, \dots, x_d) = \prod_i f_{x_i | \text{parents}(x_i)}(x_i | \text{parents}(x_i))$$

- MARKOVIAN CHAINS:

$$(X_1) \rightarrow (X_2) \rightarrow \dots \rightarrow (X_d)$$

$$f(x_1, \dots, x_d) = f_{x_1}(x_1) \cdot f_{x_2|x_1}(x_2|x_1) \cdot f_{x_3|x_2}(x_3|x_2) \cdot \dots \cdot f_{x_d|x_{d-1}}(x_d|x_{d-1})$$

## CONDITIONAL EXPECTATION:

$$E[Y|X=x] = \mu_y + \rho \frac{\sigma_y}{\sigma_x} (x - \mu_x) \quad / \rho = \frac{Cov(X, Y)}{V(X)}$$

$$V[Y|X=x] = V(Y) - \frac{Cov^2(X, Y)}{V(X)}$$

$$1. E[E[Y|X]] = E[Y]$$

$$2. E[a|X] = a \quad / a = \text{const.}$$

$$3. E[aX + bY|Z] = aE(X|Z) + bE(Y|Z)$$

$$4. E[g(X) \cdot Y|X] = g(X) \cdot E(Y|X)$$

## BAYESIAN STATISTICS:

	p.d.f.	$E(x)$	$\text{Var}(x)$
$x \sim \text{Beta}(\alpha, \beta)$	$\frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}$	$\frac{\alpha}{\alpha+\beta}$	$\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$
$x \sim \text{Gamma}(\alpha, \beta)$	$\frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp(-\beta x)$	$\frac{\alpha}{\beta}$	$\frac{\alpha}{\beta^2}$

### Examples

(1)  $\theta \sim \text{Beta}(\alpha, \beta)$  is the conjugate prior to  $X|\theta \sim \text{Bin}(n, \theta)$  with kernel:

$$g(\theta) = \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

(2)  $\theta \sim N(\mu_0, \sigma_0^2)$  is the conjugate prior to  $X|\theta \sim N(\theta, \sigma^2)$  with kernel:

$$g(\theta) = \exp\left(-\frac{(\theta - \mu_0)^2}{2\sigma_0^2}\right)$$

(3)  $\sigma^2 \sim \text{InvGamma}(\alpha, \beta)$  is the conjugate prior to  $X|\sigma^2 \sim N(\theta, \sigma^2)$  with kernel:

$$g(\sigma^2) = (\sigma^2)^{-(\alpha+1)} \exp\left(-\frac{\beta}{\sigma^2}\right)$$

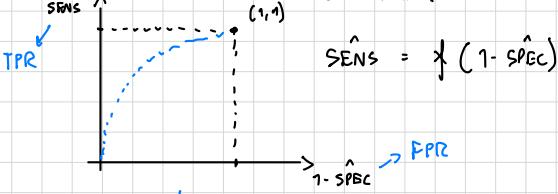
(4)  $\lambda \sim \text{Gamma}(\alpha, \beta)$  is the conjugate prior to  $X|\lambda \sim \text{Poi}(\lambda)$  with kernel:

$$g(\lambda) = \lambda^{\alpha-1} \exp(-\beta\lambda)$$

- IF WE VARY  $t \rightarrow$  CHANGE OF SENSITIVITY

AND SPECI<sub>I</sub>TY

ROC CURVE:



$$\int_{X|Y=y} f(x|y) (x_i | Y=y) = \frac{\int_X f_X(x) \cdot f_{Y|X=x}(y|x) dx}{\int_X f_X(x) dx}$$

$$X \rightarrow Y \quad \begin{matrix} w \\ z \end{matrix} \quad \begin{matrix} Y = \text{parents}(w, z) \\ Z = \text{children}(Y) \end{matrix}$$

$$\rightarrow f_{(X, Y, Z)} = f_Z(z|Y) \cdot f_{(W|Y)}(w|Y) \cdot f_{Y|X}(Y|X) \cdot f_X(X)$$

- INDEPENDENT RVs:

$$(X_1) \quad (X_2) \quad (X_d) : \quad f_{(X_1, \dots, X_d)} = \prod_i f_{X_i}(x_i)$$

$$- \text{Cov}(aX + b, Y) = a \cdot \text{Cov}(X, Y)$$

$$\cdot \text{Cov}(aX + bY, cV + dW) =$$

$$= ac \text{Cov}(X, V) + ad \text{Cov}(X, W) + bc \text{Cov}(Y, V) + bd \text{Cov}(Y, W)$$

$$V(Y|X) = E[(Y - E(Y|X))^2 | X]$$

$$5. E[V(Y|X)] = V(Y) - V(E(Y|X))$$

$$\hookrightarrow V(Y) = E[V(Y|X)] + V(E(Y|X))$$

$$\text{POSTERIOR} \quad \text{PRIOR}$$

$$\tilde{f}(\theta | \text{data}) = \frac{\tilde{f}(\text{data} | \theta) \cdot \tilde{f}(\theta)}{\int \tilde{f}(\theta) \cdot \tilde{f}(\text{data} | \theta) d\theta} \propto \tilde{f}(\text{data} | \theta) \cdot \tilde{f}(\theta)$$

$$\begin{cases} \tilde{\gamma}_1 = \frac{1}{G_1} & : \text{PRECISION} \\ \tilde{\gamma}_0 = \frac{1}{G_0} & : \text{PRIOR PRECISION} \end{cases}$$

$$\cdot B(a, b) = \int_0^\infty t^{a-1} (1-t)^{b-1} dt = \frac{T(a) T(b)}{T(a+b)} \quad / \quad T(x) = \int_0^\infty t^{x-1} e^{-t} dt = (x-1)! T(x-1)$$

1. Let  $X|\theta \sim \text{Bin}(n, \theta)$  with conjugate prior  $\theta \sim \text{Beta}(\alpha, \beta)$ , then

$$\theta|X \propto \text{Beta}(\alpha + X, \beta + n - X) : E(\pi(\theta|X)) = \frac{\alpha + X}{\alpha + \beta + n}$$

2. Let  $X|\theta \sim N(\theta, \sigma^2)$  with conjugate prior  $\theta \sim N(\mu_0, \sigma_0^2)$ , then

$$\theta|X \propto \text{Normal}(\mu_0 + \frac{\sigma_0^2}{\sigma^2} X, \frac{\sigma_0^2}{\sigma^2}) : E(\pi(\theta|X)) = \left(\frac{1}{\sigma^2} + \frac{1}{\sigma_0^2}\right)^{-1} \left(\frac{X}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right)$$

3. Let  $X|\sigma^2 \sim N(\theta, \sigma^2)$  with conjugate prior  $\sigma^2 \sim \text{InvGamma}(\alpha, \beta)$ , then

$$\sigma^2|X \propto \text{InvGamma}(\alpha + 1/2, \beta + (X - \mu)^2/2) : E(\pi(\sigma^2|X)) = \frac{\beta + (X - \mu)^2/2}{\alpha - 1/2}$$

4. Let  $X|\lambda \sim \text{Poi}(\lambda)$  with conjugate prior  $\lambda \sim \text{Gamma}(\alpha, \beta)$ , then

$$\lambda|X \propto \text{Gamma}(X + \alpha, 1 + \beta) : E(\pi(\lambda|X)) = \frac{X + \alpha}{1 + \beta}$$

$$\cdot X_1, \dots, X_m \mid p \stackrel{\text{cond}}{\sim} \text{Bernoulli}(p) \quad / \quad p \sim \text{Beta}(a, b)$$

BAYES POSTERIOR:

CONDITIONAL DENSITY:

$$\begin{aligned} f(x_1, \dots, x_m \mid p) &= \frac{1}{\prod_{i=1}^m \int_{x_i} p (x_i \mid p)} = \\ \hookrightarrow L(p, x_1, \dots, x_m) &= p^{\sum x_i} \cdot (1-p)^{n-\sum x_i} \end{aligned}$$

↳ LIKELIHOOD FUNCTION

$$\begin{aligned} \prod(p \mid x_1, \dots, x_m) &= \frac{\prod_l L(p) \cdot f(x_1, \dots, x_m \mid p)}{\int_0^1 \prod_l L(q) \cdot f(x_1, \dots, x_m \mid q) dq} \stackrel{\text{const.}}{=} \frac{\prod_l \int_0^1 t^{a-1} (1-t)^{b-1} dt \cdot p^{\sum x_i} \cdot (1-p)^{n-\sum x_i}}{\int_0^1 \prod_l \int_0^1 t^{a-1} (1-t)^{b-1} dt \cdot f(x_1, \dots, x_m \mid q) dq} \stackrel{\text{const.}}{=} \\ &= p^{a+n} (1-p)^{b+n} \cdot p^{\sum x_i} \cdot (1-p)^{n-\sum x_i} = p^{a+\sum x_i} (1-p)^{b+n-\sum x_i} \end{aligned}$$

$$\rightarrow \prod(p \mid x_1, \dots, x_m) = \frac{p^{a+\sum x_i} (1-p)^{b+n-\sum x_i}}{\int_0^1 t^{a+\sum x_i} (1-t)^{b+n-\sum x_i} dt} \rightarrow p \mid x_1, \dots, x_m \sim \text{Beta}(a + \sum x_i, b + n - \sum x_i)$$

BAYESIAN POINT ESTIMATION:

NORMAL RANDOM SAMPLE:

$$\begin{aligned} E(\mu \mid x_1, \dots, x_m) &= \frac{\tilde{\tau}_0 M_0 + n \tilde{x}}{\tilde{\tau}_0 + n \tilde{\tau}} = \frac{\tilde{\tau}_0}{\tilde{\tau}_0 + n \tilde{\tau}} M_0 + \frac{n \tilde{\tau}}{\tilde{\tau}_0 + n \tilde{\tau}} \bar{x} = \\ \text{WEIGHTED AVERAGE OF: } \cdot \tilde{\mu}_0 &: \text{prior mean} \quad \text{OR} = \left( 1 - \frac{n \tilde{\tau}}{\tilde{\tau}_0 + n \tilde{\tau}} \right) M_0 + \frac{n \tilde{\tau}}{\tilde{\tau}_0 + n \tilde{\tau}} \bar{x} \quad / \quad W = \frac{n \tilde{\tau}}{\tilde{\tau}_0 + n \tilde{\tau}} \end{aligned}$$

• ↓  $\tilde{\tau}_0$  ( $\rightarrow$  HIGHER  $\tilde{\sigma}_0^2$ : HIGHER UNCERTAINTY)  $\rightarrow$  ↑ DATA IMPORTANCE FOR ESTIMATION

BINARY DATA CASE, Beta(a, b):

$$E(p \mid x_1, \dots, x_m) = \frac{a + \sum x_i}{a + \sum x_i + b + n - \sum x_i} = \frac{a + \sum x_i}{a + b + n} = \frac{a+b}{a+b+n} \cdot \frac{a}{a+b} + \frac{m}{a+b+n} \cdot \frac{\sum x_i}{m} =$$

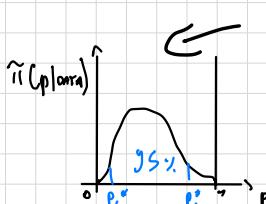
$$\text{SOME PROPERTIES AS ABOVE} \quad \text{OR} = \left( 1 - \frac{m}{a+b+n} \right) \frac{a}{a+b} + \frac{m}{a+b+n} \bar{x} \quad \hookrightarrow \hat{p}: \text{MLE OF } p$$

$$\cdot \text{IF } a=b=1: E(p \mid x_1, \dots, x_m) = \frac{1 + \sum x_i}{2 + n} = \hat{p} \quad \text{BAYES ESTIMATION}$$

↳ IT USES UNIFORM ON P: GOOD PRIOR IGNORANCE

$$\begin{cases} \frac{P_L^3 (1-P_L)^3}{\int t^3 (1-t)^3 dt} = \frac{P_0^3 (1-P_0)^3}{\int t^3 (1-t)^3 dt} \\ \int_0^1 \frac{P_L^3 (1-P)^3}{\int t^3 (1-t)^3 dt} dp = 0.95 \end{cases}$$

$$\rightarrow \begin{cases} \text{DISTRIBUTION OF A BETA} \\ \text{OR} \text{Beta}(P_L, 2, 8) = \text{Beta}(P_0, 4, 8) \\ P \text{Beta}(P_L, 2, 8) - P \text{Beta}(P_L, 4, 8) = 0.95 \end{cases}$$



CREDIBILITY INTERVAL

ML

- IF  $A$  IS SYMMETRIC  $\left\{ \begin{array}{l} \rightarrow \lambda \in \mathbb{R}, \text{ EIGENVECTORS ORTHONORMAL} \\ \rightarrow x^T A x \geq 0 \end{array} \right. \quad (\text{SPECTRAL THEOREM})$

## MATRIX PROPERTIES:

$$\cdot AB \neq BA \quad | \quad \cdot A(B+C) = AB + AC$$

$$\begin{aligned} \cdot (A, \dots, A_m)^T &= A_m^T \dots A_1^T \rightarrow \begin{cases} (AB)^T = B^T A^T \\ (A^T B^T C^T) = C^T B^T A^T \end{cases} \\ \cdot (A+B)^T &= A^T + B^T \end{aligned}$$

$$\cdot A x = B \rightarrow \text{IF } \exists A^{-1}: x = A^{-1} B$$

$$\begin{matrix} \text{ORTHOGONAL} \\ \text{MATRIX} \end{matrix} \uparrow \quad \begin{matrix} \text{DIAGONAL} \\ \text{MATRIX} \end{matrix} \uparrow \\ A = Q \begin{bmatrix} \Lambda & 0 \\ 0 & C \end{bmatrix} Q^T \end{math>$$

• ORTHONORMAL BASIS:

$$\{ \bar{v}_1, \dots, \bar{v}_n \} / \langle \bar{v}_i, \bar{v}_j \rangle = 0$$

• PAC LEARNING; VC DIM = d  $\leq \infty \rightarrow$  CAN'T LEARN FROM  $\frac{m}{2}$  SAMPLES  
VC DIM =  $\infty \rightarrow$  NOT PAC LEARNABLE

$\begin{cases} S \rightarrow P \text{ OF SUCCESS} \\ 1-S \rightarrow P \text{ OF FAILURE} \end{cases}$

$$1-S \leq \epsilon$$

PROBABLY APPROXIMATOR CORRECT LEARNING

• THE LEARNER DO NOT KNOW D AND  $\epsilon$

• THE LEARNER SHOULD OUTPUT AN HYPOTHESIS  $h$  WITH

PROBABILITY OF AT LEAST  $1-\delta$  / IT HOLDS THAT  $L_{D,\epsilon}(h) \leq \epsilon \rightarrow P_{S \sim D^m} [L_{D,\epsilon}(h) \leq \epsilon] \geq 1-\delta$

• THEOREM (NO FREE LUNCH):

$$\text{Fix } S \in (0,1), \epsilon < 1/2 :$$

$\rightarrow \forall$  LEARNER A AND TRAINING SIZE m:

$\rightarrow \exists D, \epsilon / \text{with } p \geq S$

OF m EXAMPLES, OVER TRAINING DATA S

IT HOLDS THAT  $L_{D,\epsilon}(A(S)) \geq \epsilon$

$\begin{cases} \cdot \forall \text{ LEARNER} \rightarrow \exists \text{ A FUNCTION WHICH DO NOT SUCCEED IN } L_{D,\epsilon} \leq \epsilon \\ \cdot \nexists \text{ A UNIQUE LEARNER: EVERY TIME A SPECIFIC ONE IS NEED} \\ \cdot \text{THE AVG RESULT OF ALL CLASSIFIER IS WORST THAN RANDOM GUESSING} \\ \cdot L_{D,\epsilon}(\text{RANDOM GUESING}) = 1/2 \\ \cdot \text{IF NO KNOWLEDGE IS INSERTED} \rightarrow \nexists \text{ SOMETHING BETTER THAN RANDOM GUESSING} \end{cases}$

• PRIOR KNOWLEDGE: IF LEARNER KNOWS  $H \subset \mathcal{Y}^X$   $\rightarrow$  PAC LEARNING IS POSSIBLE IF  $H \neq \mathcal{Y}^X$

• PAC LEARNING OF FINITE HYPOTHESIS CLASS: ASSUME H: FINITE HYPOTHESIS CLASS (so. acc  $h: X \mapsto \mathcal{Y}$ )

• EMPIRICAL RISK MINIMIZATION (ERM):  $\rightarrow h = \text{ERM}_H(S)$

ERM

LEARNING

RULE

$\begin{cases} \cdot \text{INPUT: TRAINING SET } S = (x_1, y_1), \dots, (x_m, y_m) \\ \cdot \text{EMPIRICAL RISK: } L_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h(x_i) \neq y_i\}} \\ \cdot \text{OUTPUT: ANY } h \in H / h \text{ MINIMIZES } L_S(h) \end{cases} \quad \begin{cases} \text{THEOREM:} \\ \text{Fix } \epsilon, S, \text{ IF } m \geq \frac{\log(|H|/\epsilon)}{\epsilon} \\ \text{OVER THE CHOICE OF S SIZE } m \\ \rightarrow \forall D, \epsilon, \text{ with } p \geq 1-\delta: L_{D,\epsilon}(\text{ERM}_H(S)) \leq \epsilon \end{cases}$

PAC LEARNABILITY:

AN HYPOTHESIS CLASS H IS PAC LEARNABLE, IF  $\exists$  FUNCTION  $m_H: (0,1)^2 \mapsto \mathbb{N}$  AND A LEARNING ALGORITHM /:

$\cdot \forall \epsilon, \delta \in (0,1) \quad | \quad \cdot \forall D \text{ OVER } X, \text{ AND } \forall \text{ LABELING FUNCTION } g: X \mapsto \{0,1\}$

WHEN RUNNING THE ALGORITHM ON  $m \geq m_H(\epsilon, \delta)$  iid EXAMPLES

GENERATED BY D AND LABELED BY  $g$  SAMPLE COMPLEXITY:  $\min_{H \in \mathcal{H}} \# \text{ INPUT DATA TO ASSURE }$

$\rightarrow$  ALGORITHM RETURNS  $h / \text{with } p \geq 1-\delta: L_{D,g}(h) \leq \epsilon$

SAMPLE COMPLEXITY FOR  $\epsilon, \delta$ :

$m_H(\epsilon, \delta) \geq \frac{\log(|H|/\delta)}{\epsilon}$

COROLLARY:

LET H BE A FINITE HYPOTHESIS CLASS:

$\rightarrow H$  IS PAC LEARNABLE WITH

$m_H(\epsilon, \delta) \leq \frac{\log(|H|/\delta)}{\epsilon} / \text{ OBTAINED USING ERM}_H \text{ LEARNING RULE}$

- VC DIMENSION: WE SAY THAT  $H$  "SHATTERS"  $C$  IF  $|H_C| = 2^{|C|}$
- $|H_C| \leq 2^{|C|}$  ↳ IF A SET OF FUNCTIONS  $\in H_C$  IS ABLE TO PRODUCE ALL COMBINATIONS OF  $S \subseteq \{1\}^{|C|}$  OUTPUTS →  $H$  shatters  $C$
- $\text{VC dim}(H) = d = \sup S |C|$ :  $H$  shatters  $\{C\}$  IT IS THE HIGHEST CARDINALITY OF A SET OF  $C/H$  SHATTERS THE SET  $C$
- IN ORDER TO SHOW THAT  $\text{VC dim}(H) = d$ :  $\text{VC dim} \leq \text{d} \quad \text{VC dim} \geq d$ 
  - JUST 1 SET IS NEEDED
  - 1. show that  $\exists$  A SET  $C$  OF SIZE  $|C| = d$  /  $H$  SHATTERS  $\rightarrow \text{VC dim}(H) \geq d$
  - 2. A SET  $C / |C| = d+1$ ,  $C$  IS NOT SHATTERED BY  $H \rightarrow \text{VC dim}(H) < d+1$

$H$  FINITE → PAC LEARNABLE  
 $H$  INFINITE →  $\exists$  PAC LEARNABLE  
 (NEED OF VC DIM <  $\infty$ )

### • THE FUNDAMENTAL THEOREM OF PAC LEARNING: TFAE: → THE FOLLOWING ARE EQUIVALENT

1.  $H$  HAS THE UNIFORM CONVERGENCE PROPERTY
2. ANY ERM RULE IS A SUCCESSFUL AGNOSTIC PAC LEARNER FOR  $H$
3.  $H$  IS AGNOSTIC PAC LEARNABLE
4.  $H$  IS PAC LEARNABLE
5. ANY ERM RULE IS A SUCCESSFUL PAC LEARNER FOR  $H$
6.  $H$  HAS A FINITE VC-DIM =  $d < \infty \rightarrow \text{VC dim}(H) = d \rightarrow H$  IS PAC LEARNABLE

### • THEOREM:

LET:  $H$  BE A FINITE HYPOTHESIS CLASS OF FUNCTIONS  $X \mapsto \{0, 1\}$  • LOSS FUNCTION: 0-1 LOSS  $\rightarrow \text{VC dim} = d < \infty$

$$\rightarrow \exists C_1, C_2 \text{ CONSTANTS, SUCH THAT: } C_1 \cdot \frac{d + \log(\frac{1}{\delta})}{\epsilon} \geq m_H(\epsilon, \delta) \leq C_2 \cdot \frac{d \log(\frac{1}{\epsilon}) + \log(\frac{1}{\delta})}{\epsilon}$$

AND  $m_H(\epsilon, \delta)$  IS ACHIEVED BY THE ERM LEARNING RULE

LEMMA: LET  $H$  BE AN HYPOTHESIS CLASS /  $\text{VC dim}(H) = d < \infty \rightarrow \forall C \subset X / |C| = m > d+1 \rightarrow |H_C| \leq \left(\frac{e^m}{d}\right)^d \approx K \cdot m^d$

AGNOSTIC PAC LEARNABILITY		PAC	Agnostic PAC
D OVER $X \times Y$	$\begin{cases} \text{0-1 loss: } l(h, (x, y)) = 1, \text{ IF } h(x) \neq y \\ \text{squared loss: } l(h, (x, y)) = (h(x) - y)^2 \end{cases}$	D over $X$	$D$ over $X \times Y$
RISK: $L_D(h) = P_{(x,y) \sim D} [h(x) \neq y] = D(\{(x, y) : h(x) \neq y\})$	SOME LOSSES	Truth	not in class or doesn't exist
"APPROXIMATELY CORRECT" NOTION: $L_D(A(S)) \leq \min_{h \in H} L_D(h) + \epsilon$	MINIMUM LOSS AMONG ALL FUNCTIONS	Training set	$L_{D,f}(h) = D(\{(x : h(x) \neq f(x)\})$ $(x_1, \dots, x_m) \sim D^m$ $\forall i, y_i = f(x_i)$
FORMAL DEF:		Goal	$L_{D,f}(A(S)) \leq \epsilon$ $L_D(A(S)) \leq \min_{h \in H} L_D(h) + \epsilon$

{ AN HYPOTHESIS CLASS  $H$  IS AGNOSTIC PAC LEARNABLE TO A SET  $Z = X \times Y$  AND A LOSS FUNCTION  $l: H \times Z \rightarrow \mathbb{R}^+$ , IF  $\exists$  A FUNCTION  $m_H: (0, 1)^2 \rightarrow \mathbb{N}$  AND A LEARNING ALGORITHM  $A$  WITH THE PROPERTY:

$$\forall \epsilon, \delta \in (0, 1), m \geq m_H(\epsilon, \delta), \text{ AND A DISTRIBUTION } D \text{ OVER } Z. D^m \left( \left\{ \sum_{z \in Z} z^m : L_D(A(S)) \leq \min_{h \in H} L_D(h) + \epsilon \right\} \right) \geq 1 - \delta$$

• A TRAINING SET  $S$  IS  $\epsilon$ -REPRESENTATIVE SAMPLE IF  $\forall h \in H, |L_S(h) - L_D(h)| \leq \epsilon$

$$\text{THEORETICAL TRUE DISTRIBUTION OF } Z \rightarrow \text{DISTRIBUTION } D^m \left( \left\{ \sum_{z \in Z} z^m : L_D(A(S)) \leq \min_{h \in H} L_D(h) + \epsilon \right\} \right) \geq 1 - \delta$$

CARDINALITY / PRODUCT / EACH  $Z \sim D$  ARE POSSIBLE  $m$ -TUPLES IN  $Z$

$\rightarrow$  EXPECTED LOSS IS AT MOST  $\epsilon$ -LARGER THAN THE EXPECTED LOSS OF THE BEST  $h$

• UNIFORM CONVERGENCE:  $\rightarrow$  IF UC  $\rightarrow H$  CAN LEARN WELL FROM  $S$

$H$  HAS THE UNIF. CONV. PROPERTY IF  $\exists$  FUNCTION  $m_H^{UC}: (0, 1)^2 \rightarrow \mathbb{N} / \forall \epsilon, \delta \in (0, 1)$  AND  $\forall d$ :

$$L_D \left( \left\{ S \in Z^m : S \text{ IS } \epsilon\text{-REPRESENTATIVE} \right\} \right) \geq 1 - \delta$$

$$\rightarrow \exists m_H^{UC} / D^m \left( \left\{ S \in Z^m : S \text{ IS } \epsilon\text{-REPRESENTATIVE} \right\} \right) \geq 1 - \delta \rightarrow UC$$

• COROLLARY:  $\rightarrow$  IF UC  $\rightarrow \exists$  BOUND (UNIF.) FOR  $m_H(\epsilon, \delta)$ :

IF  $H$  HAS UC PROPERTY WITH A FUNCTION  $m_H^{UC}$ :

• IN THAT CASE, ERM<sub>H</sub> ALGORITHM IS A SUCCESSFUL AGNOSTIC PAC LEARNER FOR  $H$

• OUTPUT OF ERM<sub>H</sub>(S)  $\rightarrow h_S = \text{ERM}_H(S)$

•  $\rightarrow$  IF  $S$  IS  $\epsilon$ -REPRESENTATIVE  $\rightarrow L_D(h_S) \leq \min_{h \in H} L_D(h) + \epsilon \rightarrow H$  IS AGNOSTICALLY PAC LEARNABLE WITH  $m_H(\epsilon, \delta) \leq m_H^{UC} \left( \frac{\epsilon}{2}, \delta \right)$

THEOREM (FINITE CLASS  $\rightarrow$  AGNOSTIC PAC LEARNABLE):

LET  $H$  BE FINITE AND  $\ell \in [0, 1]$

$\rightarrow H$  IS AGNOSTICALLY PAC LEARNABLE WITH:  $m_H(\epsilon, \delta) \leq \lceil \frac{2 \log(2|H|/\delta)}{\epsilon^2} \rceil$   $\Rightarrow \lceil \cdot \rceil : \text{CEILING FUNCTION}$   $\lceil 2.4 \rceil = 3$

DISCRETIZATION TRICK: suppose  $H$  PARAMETERIZED BY  $d$  NUMBERS / EACH NUMBER REPRESENTED BY  $b$  BITS

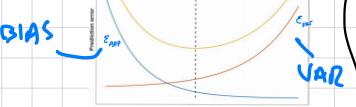
$$\rightarrow |H| \leq 2^{d \cdot b}$$

$$\rightarrow m_H(\epsilon, \delta) \leq \lceil \frac{2db + 2 \log(\frac{2}{\delta})}{\epsilon^2} \rceil$$

LEAST SQUARES:  $\left\{ \begin{array}{l} \rightarrow \text{PROBLEM: } \min_{\bar{w} \in \mathbb{R}^d} \frac{1}{2} \| \bar{w}^T \bar{w} - \bar{v} \|^2 \\ \rightarrow \text{SOL: } \bar{w} = (\bar{x} \bar{x}^T)^{-1} \bar{x} \bar{v} \end{array} \right.$

$$\left. \begin{array}{l} J_x(f \circ y) = f(y(x)) \\ = J_{y(x)}(f) \cdot J_x(y) \\ \frac{df}{dy} \cdot \frac{dy}{dx} \end{array} \right\}$$

POLYNOMIAL FITTING:  $\left\{ \begin{array}{l} \text{DEFINE } \psi: \mathbb{R} \mapsto \mathbb{R}^{m+1}, \quad \psi(x) = (1, x, x^2, \dots, x^m) \\ \text{DEFINE } \bar{a} = (a_0, a_1, \dots, a_m) \text{ AND OBSERVE: } p(x) = \sum_{i=0}^m a_i x^i = \langle \bar{a}, \psi(x) \rangle \rightarrow \text{HYP. SPACE/LIN. REGRESSION} \\ \rightarrow \text{TO FIND } \bar{a} \text{ WE CAN SOLVE LS W.R.T. } ((\psi(x_1), v_1), \dots, (\psi(x_m), v_m)) \end{array} \right.$



$$L_D(h_s) = \epsilon_{\text{app}} + \epsilon_{\text{est}}$$

↓  
CHANGING POLYNOMIAL DEGREE  
↑ n → ↓ \epsilon\_{\text{app}}

- APPX. ERROR  $\rightarrow \epsilon_{\text{app}} = \min_{h \in H} L_D(h)$
- $\rightarrow$  IT IS THE RISK WE TAKE DUE TO RESTRICTING TO  $H$
- $\rightarrow$  IT DOES NOT DEPEND ON  $S$
- $\rightarrow \uparrow$  (COMPLEXITY OF  $H$  / VC DIM)  $\rightarrow \downarrow \epsilon_{\text{app}}$

$$\cdot \text{ESTIMATION ERROR} \rightarrow \epsilon_{\text{est}} = L_D(h_s) - \epsilon_{\text{app}}$$

- $\rightarrow$  IT TELLS HOW GOOD AM I IN EXPLOITING  $H$  ON MY PROBLEM
  - $\rightarrow$  IT IS THE RESULT OF  $L_s$  BEING ONLY AN ESTIMATE OF  $L_D$
  - $\rightarrow \uparrow$  SIZE OF  $S \rightarrow \downarrow \epsilon_{\text{est}}$
  - $\rightarrow \uparrow$  COMPLEXITY OF  $H \rightarrow \uparrow \epsilon_{\text{est}}$  :  $\uparrow$  COMPLEXITY  $\rightarrow \downarrow \epsilon_{\text{app}} \rightarrow \uparrow \epsilon_{\text{est}}$
- IT CANNOT HIT ALL OTHER DATA NOT IN TRAINING SETS

VALIDATION: HOEFDING'S INEQUALITY, IF  $\ell \in [0, 1]$ :  $|L_V(h) - L_D(h)| \leq \sqrt{\frac{\log(\frac{2}{\delta})}{2 \cdot m_V}}$

$\cdot$  LET  $V$  BE A NEW VAL. SET, CHOOSE  $h^* \in \text{ERM}_H(V) \rightarrow L_D(h^*) \leq \min_{h \in H} L_D(h) + \sqrt{\frac{2 \log(\frac{2|H|}{\delta})}{|V|}}$

SVM: HALF-SPACE CLASSIFIER:  $x \mapsto \text{sign}(\langle \bar{w}, \bar{x} \rangle + b) = \pm 1, \bar{w} \in \mathbb{R}^d$

$$L = \{V: \langle \bar{w}, v \rangle + b = 0\}, \text{ol}(\bar{x}, L) = \min \{||\bar{x} - \bar{v}||: \bar{v} \in L\}$$

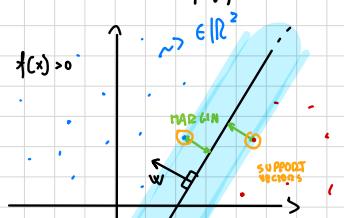
MARGIN:  $\min |\langle \bar{w}, \bar{x}_i \rangle + b| \rightarrow$  CLOSEST DISTANCE FROM  $L$  TO  $x_i, V_i$

HARD SVM:

$$\text{FIND } \bar{w}, b \text{ TO } \min_{\bar{w}, b} \left( \min_{i \in [m]} |\langle \bar{w}, \bar{x}_i \rangle + b| \right) / \forall i, V_i: \langle \bar{w}, \bar{x}_i \rangle + b \geq 0$$

$$\text{OR EQUIVALENTLY: } \max_{(\bar{w}, b): ||\bar{w}||=1} \left( \min_{i \in [m]} V_i (\langle \bar{w}, \bar{x}_i \rangle + b) \right) \text{ BEST PLANE: } \langle \bar{w}, \bar{x} \rangle + b = 0$$

$$\text{OR EQUIVALENTLY: } (\bar{w}_0, b_0) = \arg \min_{\bar{w}, b} (\|\bar{w}\|^2) / \forall i, V_i (\langle \bar{w}, \bar{x}_i \rangle + b) \geq 1$$



$$\text{MARGIN} = \frac{1}{\|\bar{w}\|}$$

2 CLASSES SEPARATION

$$\text{POINT } (x_p, y_p, z_p) \text{ PLANE } \text{ol}(P, l) = \frac{ax_p + by_p + cz_p + d}{\sqrt{a^2 + b^2 + c^2}}$$

MARGIN OF DISTRIBUTION:

$D$  IS SEPARABLE WITH A  $(\gamma, \delta)$ -MARGIN, IF  $\exists (\bar{w}^*, b^*) / \|\bar{w}^*\| = 1$

AND:  $D(\{(x, y): \|\bar{x}\| \leq \gamma\} \cap \{(\bar{x}, \bar{y}): \langle \bar{w}^*, \bar{x} \rangle + b^* \geq 1\}) = 1$

THEOREM:

$$\rightarrow \text{IF } D \text{ SEPARABLE WITH } (\gamma, \delta)-\text{MARGIN} \rightarrow m(\epsilon, \delta) = \frac{\delta}{\epsilon^2} \left( 2 \left( \frac{\gamma}{\delta} \right)^2 + \log \left( \frac{2}{\delta} \right) \right)$$

$\rightarrow$  UNLIKE VC BOUNDS,  $m(\cdot)$   $\frac{\delta}{\epsilon}$ , NOT ON  $d$

$\text{SOFT SUM} : \underset{\bar{w}, b, \xi}{\operatorname{argmin}} \left( \lambda \cdot \|\bar{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right) / \forall i, \xi_i (\langle \bar{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i, \xi_i > 0$ 
  
 $\downarrow \text{COST} \rightarrow \uparrow \text{ACC TEST}$   
 $\downarrow \# \text{FAVORS} \rightarrow \uparrow \text{ACC TEST}$

$\text{COST} : = \underset{\bar{w}, b}{\operatorname{argmin}} \left( \lambda \cdot \|\bar{w}\|^2 + L_s(\langle \bar{w}, b \rangle) \right) \quad | \quad \text{HINGE LOSS } l^{\text{Hinge}}((\bar{w}, b), (\bar{x}, y)) = \max \{ 0, 1 - \langle \bar{w}, \bar{x} \rangle + b \}$ 
  
 $C \mapsto 0 : \text{POINTS IN margin}$ 
  
 $C \mapsto \infty : \mapsto \text{HARD SUM}, \xi_i \mapsto 0 \rightarrow 1 - \xi_i \mapsto 1$

$\text{ON-AVERAGE-REPLACE-ONE-STABLE DEFINITION} : S^{(i)} = (z_1, \dots, z_{i-1}, \bar{z}', z_{i+1}, \dots, z_m) / S = (z_1, \dots, z_m)$

LET  $\varepsilon : \mathbb{N} \mapsto \mathbb{R}$  BE A MONOTONICALLY DECREASING FUNCTION.

ALGORITHM A IS ON-AVERAGE-REPLACE-ONE-STABLE WITH RATE  $\varepsilon(m)$  IF DISTRIBUTION D:

$$\underset{(S, z) \sim D}{\mathbb{E}} \left[ l(A(S^{(i)}), z_i) \right] - l(A(S), z_i) \leq \varepsilon(m)$$

• FRIE-JOHNS;  $\bar{w}_0 = \sum_i \alpha_i \bar{x}_i$   
 L > OPTIMAL PLAIN WEIGHTS ARE A LIN. COMB. OF SUPPORT VECTORS

THEOREM: IF A IS ON-AVERAGE-REPLACE-ONE-STABLE WITH RATE  $\varepsilon(m)$ :  $\rightarrow \underset{S \sim D^m}{\mathbb{E}} [L_0(A(S)) - L_s(A(S))] \leq \varepsilon(m)$

THEOREM:  $|f(x_0) - f(x_2)| \rightarrow \text{RLM RULE, WITH REGULARIZER } \lambda \|\bar{w}\|^2$ , IS ON-AVERAGE-REPLACE-ONE-STABLE TIKHONOV REGULARIZATION AS STABILIZED:  $\leq \frac{2s^2}{\lambda m}$  WITH RATE  $\varepsilon(m) = \frac{2s^2}{\lambda m}$ , AND:  $\text{LOSS BY TRUE DISTRIBUTION D OF DATA} - \text{LOSS BY TRAINING DATA}$

ASSUME THAT THE LOSS FUNCTION IS CONVEX AND S-LIPSCHITZ.  $\underset{S \sim D^m}{\mathbb{E}} [L_0(A(S)) - L_s(A(S))] \leq \frac{2s^2}{\lambda m}$

FITTING-STABILITY TRADEOFF:  $\underset{S}{\mathbb{E}} [L_0(A(S))] \leq L_0(\bar{w}^*) + \lambda \|\bar{w}^*\|^2 + \frac{2s^2}{\lambda m}$

THE REGULARIZATION PATH:

RLM RULE:  $\bar{w}(\lambda) = \underset{\bar{w}}{\operatorname{argmin}} (L_s(\bar{w}) + \lambda \|\bar{w}\|^2)$

$\cdot \bar{w}(\lambda=0) : \text{MINIMIZE JUST } L_s(\bar{w})$

$\cdot \bar{w}(\lambda=\infty) : \text{MINIMIZING } \|\bar{w}\|^2 \text{ HAS } \lambda \|\bar{w}\| \xrightarrow{\lambda \|\bar{w}\| \gg 0} 0$

SAMPLE COMPLEXITY OF SOFT-SUM.  $\left\{ \begin{array}{l} \cdot \text{IF WE SET } \lambda = \frac{1}{B} \sqrt{\frac{2s^2}{m}} \text{ AND SINCE } l^{\text{Hinge}} \text{ IS AN UPPER BOUND FOR 0-1 LOSS:} \\ \rightarrow \underset{S \sim D^m}{\mathbb{E}} [L_0^{0-1}(A(S))] \leq \min_{\|\bar{w}\|=1 \leq B} \{ L_s^{\text{Hinge}}(\bar{w}) \} + \sqrt{\frac{8s^2 B^2}{m}} \end{array} \right.$

MARGIN / NORM VS DIMENSIONALITY:  $\left\{ \begin{array}{l} \cdot \text{VC DIMENSION OF LEARNING HALF SPACES} \quad \cdot \text{DIMENSION OF} \\ \rightarrow \uparrow \text{d} \rightarrow \uparrow \text{SAMPLE COMPLEXITY} \\ \cdot \text{FOR SVMs: SAMPLE COMPLEXITY } \left( \frac{s}{\delta} \right)^2 = \left( s^2 B^2 \right) \end{array} \right.$

KERNELS:

POLYNOMIAL MAPPINGS:  $\cdot P(x) = \sum_{j=0}^K w_j x^j = \langle \bar{w}, \psi(x) \rangle / \psi(x) = (1, x, x^2, \dots, x^K)$

KERNEL TRICK:  $K(\bar{x}_1, \bar{x}_2) = \langle \psi(\bar{x}_1), \psi(\bar{x}_2) \rangle$

THE REPRESENTER THEOREM:

CONSIDER ANY LEARNING RULE,  $f : \mathbb{R}^m \mapsto \mathbb{R}$ , OF THE FORM:

$\bar{w}^* = \underset{\bar{w}}{\operatorname{argmin}} \left\{ f(\langle \bar{w}, \psi(\bar{x}_1) \rangle, \dots, \langle \bar{w}, \psi(\bar{x}_m) \rangle) + \lambda \|\bar{w}\|^2 \right\} \rightarrow \underset{\bar{a} \in \mathbb{R}^m}{\operatorname{argmin}} \left\{ f(\bar{a}^\top \bar{G} \bar{x}) + \lambda \bar{a}^\top \bar{G} \bar{a} \right\}$

$\rightarrow \exists \bar{a} \in \mathbb{R}^m / \bar{w}^* = \sum_{i=1}^m \alpha_i \cdot \psi(\bar{x}_i) \rightarrow \bar{w}^* \text{ IS A LINEAR COMBINATION OF } \psi(\bar{x}_i)$

• POSITIVE SEMI-DEFINITE  
• NOT ALWAYS (L)  
 $G = \begin{pmatrix} \langle \psi(\bar{x}_1), \psi(\bar{x}_1) \rangle & \cdots & \langle \psi(\bar{x}_1), \psi(\bar{x}_m) \rangle \\ \vdots & \ddots & \vdots \\ \langle \psi(\bar{x}_m), \psi(\bar{x}_1) \rangle & \cdots & \langle \psi(\bar{x}_m), \psi(\bar{x}_m) \rangle \end{pmatrix}$

$\cdot \|\bar{w}\|^2 = \bar{a}^\top G \bar{a} \quad b_{ij} = K(\bar{x}_i, \bar{x}_j)$

NEW INSTANCE  
 $\cdot \langle \bar{w}, \psi(\bar{x}) \rangle = \sum_i K(\bar{x}_i, \bar{x})$

DEGREE OF POLYNOMIAL  
USED TO TRANSFORM  $\bar{x}_1, \bar{x}_2$

### POLYNOMIAL KERNELS :

$$K(\bar{x}_1, \bar{x}_2) = (1 + \langle \bar{x}_1, \bar{x}_2 \rangle)^k$$

GAUSSIAN KERNELS (RBF):  $K(\bar{x}_1, \bar{x}_2) = e^{-\frac{\|\bar{x}_1 - \bar{x}_2\|}{2\sigma}}$

CAN MAP INFINITE  
FEATURES  
IT CAN LEARN  
ANY POLYNOMIAL

LEMMA (MERCER'S CONDITION): A SYMMETRIC FUNCTION  $K: X \times X \rightarrow \mathbb{R}$  IMPLEMENTS AN INNER PRODUCT IN SOME HILBERT SPACE  $\Leftrightarrow \forall \bar{x}_1, \dots, \bar{x}_m : K_{ij} = K(\bar{x}_i, \bar{x}_j)$  IS A POSITIVE, SEMI-DEFINITE MATRIX

### CURSE OF DIMENSIONALITY

### DIMENSIONALITY REDUCTION:

$O(nd^2 + d^3)$

PCA: WE WANT TO BE ABLE TO APPROXIMATELY RECOVER  $\bar{x}$  FROM  $\tilde{x} = W\bar{x}$

IN PCA:  $\tilde{x} = U\bar{x} = UW\bar{x}$  LINEAR  
RECOVERY

RESULTS "APPROXIMATE RECOVERY", SOLVE:

$$\underset{W \in \mathbb{R}^{d \times d}, U \in \mathbb{R}^{d \times n}}{\text{argmin}} \left\{ \sum_{i=1}^m \| \bar{x}_i - UW\bar{x}_i \|^2 \right\}$$

$\begin{cases} \text{ORIGINAL FEATURES} \\ \text{RECOVERED FEATURES} \end{cases}$

$\begin{cases} \sim m: \\ \# \text{SAMPLES} \end{cases}$

### THEOREM:

LET  $A = \sum_{i=1}^m \bar{x}_i \bar{x}_i^T$ , USE  $\bar{u}_1, \dots, \bar{u}_n$  AS THE  $n$  ORTHONORMAL EIGENVECTORS OF  $A$ :

$\rightarrow$  SOLUTION OF PCA:  $\begin{cases} \text{COLUMNS OF } U \text{ SET TO } \bar{u}_1, \dots, \bar{u}_n \\ W = U^T \end{cases}$

N.B.:  $\frac{1}{m-1} (\bar{x} - \mu)^T (\bar{x} - \mu) = \Sigma = \text{COVARIANCE MATRIX}$

### PCA AND SVD:

IN PCA:  $\Sigma = \frac{1}{n-1} X^T X$ , ASSUMING DATA IN  $X$  ARE ZERO-CENTRED

$\Sigma$  IS SYMMETRIC  $\rightarrow \Sigma$  CAN BE DIAGONALIZED.  $\Sigma = V \cdot L \cdot V^T$

WHERE:

$$V = \begin{bmatrix} | & | & | \\ \bar{v}_1 & \dots & \bar{v}_n \end{bmatrix}$$

EIGENVECTOR  $\bar{v}_i$

EIGENVALUE  $\lambda_i$

DIAGONAL MATRIX

$L = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \lambda_n \\ 0 & & \lambda_d \end{bmatrix}$

$/ \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$

$\bar{v}_i$ : PRINCIPAL AXES / PROJECTION OF DATA ON  $\bar{v}_i \rightarrow$  PC;

$$PC_i = \begin{bmatrix} X \\ \vdots \\ \bar{v}_i \end{bmatrix}$$

COORDINATES OF  $i$ -TH POINT

$\rightarrow X_{i,1:n} = [x_i] \begin{bmatrix} V \end{bmatrix}$  ON NEW PC SPACE

PROPORTION OF VARIANCE EXPLAINED:  $PVE = \frac{\lambda_i}{\sum_{j=1}^d \lambda_j} = \frac{\lambda_i}{\text{trace}(\Sigma)}$

① DO NOT DISTORT DISTANCES:

$$\|\bar{x}_i - \bar{x}_j\| \approx \|\bar{v}_i - \bar{v}_j\|, \forall i, j$$

$$\rightarrow \frac{\|W\bar{x}_i - W\bar{x}_j\|}{\|\bar{x}_i - \bar{x}_j\|} \approx 1, \forall i, j$$

②  $\forall \bar{x} \in Q, Q = \{\bar{x}_i - \bar{x}_j : i, j \in [m]\}$ , WE HAVE  $\frac{\|Wx\|}{\|\bar{x}\|} \approx 1$

RANDOM PROJECTIONS:  $\sim O(nd \log n)$

DATA PROJECTED IN  $k$ -DIMENSIONAL SPACE.

$W$  IS A RANDOM MATRIX:  $W_{i,j} \sim N(0, \frac{1}{n})$

LET  $\bar{w}_i$  BE THE  $i$ -TH ROW OF  $W$ :

IN FACT:

$$\rightarrow E[\|W\bar{x}\|^2] = \|\bar{x}\|^2 \text{ PROPERTY SATISFIED} \quad \rightarrow \|W\bar{x}\|^2 \sim \chi^2_n$$

$$\rightarrow P\left[\left|\frac{\|W\bar{x}\|^2}{\|\bar{x}\|^2} - 1\right| > \epsilon\right] \leq 2e^{-\frac{\epsilon^2 n}{6}} \text{ MOFFONI'S INEQUALITY}$$

LEMMA (JOHNSON-LINDSTROMSSON LEMMA): LET  $Q$  BE A FINITE SET OF VECTORS IN  $\mathbb{R}^d$

$$\rightarrow \epsilon = \sqrt{\frac{6 \cdot \log(2/\delta)}{n}} \leq 3$$

PROPERTY ② STATISFIED, WITH C.I. 1-3

## COMPRESSED SENSING:

- PRIOR ASSUMPTION:  $\mathbf{U}^* = \mathbf{V}^*$ ,  $\langle \bar{\mathbf{v}}_i, \bar{\mathbf{u}}_j \rangle = 0$ ,  $\|\bar{\mathbf{v}}_i\|_2^2 = 1$
- $\bar{\mathbf{x}} \approx \mathbf{U} \bar{\mathbf{z}}$  /  $\mathbf{U}$  is ORTHONORMAL |  $\|\bar{\mathbf{z}}\|_0 := |\{i : \alpha_i \neq 0\}| \leq s$ , small
- ANY SPARSE SIGNAL CAN BE FULLY RECONSTRUCTED IF IT WAS COMPRESSED BY  $\bar{\mathbf{x}} \mapsto \mathbf{W} \bar{\mathbf{x}}$  /  $\mathbf{W}$  SATISFIES RESTRICTED ISOPERIMETRIC PROPERTY (RIP)
- RESTRICTED ISOPERIMETRIC PROPERTY:

A matrix  $\mathbf{W} \in \mathbb{R}^{n,d}$  is  $(\epsilon, s)$ -RIP if  $\forall \bar{\mathbf{x}} \neq 0 / \|\bar{\mathbf{x}}\|_0 \leq s$ :  $\left| \frac{\|\mathbf{W} \bar{\mathbf{x}}\|_2^2}{\|\bar{\mathbf{x}}\|_2^2} - 1 \right| \leq \epsilon \rightarrow 1 - \epsilon \leq \frac{\|\mathbf{W} \bar{\mathbf{x}}\|_2^2}{\|\bar{\mathbf{x}}\|_2^2} \leq 1 + \epsilon$

## PCA VS RANDOM PROJECTIONS:

- IF DATA IS  $\bar{\mathbf{c}}_1, \dots, \bar{\mathbf{c}}_d$ :
- IF  $d \gg n$ , DATA EXACTLY IN  $n$ -DIM SUBSPACE:

{ RANDOM PROJECTIONS  $\rightarrow$  **PERFECT**  
PCA  $\rightarrow$  **PAI**L }

{ RANDOM PROJECTIONS  $\rightarrow$  **MIGHT FAIL**  
PCA  $\rightarrow$  **PERFECT**

$$\bar{\mathbf{x}} \mapsto \mathbf{W} \bar{\mathbf{x}}$$

↳ 2D: DATA IN  $\mathbb{R}^2$  AND 1 PC

(BUT WITH NON-LINEAR RECONSTRUCTION)

• PCA: OPTIMAL IF LINEAR RECONSTRUCTION AND ERROR IS SQUARED DISTANCE

• RANDOM PROJECTIONS: PRESERVE DISTANCES, EXACT RECONSTRUCTION FOR SPARSE VECTORS

## DECISION TREES:

- PREDICTION FUNCTION:  $y(\bar{\mathbf{x}}) = \sum_{w \in W} y^w(\bar{\mathbf{x}}) \cdot \mathbb{1}_{[\bar{\mathbf{x}} \in R_w]}$

### TRAINING LOSS:

GIVEN  $y(\bar{\mathbf{x}})$  AS ABOVE, TRAINING SET  $\tilde{T} = \{(\bar{\mathbf{x}}_i, y_i)\}_{i=1}^n$ :

$$l_{\tilde{T}}(y) = \sum_{w \in W} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[\bar{\mathbf{x}}_i \in R_w]} \text{Loss}(y_i, y^w(\bar{\mathbf{x}}_i))$$

CONTRIBUTION OF REGIONAL  $y^w$  TO OVERALL TRAINING LOSS

### REGIONAL PREDICTION FUNCTION

GIVEN  $P_z^w = \frac{1}{n_w} \sum_{\{(\bar{\mathbf{x}}_i, y_i) \in \tilde{T} : \bar{\mathbf{x}}_i \in R_w\}} \mathbb{1}_{[y_i = z]}$

CROSS  $y^w(\bar{\mathbf{x}}) = \text{any max } \{P_z^w\}_{z \in \{0, \dots, c-1\}}$

REF  $y^w(\bar{\mathbf{x}}) = \bar{y}_{R_w} := \frac{1}{n_w} \sum_{\{(\bar{\mathbf{x}}_i, y_i) \in \tilde{T} : \bar{\mathbf{x}}_i \in R_w\}} y_i$

### OPTIMAL SPLITTING RULES:

REC  $\frac{1}{n} \sum_{(\bar{\mathbf{x}}_i, y_i) \in \tilde{T} : x_j \in \mathcal{E}} (y_i - \bar{y}_e)^2 + \frac{1}{n} \sum_{(\bar{\mathbf{x}}_i, y_i) \in \tilde{T} : x_j \in \mathcal{S}} (y_i - \bar{y}_s)^2$

CROSS  $\frac{1}{n} \sum_{(\bar{\mathbf{x}}_i, y_i) \in \mathcal{E}_i} \mathbb{1}_{[y_i \neq y_e^*]} + \frac{1}{n} \sum_{(\bar{\mathbf{x}}_i, y_i) \in \mathcal{S}_i} \mathbb{1}_{[y_i \neq y_s^*]} / y^* : \text{most relevant class (majority vote)}$

- IMPURITY:
- OTHER IMPURITY

ENTROPY:  $-\sum_{z=0}^{c-1} P_z \cdot \log_2(P_z)$

GINI:  $\frac{1}{2} \left( 1 - \sum_{z=0}^{c-1} P_z^2 \right)$

... BUT: HIGH CORRELATED PREDICTION

$\Rightarrow \downarrow G^2$

$\Rightarrow$  BOOTSTRAPPED

WE CAN SUBSTITUTE THEM WITH BOOTSTRAPPED ONES: BOOTSTRAPPING:

WE CAN OBTAIN  $\tilde{T}_1, \dots, \tilde{T}_B$  BY REAMPLING FROM A

SINGLE FIXED TRAINING SET  $\tilde{T}$ , AND USE THEM TO TRAIN B SEPARATE MODELS

WE OBTAIN A BAGGED ESTIMATOR OF THE FORM:

$$g_{\text{bag}}(\bar{\mathbf{x}}) = \frac{1}{B} \sum_{b=1}^B \bar{y}_{\tilde{T}_b}(\bar{\mathbf{x}})$$

RANDOM FOREST  $\rightarrow$  ONLY A SUBSET OF FEATURES CONSIDERED DURING TREE CONSTRUCTION

BOOSTING:  $\rightarrow$  IF WEAK LEARNER  $\rightarrow$  **1 ACCURACY**

2. "BOOST"  $y_0$  TO A NEW LEARNER:  $y_1 := y_0 + h_1$ ,

$$\int h_1 = \underset{h \in H}{\text{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, y_0(\bar{\mathbf{x}}_i) + h(\bar{\mathbf{x}}_i)) \right\}$$

### GRADIENT BOOSTING:

$$y_B(\bar{\mathbf{x}}) = y_0(\bar{\mathbf{x}}) + \sum_{b=1}^B h_b(\bar{\mathbf{x}})$$

$$\frac{z = y_{b+1}(\bar{\mathbf{x}})}{\rightarrow - \frac{\delta \text{Loss}(y_i, z)}{\delta z} = \frac{\delta(y_i - z)^2}{\delta z}} = 2(y_i - y_{b+1}(\bar{\mathbf{x}})) = 2e_i^{(b)}$$

$$\therefore \sqrt{\text{BOOSTING OPERATION:}}$$

$$(a_b, c_b) = \underset{a \geq 0, c \in C}{\text{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, y_{b+1}(\bar{\mathbf{x}}_i) + a \bar{c}(\bar{\mathbf{x}}_i)) \right\}$$

### ADA BOOST: