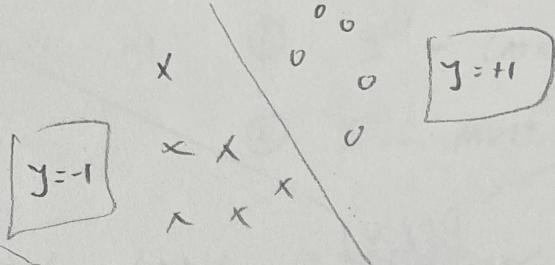


L3-1

Lesson 02 Recap

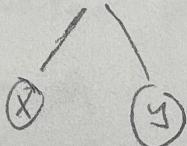
LINEAR DECISION Boundary Learning:



INFERENCE :

$$y = \text{Sign}(xw^T + b)$$

$P_d(x, y)$  ← unobserved, intangible to Learn  
Directly!

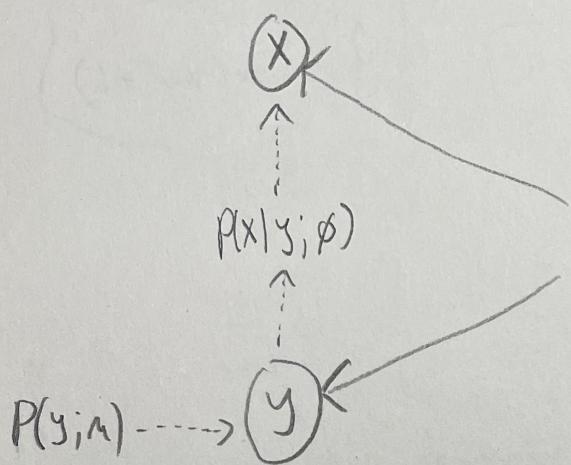


Generative MACHINE LEARNING IS AT GENERAL CLASS OF PROBLEMS OF THE FORM:

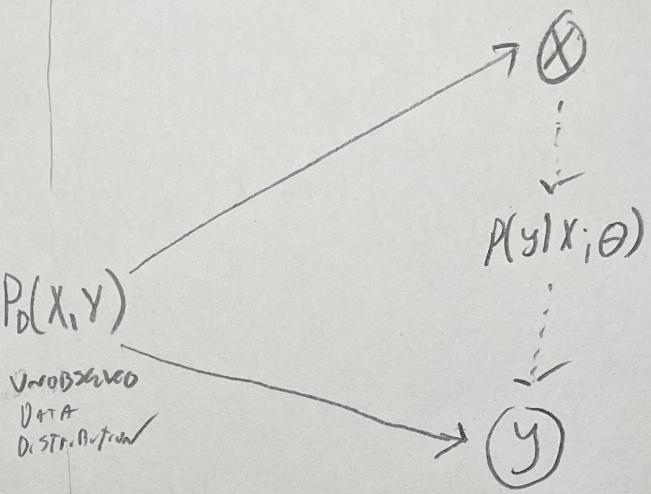
Given some input  $x^i$ , predict the most likely output  $y^i$  using the information from  $M$  independent draws from  $X^{(i)}, Y^{(i)} \sim P_d(x, y)$

## STATISTICAL ML TAXONOMY

### GENERATIVE MODELING



### DISCRIMINATIVE Modeling



### PARAMETER ESTIMATION

$$a) \hat{M} = \underset{M}{\operatorname{argmax}} \sum_{i=1}^m -\log P(y^{(i)}; M)$$

$$b) \hat{\phi} = \underset{\phi}{\operatorname{argmax}} \sum_{i=1}^m -\log P(x^{(i)}|y^{(i)}; \phi)$$

### INFERENCE

$$P(y|x; \hat{M}, \hat{\phi}) \propto P(x|y; \hat{\phi}) P(y; \hat{M})$$

↳ Bayes' rule

~~Bayes' rule~~

$$\hat{y} = \underset{y}{\operatorname{argmax}} \log P(x|y; \hat{\phi}) + \log P(y; \hat{M})$$

### PARAMETER ESTIMATION

$$a) \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log P(y^{(i)}|x^{(i)}; \theta)$$

### INFERENCE

$$\hat{y} = \underset{y}{\operatorname{argmax}} P(y|x; \hat{\theta})$$

# GENERATIVE MODEL: NAIVE BAYES

L3-3

THE DATA GENERATION PROCESS:

- $\forall i \in \{1, \dots, n\}$
- ①  $y^{(i)} \sim \text{CATEGORICAL } (\pi)$  LABEL
  - ②  $x^{(i)} \sim \text{MULTINOMIAL } (\phi)$  FEATURE

① CATEGORICAL LABEL DISTRIBUTION

$$P(y_i; \pi) = [\pi_1, \dots, \pi_K], \quad \sum_{k=1}^K \pi_k = 1$$

↳ ARBITRARY DISTRIBUTION OVER K CLASSES

② MULTINOMIAL FEATURE DISTRIBUTION

$X = \text{"CATS Rule, Dogs Drool" ... Sequence length: T = 4}$

$P(x|y, \phi)$ : MODELS FIRST-THIRD OF AN UNCONDITIONAL

$$= P(x_1, \dots, x_n | y; \phi) \rightarrow \text{intractable}$$

$$= P(x_1 | y) \cdot P(x_2 | x_1, y), \dots, P(x_n | x_1, \dots, x_{n-1}, y) \Rightarrow \text{intractable}$$

$\Rightarrow$  NAIVE BAYES ASSUMPTION: WORD COUNTS ARE INDEPENDENT GIVEN A CLASS LABEL Y.

$$= P(x_1 | y) \cdot P(x_2 | y), \dots, P(x_n | y) \Rightarrow \text{tractable!}$$

NAIVE BAYES Cont.. -

$P(x|y; \phi)$  is a multinomial distribution

$$= B(x) \prod_{j=1}^N \phi_{kj}^{x_j}, \quad B(x) = \frac{(\sum_{j=1}^N x_j)!}{\prod_{j=1}^N x_j!}$$

where  $\phi \in [0, 1]^{K \times N}$ ,  $\forall k \sum_{j=1}^N \phi_{kj} = 1$

NAIVE BAYES: ~~Parameter~~ ESTIMATION

$$\hat{\mu} = \underset{\mu}{\operatorname{ARGMAX}} \underbrace{\sum_{i=1}^M \log \mu_{y^{(i)}}}_{\ell(\mu)} \quad \text{s.t.} \quad \underbrace{\sum_{k=1}^K \mu_k = 1}_{C(\mu)}$$

Solve using method of Lagrange Multipliers:  $\forall k \nabla_{\mu_k} \ell(\mu) = \lambda \nabla_{\mu_k} C(\mu)$

$$\text{LHS } \forall k \nabla_{\mu_k} \left[ \sum_{i=1}^M \log \mu_{y^{(i)}} \right] = \sum_{i=1}^M \delta(y^{(i)}=k) \cdot \frac{1}{\mu_k} = \frac{M_k}{\mu_k}$$

$$\text{RHS } \forall k \lambda \nabla_{\mu_k} \left[ \sum_{k=1}^K \mu_k \right] = \lambda$$

∴  $\forall k \mu_k = \frac{M_k}{\lambda} \rightarrow$  Solve for  $\lambda$

$$\sum_{k=1}^K \mu_k = \frac{1}{\lambda} \sum_{k=1}^K M_k \iff 1 = \frac{M}{\lambda} \Rightarrow \lambda = M$$

∴  $\forall k \in \{1, \dots, K\} \quad \hat{\mu}_k = \frac{\sum_{i=1}^M \delta(y^{(i)}=k)}{M}$

↳ LABEL DISTRIBUTION IS THE RELATIVE FREQUENCY OF EACH CLASS IN THE DATA

3-5

NAIVE BAYES:  $\phi$  ESTIMATION

$$\hat{\phi} = \underset{\phi}{\operatorname{argmax}} \sum_{i=1}^M \log P(x_i) + \sum_{j=1}^N x_j^{(i)} \log \phi_{y^{(i)}, j}$$

 $\ell(\phi)$ 

$$\text{s.t. } \forall k \quad \sum_{j=1}^N \phi_{kj} = 1$$

 $c(\phi)$ (D) METHOD OF LAGRANGE MULTIPLIES:  $\forall k \quad \nabla_{\phi_{kj}} \ell(\phi) = \lambda \nabla_{\phi_{kj}} c(\phi)$ 

$$\text{RHS: } \forall k \quad \lambda \nabla_{\phi_{kj}} \left( \sum_{j=1}^N \phi_{kj} \right) = 1$$

$$\text{LHS: } \forall k \quad \nabla_{\phi_{kj}} \left( \sum_{i=1}^M \delta(y^{(i)}=k) \sum_{j=1}^N x_j^{(i)} \log \phi_{kj} \right) = \sum_{i=1}^M \delta(y^{(i)}=k) \cdot \frac{x_j^{(i)}}{\phi_{kj}}$$

(D) SOLVE FOR  $\lambda$ :

$$\forall_{k,j} \quad \lambda = \sum_{i=1}^M \delta(y^{(i)}=k) \frac{x_j^{(i)}}{\phi_{kj}}$$

$$\forall_k \quad \lambda \sum_{j=1}^N \phi_{kj} = \sum_{i=1}^M \delta(y^{(i)}=k) \sum_{j=1}^N x_j^{(i)}$$

$$\forall_k \quad \lambda = \sum_{i=1}^M \delta(y^{(i)}=k) \sum_{j=1}^N x_j^{(i)}$$

$\lambda_k$  = TOTAL  
NUMBER  
OF WORDS  
IN DOCUMENTS  
WITH LABEL  
 $y=k$

(D) SOLVE FOR  $\hat{\phi}$ :

$$\forall_{k,j} \quad \sum_{i=1}^M \delta(y^{(i)}=k) x_j^{(i)} = \lambda \phi_{kj}$$

$$\therefore \hat{\phi}_{kj} = \frac{\sum_{i=1}^M \delta(y^{(i)}=k) x_j^{(i)}}{\sum_{i=1}^M \delta(y^{(i)}=k) \sum_{j=1}^N x_j^{(i)}} = \frac{\text{count}(y=k, x=j)}{\text{count}(y=k, \text{All words})}$$

## NAIVE BAYES INFERENCE

WE USE BAYES' RULE TO COMPUTE  $P(y|x)$ :

$$P(y|x) \propto P(x|y; \hat{\phi}) P(y; \hat{\lambda})$$

$$\hat{y} = \underset{y}{\operatorname{argmax}} \log P(x|y; \hat{\phi}) + \log P(y; \hat{\lambda})$$

$$\hat{y} = \underset{k}{\operatorname{argmax}} \sum_{j=1}^N x_j \log \hat{\phi}_{kj} + \log \hat{\lambda}_k$$

NOTE: THIS IS A LINEAR MODEL!

$$\hat{y} = \underset{k}{\operatorname{argmax}} \vec{x} (\log \hat{\phi}_k)^T + \log \hat{\lambda}_k$$

Smoothing: MANY WORDS WILL NEVER APPEAR IN DOCUMENTS AT A GIVEN CLASS, AND THEREFORE MUCH OF  $\phi$  WILL BE ZERO IN REAL WORLD PROBLEMS.  
THIS IS PROBLEEM BECAUSE  $\log(0)$  IS UNDEFINED.

LAPLACE SMOOTHING:  $\hat{\phi}_{kj} = \frac{\alpha + \text{count}(y=k, x=j)}{\alpha N + \text{count}(y=k, \text{all } j)}$

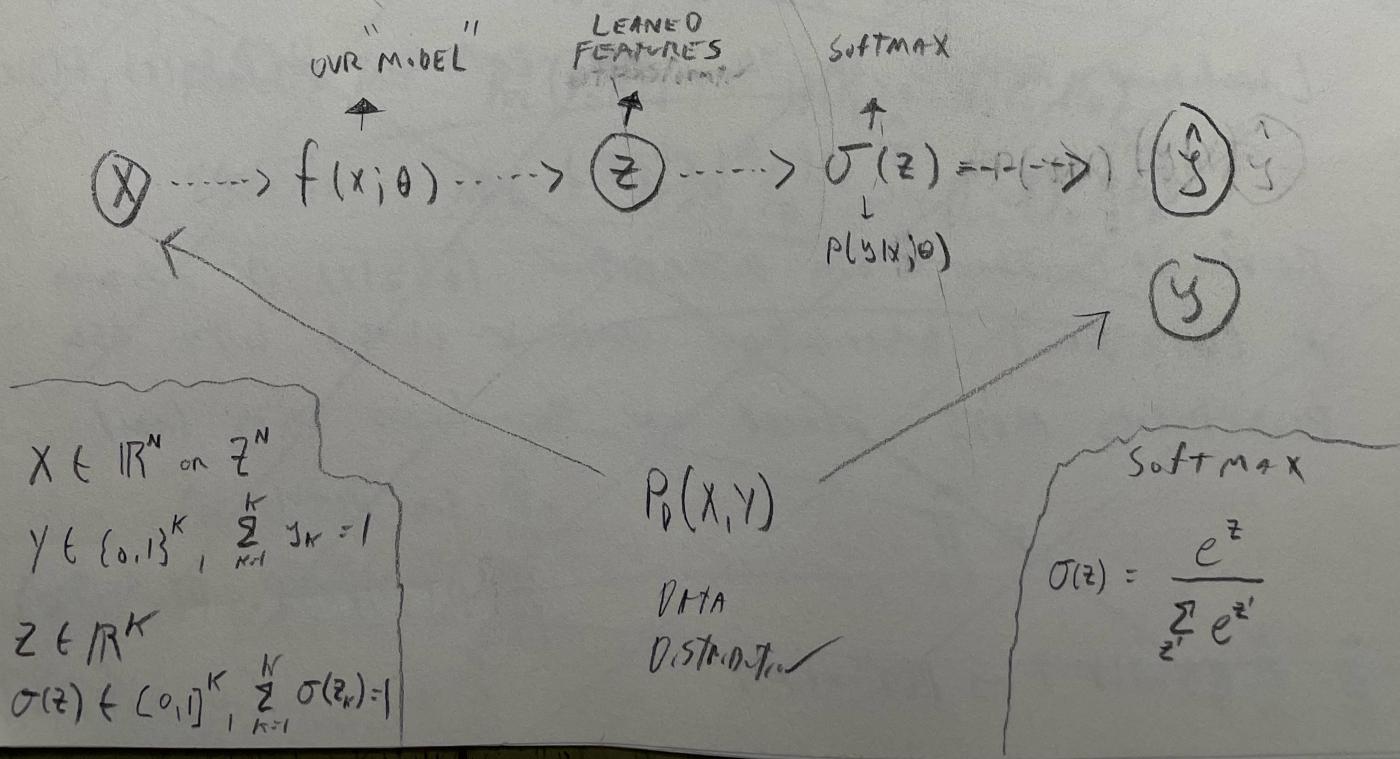
Where  $\alpha \in \mathbb{R}$ ,  $\alpha \geq 0$

## DISCRIMINATIVE LEARNING

L3-7

our APPROACH TO DISCRIMINATIVE LEARNING  
 IS MOST OFTEN DONE BY TRANSFORMING  
 OUR FEATURE SPACE onto our LABEL  
 SPACE, CONVERTING THOSE TRANSFORMED FEATURES  
 INTO A PROBABILITY DISTRIBUTION over the  
 CLASS LABELS, AND PERFORMING GRADIENT  
 DESCENT TO LEARN THE TRANSFORMATION  
 THAT YIELDS HIGH LIKELIHOOD  $P(Y; \theta)$ .  
 THIS ENCOMPASSES NEARLY ALL OF DEEP  
 LEARNING.

THE HIGH LEVEL MACHINERY IS:



# Discriminative Learning: $\hat{\theta}$ Estimation

L3-8

Negative Log-Likelihood:

$$NLL = - \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta)$$

REMEMBER THAT OUR EMPIRICAL LABEL DISTRIBUTION IS  
A ONE-HOT ENCODED VECTOR WITH ALL PROBABILITY  
MASS PLACED ON THE GROUND TRUTH LABEL

$$\begin{aligned} NLL &= - \sum_{i=1}^m \sum_{k=1}^K p_0(y^{(i)}=k | x^{(i)}) \log p(y^{(i)}=k | x^{(i)}; \theta) \\ &= - \sum_{i=1}^m E_{y^{(i)} \sim p_0(y|x)} [\log p(y^{(i)} | x^{(i)}; \theta)] \rightarrow \text{(cross entropy } H(p_0, p_\theta)) \end{aligned}$$

D RECALL THAT  $H(p_0) = 0$  BY DEFINITION (ONE-HOT ENCODING)

TRECALL THAT  $D_{KL}(p_0 || p_\theta) = H(p_0) + H(p_0, p_\theta)$

therefore  $NLL = H(p_0, p_\theta) = D_{KL}(p_0 || p_\theta)!$

## Softmax Regression

$$\Theta \rightarrow f(x_i; \Theta) \rightarrow z \rightarrow \sigma(z)$$

L3-

$$z = f(x; \theta) = x w^T + b \Rightarrow \theta = \{w, b\}$$

where  $w \in \mathbb{R}^{K \times N}$ ,  $b \in \mathbb{R}^K$

$$\sigma(z) = \frac{e^z}{\sum_{z'} e^{z'}} = \frac{e^{x w_k^T + b_k}}{\sum_{k'=1}^K e^{x w_{k'}^T + b_{k'}}}$$

$$\hat{y} = \underset{k}{\operatorname{argmax}} \sigma(z)$$

]

SOFTMAX  
AKA  $P(y|x; \theta)$

]

DECISION RULE

## PARAMETER ESTIMATION

$$NLL = - \sum_{i=1}^m \sum_{k=1}^K P_0(y_k^{(i)} | x^{(i)}) \log P(y_k^{(i)} | x^{(i)}; \theta)$$

$$= - \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log (P(y_k^{(i)} | x^{(i)}; \theta))$$

$$= - \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log \left[ \frac{e^{x^{(i)} w_k^T + b_k}}{\sum_{k'=1}^K e^{x^{(i)} w_{k'}^T + b_{k'}}} \right]$$

$$= - \sum_{i=1}^m \left[ \sum_{k=1}^K y_k^{(i)} (x^{(i)} w_k^T + b_k) - y_k^{(i)} \log \left[ \sum_{k'=1}^K e^{x^{(i)} w_{k'}^T + b_{k'}} \right] \right]$$

$$= - \sum_{i=1}^m \left[ \left( \sum_{k=1}^K y_k^{(i)} (x^{(i)} w_k^T + b_k) \right) - \log \left[ \sum_{k=1}^K e^{x^{(i)} w_k^T + b_k} \right] \right]$$

# SOFTMAX REGRESSION (CONTINUED)

L3-10

WE COMPUTE  $\hat{\theta}$  USING GRADIENT DESCENT

STEP 1: COMPUTE  $\nabla_{\theta} \text{NLL}(\theta | D)$

STEP 2: UPDATE  $\theta$ :  $\theta \leftarrow \theta - \gamma \nabla_{\theta} + E + R$

## $W$ ESTIMATION

$$\begin{aligned}\nabla_{w_k} \text{NLL}(w, b | D) &\in \mathbb{R}^N \\ &= - \sum_{i=1}^m y_k^{(i)} x^{(i)} - x^{(i)} \cdot \frac{e^{x^{(i)} w_k^T + b_k}}{\sum_{k'=1}^C e^{x^{(i)} w_{k'}^T + b_{k'}}} \\ &= - \sum_{i=1}^m x^{(i)} \left( P(y_k^{(i)} | x^{(i)}; \theta) - y_k^{(i)} \right)\end{aligned}$$

## $b$ ESTIMATION

$$\begin{aligned}\nabla_{b_k} \text{NLL}(w, b | D) &\in \mathbb{R} \\ &= - \sum_{i=1}^m y_k^{(i)} - \frac{e^{x^{(i)} w_k^T + b_k}}{\sum_{k'=1}^C e^{x^{(i)} w_{k'}^T + b_{k'}}} \\ &= \sum_{i=1}^m P(y_k^{(i)} | x^{(i)}; \theta) - y_k^{(i)}\end{aligned}$$

## GRADIENT DESCENT VERSIONS

"VANILLA":  $\nabla_{\theta} \sum_{i=1}^m -\log P(y^{(i)} | x^{(i)}; \theta)$

"STOCHASTIC":  $\nabla_{\theta} -\log P(y^{(i)} | x^{(i)}; \theta)$

"MINI-BATCH":  $\nabla_{\theta} \sum_{i=1}^{M'} -\log P(y^{(i)} | x^{(i)}; \theta)$   
where  $M' \ll m$

IN PRACTICE we typically USE MINI-BATCH GRADIENT DESCENT.