

Lecture 5 - scribbles

Tuesday, September 19, 2017
14:37

aside on proofs:

assumptions \rightsquigarrow use logical laws like $A \Rightarrow B$ \rightsquigarrow conclusion

example of thm: $X \perp\!\!\!\perp Y \Leftrightarrow p(x|y) = p(x) \quad \forall x, y$

$X \perp\!\!\!\perp Y$
(assumption)

by definition \Rightarrow (*) $p(x, y) = p(x)p(y) \quad \forall x, y$

suppose y st. $p(y) \neq 0$
by definition $p(x|y) = \frac{p(x, y)}{p(y)} \stackrel{\text{by (*)}}{=} \frac{p(x)p(y)}{p(y)} = p(x)$

\Downarrow only if ("equivalent")

(proof by contradiction: use fact that $(A \Rightarrow B) \Leftrightarrow (\neg B \Rightarrow \neg A)$)

today: statistical decision theory

(frequentist) statistical decision theory

unknown distribution which models the "world"

formal setup: • random observation $D \sim \mathcal{D}$ (perhaps P_0)

• action space \mathcal{A}

• loss: $L(D, a) =$ loss of choosing action a $\%$ describes the

when "true world" is $p \in \mathcal{P}$ goal/task

(if you have parametric model in mind; often write $L(\theta, a)$ where θ is s.t. $p = p_\theta$)

• $\mathcal{S} : \mathcal{D} \rightarrow \mathcal{A}$ "decision rule"

\mathcal{D} set of parameters for \mathcal{P}_θ

examples: a) if $\mathcal{A} = \mathcal{H}$ for a parametric family \mathcal{P}_θ

\mathcal{S} is then parameter estimator from data

typical loss $L(\theta, a) = \|\theta - a\|_2^2$ "squared loss"

[more specifically $\mathcal{D} = (X_1, \dots, X_n)$ where $X_i \stackrel{iid}{\sim} p_\theta$ (θ is unknown)]

$\mathcal{S}(\mathcal{D}) = \hat{\theta}$ $L(\theta, \mathcal{S}(\mathcal{D})) = \|\theta - \hat{\theta}\|_2^2$

suppose \mathcal{P}_θ is Gaussian family $\{N(\mu, 1) : \mu \in \mathbb{R}\}$

$\mathcal{H} = \mathbb{R}$ (set of possible means)

$\mathcal{S}(\mathcal{D}) = \frac{1}{n} \sum X_i$

b) $\mathcal{A} = \{0, 1\}$; this is basically hypothesis testing

here \mathcal{S} describes a statistical test

c) machine learning is learning a prediction rule

here $\mathcal{D} = ((X_i, Y_i))_{i=1}^n$

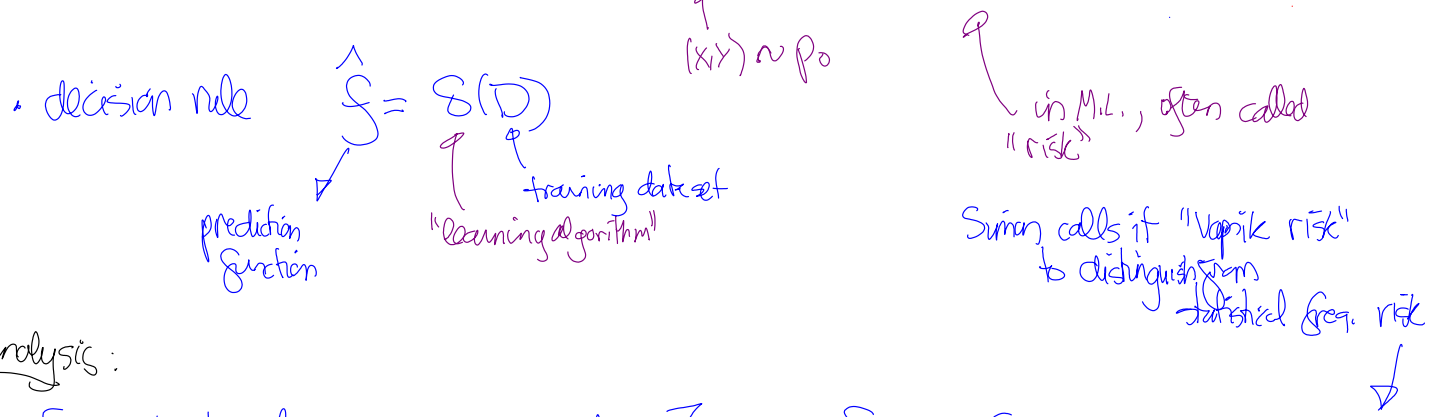
$X_i \in \mathcal{X}$ (input space)
 $Y_i \in \mathcal{Y}$ (output space)

e.g. $\mathcal{Y} = \{0, 1\}$ for binary classification

if p_θ gives joint on (X_i, Y_i) ; then $\mathcal{D} \sim \mathcal{D}$ where $\mathcal{D} = \mathcal{X}^n \times \mathcal{Y}^n = \mathcal{P}_1(\mathcal{X}) \times \mathcal{P}_1(\mathcal{Y}) \times \dots \times \mathcal{P}_1(\mathcal{X} \times \mathcal{Y})$

then $D \sim P$ where $P = P_0^{\otimes n} = \underbrace{P_0 \otimes P_0 \otimes \dots \otimes P_0}_{n \text{ times}}$ (i.i.d. model)

$\mathcal{A} = \mathcal{Y}^{\mathcal{X}}$ (set of functions from $\mathcal{X} \rightarrow \mathcal{Y}$)
 in machine learning, $L(P, \mathcal{S}) \triangleq \mathbb{E}_{P_0} [l(Y, \mathcal{S}(X))]$ (prediction loss)
 "generalization error"

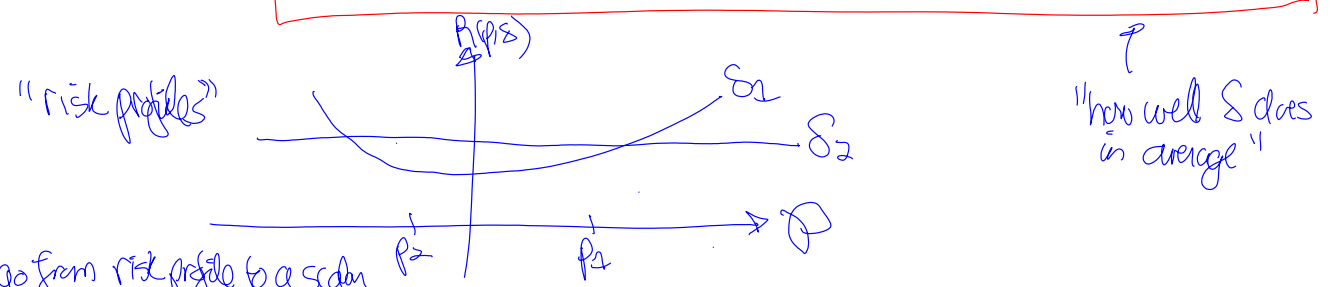


procedure analysis:

given this framework, how do we compare procedures? e.g. \mathcal{S}_1 vs \mathcal{S}_2

first property: (frequentist) risk $R(P, \mathcal{S}) \triangleq \mathbb{E}_P [L(P, \mathcal{S}(D))]$

\uparrow
 $D \sim P$



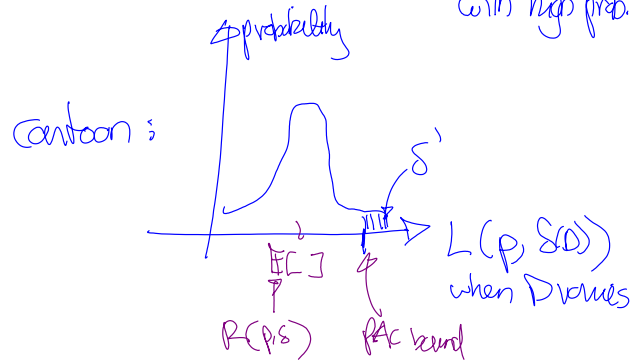
how to go from risk profile to a scalar number

* "minimax"
 \rightarrow look at $\max R(P, \mathcal{S})$

alternative in ML theory is PAC theory \rightarrow looking at tail bounds

$\theta \in \Theta$ " /
 * weighting = "prior" over Θ
 $\int_{\Theta} R(\theta, \delta) \pi(\theta) d\theta$

$P\{L(p, \delta(D)) \leq \delta'\}$
 "with high prob." statement



⊗ Bayesian decision theory

→ condition on data D

Bayesian posterior risk

$$R_B(a|D) = \int_{\Theta} L(\theta, a) \underbrace{p(\theta|D)}_{\text{posterior} \propto p(\theta)p(D|\theta)} d\theta$$

Bayesian optimal action

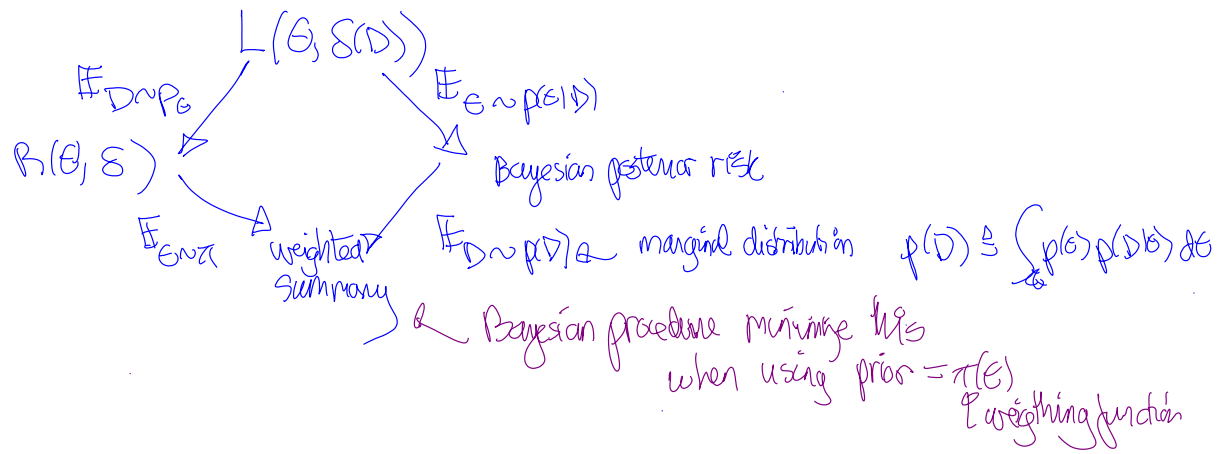
$$S_{\text{Bayes}}(D) = \underset{a \in A}{\text{argmin}} R_B(a|D)$$

example: if $A = \Theta$ ("estimation")

$$L(\theta, a) = \|\theta - a\|_2^2$$

then (exercise) $S_{\text{Bayes}}(D) = E[B|D]$ (posterior mean)

$$L(\theta, \delta(D))$$



examples of estimators: $\Theta \ni \mathcal{D} \rightarrow \Theta$

• MLE

• another: MAP, given a prior $p(\theta)$

'maximum a posteriori'

then pick $\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} p(\theta|D)$
 $\propto p(D|\theta)p(\theta)$

• method of moments

idea: find an injective mapping from Θ to 'moments'

and surjective on 'possible moments'

and then invert it from empirical moments

$\mathbb{E}[X]$
 $\mathbb{E}[X^2]$
 etc...

$$\hat{\mathbb{E}}[X] \triangleq \frac{1}{n} \sum_{i=1}^n X_i$$

example: for Gaussian $X \sim N(\mu, \sigma^2)$

$$\mathbb{E}[X] = \mu$$

$$\mathbb{E}[X] = \mu + \mu^2$$

$$\begin{pmatrix} \mathbb{E}[X] \\ \mathbb{E}[X^2] \end{pmatrix} = f(\mu, \sigma^2)$$

$$\begin{pmatrix} \hat{\mu} \\ \hat{\sigma}^2 \end{pmatrix} \triangleq f^{-1} \left(\begin{pmatrix} \mathbb{E}[X] \\ \mathbb{E}[X^2] \end{pmatrix} \right)$$

(here, this same as MLE
(property of exponential family))

⊗ this is useful for latent variable
models (e.g. mixture of Gaussians)
("spectral methods")
e.g.

• in the context of prediction $\mathcal{F} = \{f: \mathcal{X} \rightarrow \mathcal{Y}\}$ \mathcal{X} input space
 \mathcal{Y} output space

example of $\mathcal{S}: \mathcal{D} \rightarrow \mathcal{F}$

is using empirical "risk" minimization (ERM)
"Vapnik risk" i.e. generalization error

recall: $L(p, f) = \mathbb{E}_{(x,y) \sim p} [l(y, f(x))]$

replace $\mathbb{E} [l(y, f(x))] = \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i))$

$\hat{f}_{\text{ERM}} = \underset{f \in \mathcal{F}}{\text{argmin}} \mathbb{E} [l(y, f(x))]$
 $f \in \mathcal{F}$ hypothesis class