

Lecture 13: 3/1/17

Generative vs. Discriminative

Disc. model is a conditional dist.

$$p(y|\vec{x}, \vec{\Theta})$$

Gen. model is a joint dist.

$$p(\vec{x}, y|\vec{\Theta}) = p(y|\vec{x}) p(\vec{x})$$

Disc. model

Model of the Data instances

usually we write joint as $p(x|y)p(y)$

\Rightarrow Gen. vs. Disc. tradeoff can be understood as choosing whether or not to model $p(\vec{x})$

$$p(\vec{x}) = \sum_y p(x|y) p(y)$$

$$p(y|\vec{x}) = \dots \quad \text{by Bayes rule}$$

Bayes Classifier

← aka. "Bayes Optimal Classifier"
 ← aka. "Minimum Bayes Risk Decoder"

Two problems we care about:

① Density Estimation

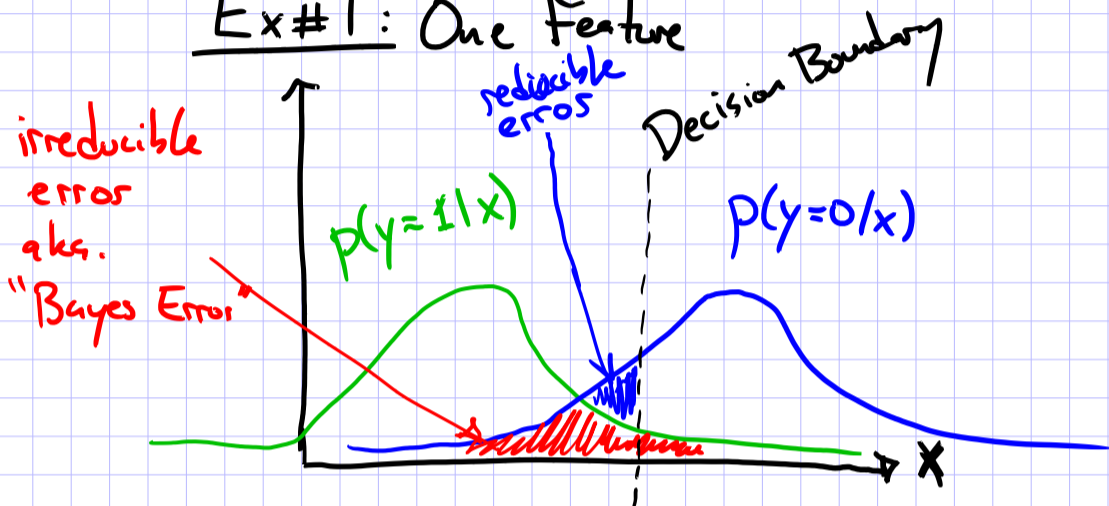
What does the distribution $p(x, y)$ look like?

② Choosing a Decision Function

How do we predict $\hat{y} = h(\vec{x})$? What is h ?

Not the same problem!

Ex#1: One Feature



Assume:

- Instances $x \in \mathcal{X}$ and labels $y \in \mathcal{Y}$
- Given a probability distribution $p(x, y)$
- Given a loss function $l(y, y')$

- Ex: 0-1 loss (for discrete y)

$$l(y, y') = \mathbb{I}(y \neq y') = \begin{cases} 1 & \text{if } y \neq y' \\ 0 & \text{otherwise} \end{cases}$$

- Ex: quadratic loss (for continuous y)

$$l(y, y') = (y - y')^2$$

Question: Given a new instance \vec{x} , what is the optimal prediction \hat{y} ?

Answer: $\hat{y} = h(x) = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} p(x, y) = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} p(y|x)$

Def: The expected loss $\text{risk}(h)$ of a classifier $h(\vec{x})$ is:

$$\begin{aligned} \text{risk}(h) &= E_{x, y \sim p(x, y)} [l(y, h(x))] \\ &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) l(y, h(x)) \end{aligned}$$

a functional
 a function of
 a function

Def: The Bayes Classifier h_{BC} is the h that minimizes the Bayes Risk, $risk(h)$

$$h_{BC} = \underset{h \in H}{\operatorname{argmin}} risk(h)$$

Def: The Bayes Error is $risk(h_{BC})$

↑ best we could possibly do.

Ex: Classification with 0/1 loss

Q: What is the Bayes Classifier?

$$\begin{aligned} \underline{h_{BC}} &= \underset{h}{\operatorname{argmin}} risk(h) \\ &= \underset{h}{\operatorname{argmax}} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) l(y, h(x)) \\ &= \underset{h}{\operatorname{argmax}} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \mathbb{I}(y \neq h(x)) \\ &= \underset{h}{\operatorname{argmin}} \sum_{x \in \mathcal{X}} G_x(h(x)) \end{aligned}$$

⇒ We want some h for each x , and it should return some \hat{y} that minimizes $G_x(\hat{y})$

$$\begin{aligned} h_{BC}(x) &= \underset{\hat{y}}{\operatorname{argmin}} G_x(\hat{y}) \\ &= \underset{\hat{y}}{\operatorname{argmin}} \sum_{y \in \mathcal{Y}} p(x,y) \mathbb{I}(y \neq \hat{y}) \quad \text{this term is 0 when } y = \hat{y} \\ &= \underset{\hat{y}}{\operatorname{argmin}} \sum_{y \neq \hat{y}} p(x,y) p(y|x) p(x) \\ &= \underset{\hat{y}}{\operatorname{argmin}} p(x) \sum_{y \neq \hat{y}} p(y|x) \\ &= \underset{\hat{y}}{\operatorname{argmin}} p(x) (1 - p(\hat{y}|x)) \\ &= \underset{\hat{y}}{\operatorname{argmin}} - p(\hat{y}|x) p(x) \\ &= \underset{\hat{y}}{\operatorname{argmax}} p(\hat{y}, x) \end{aligned}$$

⇒ $h(x) = \underset{y}{\operatorname{argmax}} p(x,y)$ is the Bayes Classifier!

Question: Where does the distribution $p(x,y)$?

Answer: It's usually unknown.

So we try to learn it from data.

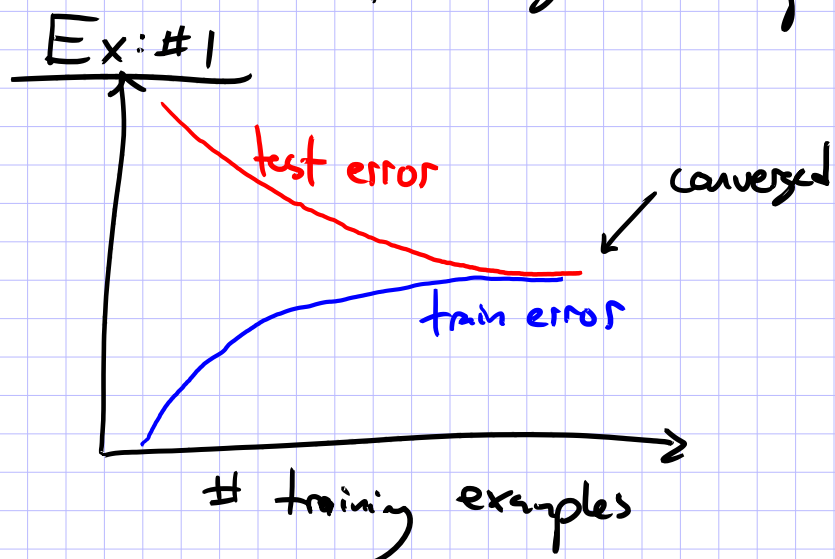
Maximum Likelihood Estimation

Question: Why should we use parameters that maximize likelihood?

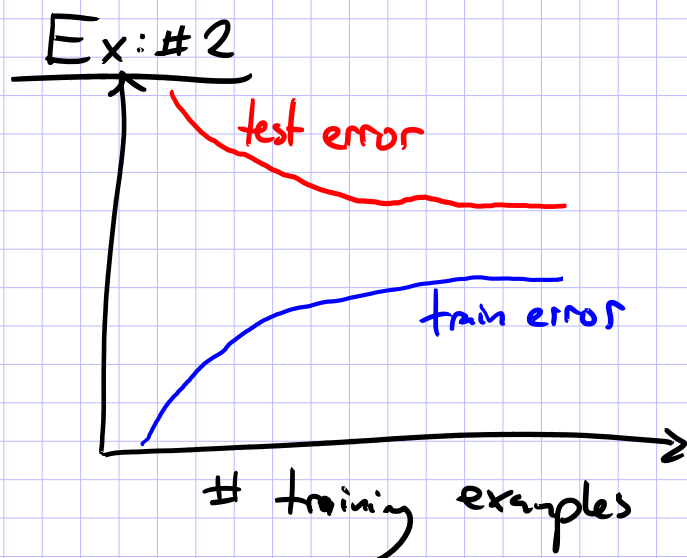
Answer: Because the MLE Θ^{MLE} is a consistent estimate of the true parameters Θ^*

Assume: $x^{(i)}, y^{(i)} \sim p^*(x, y | \Theta^*)$

Def: A learning method is consistent if model error on new samples converges to model error on the original sample (as the size of original sample increases)



Consistent!



Not Consistent!

Note: The average likelihood converges almost surely to the expected log-likelihood by the strong law of large numbers.

$$\frac{1}{N} \sum_{i=1}^N \log p(x^{(i)}, y^{(i)}) \xrightarrow{\text{a.s.}} E_{x, y \sim p^*} [\log p(x, y)]$$

Def: almost surely convergence

$$Pr(\lim_{n \rightarrow \infty} X_n = X) = 1$$

Def: Strong law of large numbers

$$\bar{X}_n \xrightarrow{\text{a.s.}} E[X]$$

as $n \rightarrow \infty$

So not too surprising that given

$$\hat{\Theta}_{MLE}^{(N)} = \operatorname{argmax}_{\Theta} \frac{1}{N} \sum_{i=1}^N \log p(x^{(i)}, y^{(i)} | \Theta)$$

Θ^* = true parameters

We have that $\hat{\Theta}_{MLE}^{(N)} \xrightarrow{\text{a.s.}} \Theta^*$ (under some pairwise indep.)

as $N \rightarrow \infty$

(Proof just requires KL divergence, and some prop.)