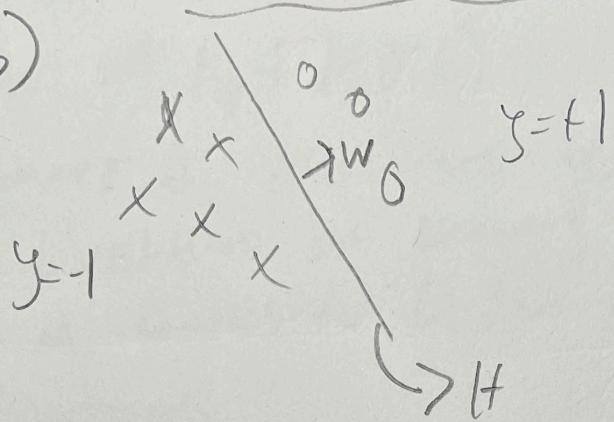


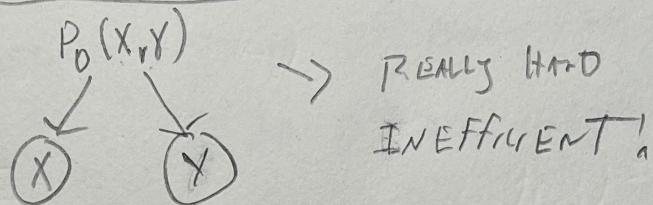
RECAP Lec 02 : Decision Boundary Learning

L3-1

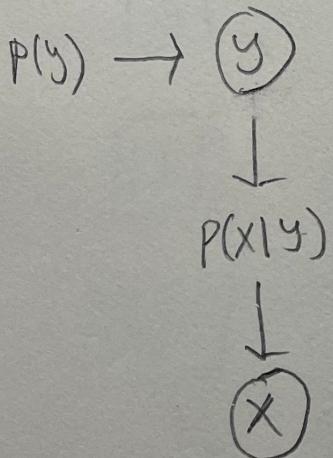
$$\hat{y} = \text{Sign}(x^T w + b)$$



STATISTICAL ML APPROACHES



GENERATIVE MODELING



PARAMETER ESTIMATION STEP:

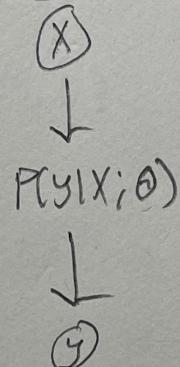
- (1) ESTIMATE $p(y; \theta)$
- (2) ESTIMATE $p(x|y; \theta)$

INFERENCE STEP: BAYES RULE

$$p(y|x) \propto p(x|y; \hat{\theta}) p(y; \hat{M})$$

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

DISCRIMINATIVE MODELING



PARAMETER ESTIMATION STEP:

- (1) ESTIMATE $p(y|x; \theta)$

INFERENCE:

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

PRINCIPLE OF MAXIMUM LIKELIHOOD

L3-2

- UNOBSERVED DATA DIST: $P_D(D) \sim \begin{cases} P(y) \\ P(y|x) \\ P(x|y) \end{cases}$
- APPROXIMATION OF P_D : $P(D;\theta)$
- GOAL: ESTIMATE, $\hat{\theta}$, THAT MAXIMIZES THE LIKELIHOOD OF DRAWING THE M SAMPLES IN D FROM $P(D;\theta)$.
- LIKELIHOOD FUNCTION:

$$L(\theta | D) = \prod_{i=1}^M P(D^{(i)}; \theta)$$

ASSUMPTIONS: ① All $D^{(i)}$ drawn from P_D
 i.i.D ② All $D^{(i)}$ are INDEPENDENT

- MAXIMUM LIKELIHOOD ESTIMATION (MLE):

L

$$\hat{\theta} = \underset{\theta}{\operatorname{ARGMAX}} \prod_{i=1}^M P(D^{(i)}; \theta)$$

LL

$$\hat{\theta} = \underset{\theta}{\operatorname{ARGMAX}} \sum_{i=1}^M \log P(D^{(i)}; \theta)$$

NLL

$$\hat{\theta} = \underset{\theta}{\operatorname{ARGMAX}} - \sum_{i=1}^M -\log P(D^{(i)}; \theta)$$

↳ $\hat{\theta} = \text{ML ESTIMATE OF } \theta$

GENERATIVE MODELING EXAMPLE: NB

L3-3

DATA GENERATING PROCESS:

$$Y_i \in \{1, \dots, m\} \quad \begin{array}{l} \textcircled{1} \quad y^{(i)} \sim \text{CATEGORICAL } (\mu) \\ \textcircled{2} \quad x^{(i)} \sim \text{MULTINOMIAL } (\phi) \end{array} \quad \begin{array}{l} \text{LABEL DIST} \\ \text{FEATURE DIST} \end{array}$$

EXPLANATION:

- $P(y; \mu) \Rightarrow$ OUR LABEL DIST CAN BE ANY ARBITRARY DMF OVER K -CLASSES

$$P(y; \vec{\mu}) = \left\{ \vec{\mu}_1, \vec{\mu}_2, \dots, \vec{\mu}_K \right\} \quad \sum_{k=1}^K \mu_{k,r} = 1$$

INTUITION: Single
Single FLIP of A K-SIDED COIN

$$\bullet P(\vec{x}|y; \phi)$$

INTUITION: SUM(x) NUMBER OF FLIPS OF N-SIDED COIN

$$P(\vec{x}|y; \phi) = B(x) \prod_{j=1}^N \phi_{kj}^{x_j} \quad B(x) = \frac{\left(\sum_{j=1}^N x_j \right)!}{\prod_{j=1}^N x_j!}$$

$$\triangleright P(\vec{x}|y; \phi) = P(x_1, \dots, x_n | y; \phi)$$

$$= P(x_1|y; \phi_1) P(x_2|x_1, y; \phi_2) \cdots P(x_n|x_1, \dots, x_{n-1}, y; \phi_n)$$

ϕ INTRACTABLE

▷ NB ASSUMPTION:

WORD COUNTS ARE INDEPENDENT CONDITIONED ON
SOME CLASS LABEL y .

$$\triangleright P(\vec{x}|y; \phi) = P(x_1|y; \phi) P(x_2|y; \phi) \cdots P(x_n|y; \phi)$$

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

NB PARAMETER ESTIMATION

L3-4

- $P(y; \lambda) \Rightarrow$ Intuitive Guessimate:

$$\forall_k \hat{\mu}_k = \frac{\text{Number of times } y=k \text{ in } D}{M}$$

$$= \frac{\sum_{i=1}^M \delta(y^{(i)}=k)}{M}$$

↳ LETS Prove this w/ MLE

$$\text{LL}(\lambda | D) = \sum_{i=1}^M \log \lambda y^{(i)} \quad \text{s.t. } \sum_{k=1}^K \mu_k = 1$$

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

→ Method of Lagrange Multipliers: $\forall_k \nabla_{\mu_k} \ell(\lambda) = \lambda \nabla_{\mu_k} (\ell(\lambda))$

$$\forall_k \nabla_{\mu_k} \left[\sum_{i=1}^M \log \lambda y^{(i)} \right] = \sum_{i=1}^M \delta(y^{(i)}=k) \cdot \frac{1}{\mu_k} = \frac{M_{y=k}}{\mu_k}$$

$$\forall_k \nabla_{\mu_k} \left[\sum_{k=1}^K \mu_k \right] = 1$$

$$\hookrightarrow \forall_k \frac{M_{y=k}}{\hat{\mu}_k} = \lambda \Rightarrow \boxed{\forall_k \hat{\mu}_k = \frac{M_{y=k}}{\lambda}} \Rightarrow \sum_{k=1}^K \hat{\mu}_k = \frac{1}{\lambda} \sum_{k=1}^K M_{y=k}$$

$$\hookrightarrow 1 = \frac{M}{\lambda} \Rightarrow \boxed{\lambda = M}$$

$$\hookrightarrow \boxed{\forall_k \hat{\mu}_k = \frac{\text{Count}(y=k)}{M} = \frac{\sum_{i=1}^M \delta(y^{(i)}=k)}{M}}$$

↳ MATCHES GUESSEIMATE!

(D) Estimate $P(\vec{x}|y; \phi)$ use an conditional INDEPENDENCE ASSUMPTION

$$\Rightarrow P(\vec{x}|y; \phi) = \prod_{j=1}^N P(x_j|y; \phi)$$

$$L(\phi|D) = \prod_{i=1}^m P(\vec{x}^{(i)}|y^{(i)}; \phi) = \prod_{i=1}^m B(x) \prod_{j=1}^N \phi_{y^{(i)}, j}^{x_j^{(i)}}$$

$$LL(\phi|D) = \sum_{i=1}^m \log B(x) + \sum_{j=1}^N x_j^{(i)} \log \phi_{y^{(i)}, j}$$

$$(D) \hat{\phi} = \underset{\phi}{\text{argmax}} \underbrace{LL(\phi|D)}_{\ell(\phi)} \quad \text{s.t.} \underbrace{\forall k \sum_{j=1}^N \phi_{kj}}_{C(\phi)} = 1$$

\Rightarrow Method of Lagrange Multipliers:

$$\forall k \nabla_{\phi_{kj}} \ell(\phi) = \lambda \nabla_{\phi_{kj}} C(\phi)$$

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

$$\forall k \nabla_{\phi_{kj}} \left(\sum_{i=1}^m \delta(y^{(i)}=k) \sum_{j=1}^N x_j^{(i)} \log \phi_{kj} \right) = \sum_{i=1}^m \delta(y^{(i)}=k) \frac{x_j^{(i)}}{\phi_{kj}}$$

$$(D) \forall k \lambda = \sum_{i=1}^m \delta(y^{(i)}=k) \frac{x_j^{(i)}}{\phi_{kj}} \Rightarrow \forall k \lambda \hat{\phi}_{kj} = \sum_{i=1}^m \delta(y^{(i)}=k) x_j^{(i)}$$

$$(D) \forall k \lambda \sum_{j=1}^N \hat{\phi}_{kj} = \sum_{i=1}^m \delta(y^{(i)}=k) \sum_{j=1}^N x_j^{(i)} \Rightarrow \text{BUT } \forall k \sum_{j=1}^N \phi_{kj} = 1$$

$$(D) \boxed{\forall k \lambda = \sum_{i=1}^m \delta(y^{(i)}=k) \sum_{j=1}^N x_j^{(i)}} \Rightarrow \boxed{\lambda_k = \text{TOTAL Number of words in documents with Label } y=k}$$

NB Parameter Estimation cont..

13-6

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

$$\nabla_k \lambda = \sum_{i=1}^M \delta(y^{(i)} = k) x_j^{(i)} = \lambda \hat{\phi}_{kj}$$

$$\nabla_k \hat{\phi}_{Nj} = \frac{\sum_{i=1}^M \delta(y^{(i)} = k) x_j^{(i)}}{\sum_{i=1}^M \delta(y^{(i)} = k) \sum_{j=1}^N x_j^{(i)}} = \frac{\text{Count(class=k, word=j)}}{\text{Count(class=k, all words)}}$$

NB INFERENCE: : now that we have $\hat{\lambda}, \hat{\phi}$, we can compute $P(y|x)$ using BAYES RULE:

$$P(y_{te} | x_{te}) \propto P(x_{te} | y_{te}; \hat{\phi}) P(y_{te}; \hat{\lambda})$$

$$\begin{aligned} \hat{y}_{te} &= \underset{j_{te}}{\operatorname{argmax}} \log P(x_{te} | y_{te}; \hat{\phi}) + \log P(y_{te}; \hat{\lambda}) \\ &= \underset{k}{\operatorname{argmax}} \sum_{j=1}^N x_{te;j} \log \hat{\phi}_{kj} + \log \hat{\lambda}_k \end{aligned}$$

$$= \underset{k}{\operatorname{argmax}} \left\{ \vec{x}_{te} (\log \hat{\phi}_{kj})^T + \log \hat{\lambda}_k \right\}$$

BUT WAIT: Many words will never appear in documents of a given class, therefore MANY $\hat{\phi}_{kj}$ ARE UNDEFINED.

$$\hat{\phi}_{kj} \leftarrow \frac{\alpha + \text{Count}(w_{ij})}{\alpha N + \text{Count}(K, \text{all } i)}$$

α can be any positive real number

NB INFERENCE

L3-7

- Given $\hat{m}, \hat{\phi}$: we can now use Bayes rule to compute the posterior $P(y|x)$:

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

$$\underset{y_{te}}{\text{LD}} = \underset{y_{te}}{\text{ARGMAX}} \log P(x_{te}|y_{te}; \hat{\phi}) + \log P(y_{te}; \hat{m})$$

$$= \underset{y_{te}}{\text{ARGMAX}} \sum_{j=1}^K x_{te,j} \log \hat{\phi}_{kj} + \log \hat{m}_k$$

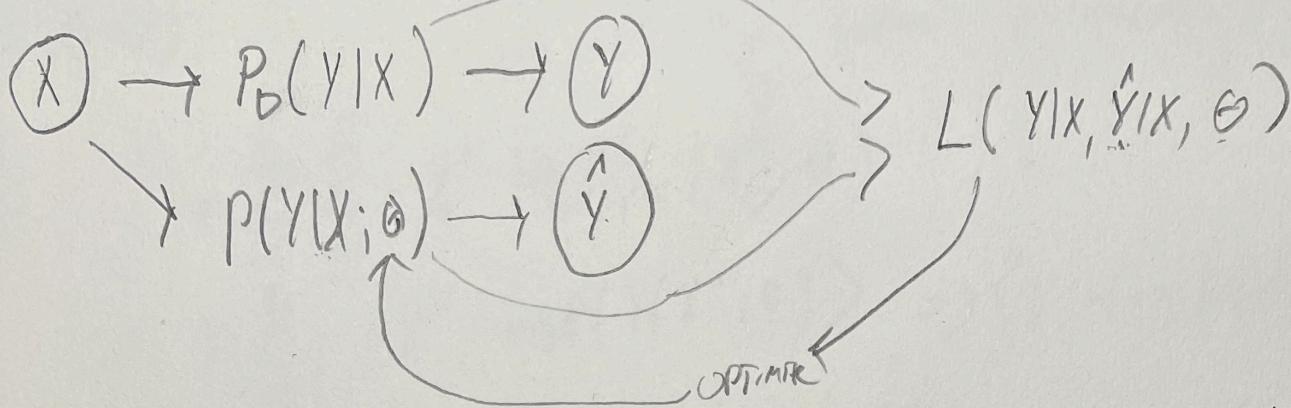
$$= \underset{k}{\text{ARGMAX}} \left[\vec{x}_{te} \cdot (\log \hat{\phi}_k)^T + \log \hat{m}_k \right]$$

\hookrightarrow Multinomial NB is a linear classifier

But wait: many k, j combinations will NOT have support in the data! $\Rightarrow \log 0$

Laplace Smoothing:

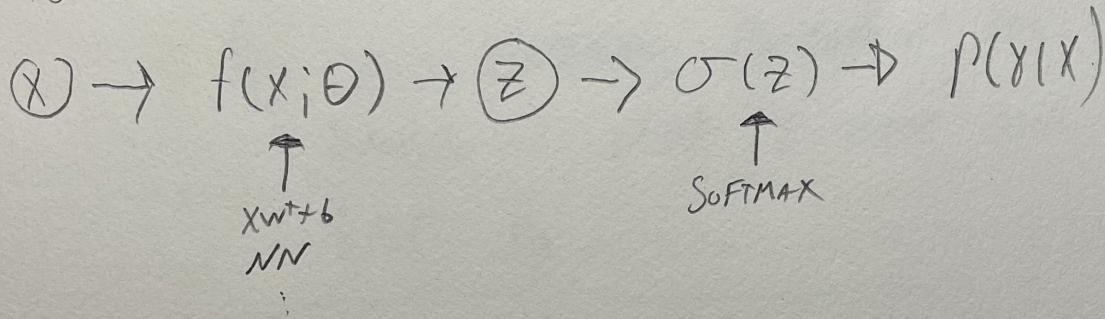
$$\hat{\phi}_{kj} = \frac{\alpha + \text{Count}(k, j)}{\alpha \cdot N + \text{Count}(k, \text{all } j)}$$



$P_D(Y|X)$ = UNOBSERVED DISTRIBUTION THAT Generates $Y|X$

$P(Y|X;\theta)$ = APPROX TO P_D PARAMETERIZED BY θ

GENERAL FRAMEWORK FOR DISCRIMINATIVE CLASSIFIERS:



$$f(x;\theta) : \mathbb{R}^N \rightarrow \mathbb{R}^K$$

$$\sigma(z) : \mathbb{R}^K \rightarrow [0,1]^K \Rightarrow \sigma_{\text{softmax}}(z_K) = \frac{e^{z_K}}{\sum_{k=1}^K e^{z_k}}$$

$\sum_{k=1}^K \sigma(z_k) = 1$ ← SATISFIES

$$\vec{y} \in \{0,1\}^K \rightarrow \text{Set, NOT A RANGE}$$

$$\sum_{k=1}^K y_k = 1$$

$$\prod_{k=1}^K \prod_{j=1}^{n_k} \prod_{i=1}^{m_{kj}} \prod_{l=1}^{p_{ijk}}$$

$$[0.2 \ 0.1 \ 0.6 \ 0.1]$$

↳ E.g.: $\vec{y} = [0, 0, 1, 0] \in \{0,1\}^4 \quad y_4 = 1, y_{i \neq 4} = 0$

MLE IN DISCRIMINATIVE LEARNING

L3-9

$$\begin{aligned}
 \text{NLL}(\theta | D) &= \sum_{i=1}^m -\log P(Y_k | X^{(i)}; \theta) \quad \text{where } K = \text{ARGMAX } Y^{(i)} \\
 &= \sum_{i=1}^m \sum_{k=1}^K -y_k^{(i)} \log P(Y_k | X^{(i)}; \theta) \\
 &= \sum_{i=1}^m -\langle y_i^{(i)}, \log P(Y | X^{(i)}; \theta) \rangle \Rightarrow \langle Y, -\log P(Y | X; \theta) \rangle \\
 &= \sum_{i=1}^m E_{Y|X \sim P_0} [-\log P(Y | X^{(i)}; \theta)] \quad \leftarrow E_{Y|X \sim P_0} [-\log] \\
 &= \sum_{i=1}^m H(P_0, P(Y | X^{(i)}; \theta))
 \end{aligned}$$

↳ CROSS ENTROPY Between the observed
EMPIRICAL DIST Y AND our Model
DIST P_θ .

Softmax Regression:

$$f(x|\theta) \Rightarrow XW^T + b \Rightarrow \theta = \{W, b\}$$

$$\sigma(z) \Rightarrow \text{Softmax Function} \quad W \in \mathbb{R}^{K \times N} \quad b \in \mathbb{R}^K$$

$$\sigma(f(x; w, b))_k = \frac{e^{xw_k^T + b_k}}{\sum_{k'=1}^K e^{xw_{k'}^T + b_{k'}}}$$

SoftMax Regression

13-10

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

▷ Plug in σ_{softmax} :

$$NLL(w, b | D) = -\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log \frac{e^{x^{(i)} w_k^T + b_k}}{\sum_{k'=1}^K e^{x^{(i)} w_{k'}^T + b_{k'}}}$$

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

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

▷ GRADIENT DESCENT:

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

(2) UPDATE $\theta = \theta - \eta \nabla_{\theta}$

▷ $\nabla_{w_k} NLL(w, b | D) \in \mathbb{R}^N$

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

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

$$\Rightarrow \boxed{\nabla_{w_k} = \sum_{i=1}^m x^{(i)} (P(Y_k | X^{(i)}; w, b) - y_k^{(i)})} \Rightarrow \text{GRADIENT DESCENT}$$

$$\boxed{\nabla_{w_k} = x^{(i)} (P(Y_k | X^{(i)}; w, b) - y_k^{(i)})} \Rightarrow \text{STOCHASTIC GRADIENT DESCENT}$$

Softmax Regressor (Continued)

$$NLL(w, b | D) = - \sum_{i=1}^M Y_k^{(i)} - \frac{e^{x^{(i)} w_k^T + b_k}}{\sum_{k'=1}^K e^{x^{(i)} w_{k'}^T + b_{k'}}}$$

↳ $\nabla_{b_k} NLL(w, b | D) =$

$$\nabla_{b_k} = - \sum_{i=1}^M P(Y_k | x_k^{(i)}; w, b) - Y_k^{(i)}$$

\Rightarrow GRADIENT DESCENT

$$\nabla_{b_k} = P(Y_k | x_k^{(i)}; w, b) - Y_k^{(i)}$$

\Rightarrow STOCHASTIC GRADIENT DESCENT

↳ MINI-BATCH GRADIENT DESCENT:

$$1 \leq \text{Batch-size} < M$$