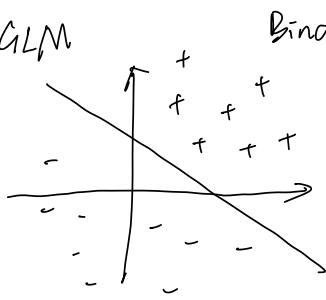


GLM



Binary classification

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$P(y \in \mathcal{Y}^+ | x) = \phi$$

$$P(y \in \mathcal{Y}^- | x) = 1 - \phi$$

$$\mathcal{Y}^+ = \{y | y > \text{Threshold}\}$$

$$\mathcal{Y}^- = \{y | y \leq \text{Threshold}\}$$

$$(y | x, \phi) \sim \text{Bernoulli}(\phi)$$

Exponential Family.

$$p(y; \eta) = b(y) e^{\eta^T T(y) - a(\eta)}$$

$$p(y; \phi) = \phi^y (1-\phi)^{1-y} = e^{y \log \phi + (1-y) \log(1-\phi)}$$

$$= e^{y(\log \phi + (1-y) \log(1-\phi))}$$

$$= e^{y(\log \phi - \log(1-\phi)) + \log(1-\phi)}$$

$$\text{Bernoulli distribution} = \frac{e^{y(\log \frac{\phi}{1-\phi}) - (-\log(1-\phi))}}{\log \frac{\phi}{1-\phi}}$$

$$b(y) = 1$$

$$T(y) = y$$

$$a(\eta) = -\log(1-\phi)$$

$$g(\eta) = \mathbb{E}[T(y) | \eta]$$

① canonical link function

$$\eta = g^{-1}(\mathbb{E}[T(y) | \eta])$$

② canonical response function

$$\eta = g^{-1}(\mathbb{E}[T(y) | \eta])$$

Generalized linear model

1.  $y | x; \theta \sim \text{Exponential Family}$

1.  $y | x; \phi \sim \text{Bernoulli}(\phi)$

2.  $h(x) = \mathbb{E}[T(y) | x]$

2.  $T(y) = y \quad h(x) = \mathbb{E}[y | x] = \phi$

3.  $\eta$  varies linearly with  $x$

3.  $\eta = \theta^T x$

$$\eta = \theta^T x$$

$$h(x) = \frac{1}{1 + e^{-\theta^T x}}$$

= sigmoid  $\eta$

$$= \frac{1}{1 + e^{-\eta}}$$

$\{ (x_i, y_i) \}_{i=1}^m$   
 $\uparrow$   $\rightarrow$  Label.  
 complex world

$y' = \underset{y}{\operatorname{argmax}} P(y|x) \rightarrow$  Discriminative Model.  
 $= \underset{y}{\operatorname{argmax}} P(x|y) \cdot P(y) \rightarrow$  Generative Model.  
 $\hookrightarrow$

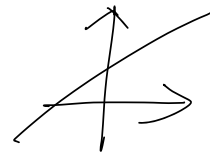
Basic knowledge.

• Discriminative Model:  $P(y|x; \theta)$

- Regression: Linear Regression  $y|x \sim \mathcal{N}(h_\theta(x), \sigma^2)$ .
- Binary class: Logistic Regression  $y|x \sim \operatorname{Bern}(h_\theta(x))$
- Multi-class ...: Softmax ...  $y|x \sim \operatorname{Multinomial}(h_\theta(x))$ .

▷ Linear Model: Given  $x$ ,  $y = \theta^T x + b$

▷ Given  $x$ ,  $h_\theta(x) = \begin{bmatrix} P(y=1|x) \\ \vdots \\ P(y=k|x) \end{bmatrix}$   
 $= \begin{bmatrix} 0.1 \\ 0.2 \\ 0.7 \end{bmatrix} \leftarrow$



• Generative Model:  $P(x, y) = P(x|y) \cdot P(y)$

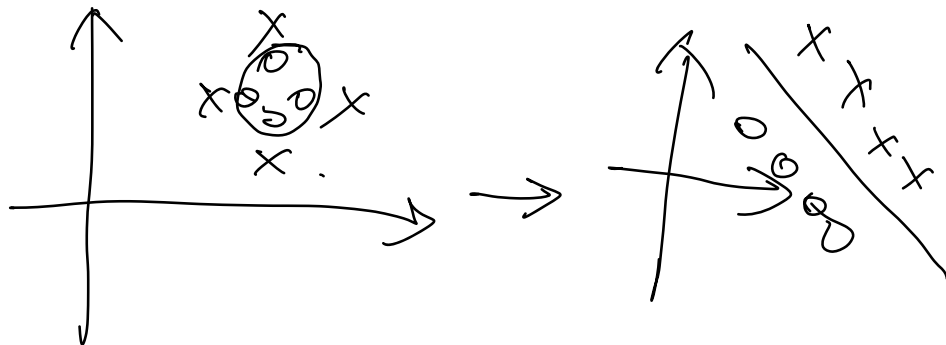
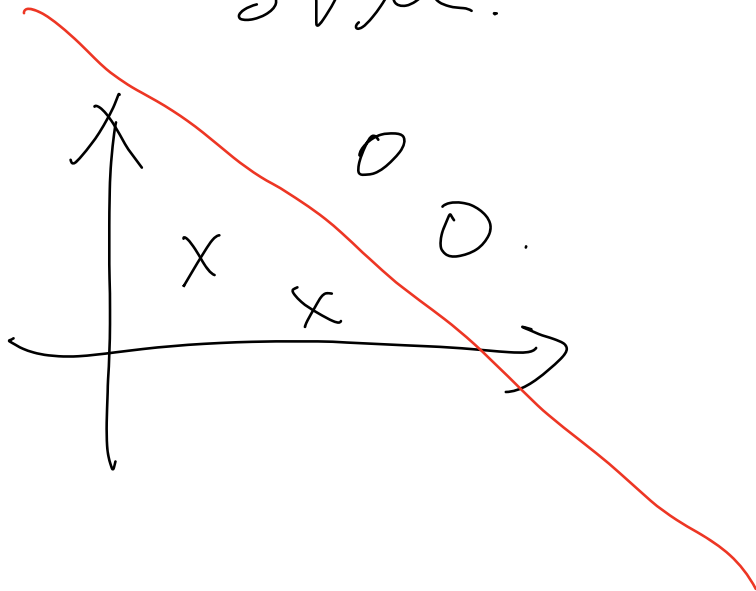
- continuous Input: GDA  $\begin{cases} y \sim \operatorname{Bern}(\phi) \\ x|y=b \sim \mathcal{N}(\mu_b, \Sigma) \end{cases}$
- Discrete Input: NB  $\begin{cases} P(x_1, \dots, x_n|y) = \prod_{i=1}^n P(x_i|y) \\ y \sim \operatorname{Bern}(\phi) \\ x_i|y=b \sim \operatorname{Bern}(\phi_{i|y=b}) \end{cases}$

▷ Bayes Rule.

$\uparrow$  Statistical

↓ Optimization.

SVM.



① Linear Regression.

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

$\epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$

LR = MLE + Linear Model + Gaussian Ass

Laplacian.

②  $\theta \sim \mathcal{N}(0, \frac{1}{2\lambda} I)$ .     $\text{MAP}$ .

18:41

③ Softmax function with

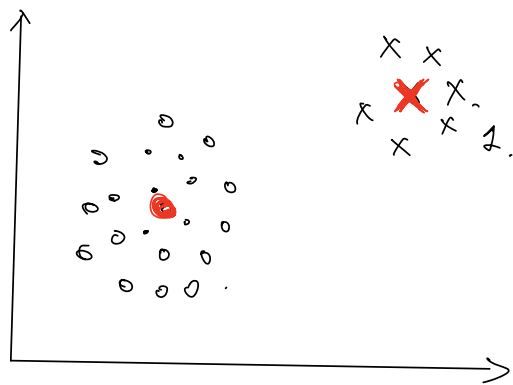
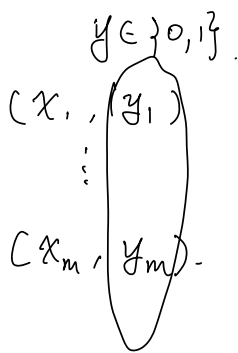
$k=2$

$$P_{Y|X}(1|x) = \frac{e^{\theta_1^T x}}{e^{\theta_1^T x} + e^{\theta_2^T x}} = \frac{1}{1 + e^{-10, -0.2} x}} = \underbrace{\sigma}_{\text{Sigmoid}}((\theta_1 - \theta_2)^T x)$$

$$\begin{aligned} \theta &= \operatorname{argmax}_{\theta} P\{\theta | \{ \} \} \\ &= \operatorname{argmax}_{\theta} P\{\{x, y\} | \theta\} \cdot P\{\theta\} \end{aligned}$$

Gaussian Discriminative Analysis (GDA)  $\rightarrow$  QDA

- Assume:  $y \sim \text{Bern}(\phi)$
- $x|y=0 \sim \mathcal{N}(\underline{\mu}_0, \underline{\Sigma}_0)$
- $x|y=1 \sim \mathcal{N}(\underline{\mu}_1, \underline{\Sigma}_1)$



$$\phi = \frac{\sum_{i=1}^m \mathbb{1}\{y_i=0\}}{m}$$

$$\mu_b = \frac{\sum_{i=1}^m \mathbb{1}\{y_i=b\} \cdot x_i}{\sum_{i=1}^m \mathbb{1}\{y_i=b\}}, \quad b=0, 1$$

$$\Sigma_b = \frac{\sum_{i=1}^m \mathbb{1}\{y_i=b\} (x^{(i)} - \mu_b)(x^{(i)} - \mu_b)^T}{\sum_{i=1}^m \mathbb{1}\{y_i=b\}}$$

$$\begin{aligned} \ell(\phi, \mu_0, \mu_1, \Sigma_0, \Sigma_1) &\triangleq \log \prod_{i=1}^m P(x^{(i)}, y^{(i)}) \\ &= \sum_{i=1}^m (\log P(x^{(i)} | y^{(i)}) + \log P(y^{(i)})). \quad \leftarrow \end{aligned}$$

$$\begin{aligned} \bullet y \sim \text{Bern}(\phi) &\Rightarrow P(y) = \phi^y (1-\phi)^{1-y}, \quad y=0,1. \\ &= \begin{cases} \phi^0 (1-\phi)^1, & y=0. \\ \phi^1 (1-\phi)^0, & y=1 \end{cases} \end{aligned}$$

$$\log P(y) = y \cdot \log \phi + (1-y) \log(1-\phi). \quad \leftarrow$$

$$\begin{aligned} \bullet x|y=0 \sim \mathcal{N}(\mu_0, \Sigma_0) &\Rightarrow P(x|y=0) = \left( (2\pi)^{\frac{n}{2}} |\Sigma_0|^{-\frac{1}{2}} \right)^{-1} \cdot \exp\left(-\frac{1}{2} (x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0)\right). \\ \log P(x|y=0) &= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_0| \\ &\quad - \frac{1}{2} (x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0). \end{aligned}$$

$$\begin{aligned} x|y=1 \sim \mathcal{N}(\mu_1, \Sigma_1) &\Rightarrow P(x|y=1) = \left( (2\pi)^{\frac{n}{2}} |\Sigma_1|^{-\frac{1}{2}} \right)^{-1} \cdot \exp\left(-\frac{1}{2} (x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1)\right). \\ \log P(x|y=1) &= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_1| \\ &\quad - \frac{1}{2} (x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1). \end{aligned}$$

eg.  $x \in \mathbb{R}^{n+1}$ ,  $\nabla_{\theta} \ell(\theta)$ ,  $\nabla_{\theta} \ell$

$$\begin{aligned} \bullet \ell &= \sum_{i=1}^m (\log P(x^{(i)} | y^{(i)}) + \log P(y^{(i)})) \\ &= \sum_{i=1}^m \left[ -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_{y^{(i)}}| - \frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^T \Sigma_{y^{(i)}}^{-1} (x^{(i)} - \mu_{y^{(i)}}) \right] \\ &\quad + \sum_{i=1}^m y^{(i)} \log \phi + (1-y^{(i)}) \log(1-\phi). \\ &= \sum_{i=1}^m \sum_{b=0}^1 \mathbb{1}\{y^{(i)}=b\} \left[ -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_b| - \frac{1}{2} (x^{(i)} - \mu_b)^T \Sigma_b^{-1} (x^{(i)} - \mu_b) \right] \\ &\quad + \sum_{i=1}^m y^{(i)} \log \phi + (1-y^{(i)}) \log(1-\phi). \end{aligned}$$

$$\bullet \frac{\partial l}{\partial \phi} = 0, \quad \frac{\partial l}{\partial \mu_b} = 0, \quad \frac{\partial l}{\partial \Sigma_b}$$

$$\frac{\partial a^T x}{\partial x} = a$$

$$\bullet \mu_0 = \frac{\sum_{i=1}^m \mathbb{1}\{y^{(i)}=0\}}{\sum_{i=1}^m \mathbb{1}\{y^{(i)}=0\} + \sum_{i=1}^m \mathbb{1}\{y^{(i)}=1\}} x^{(i)}$$

• Relation to Logistic Regression.



$$\bullet P(y=1|x) = \frac{1}{1 + e^{-a^T x}}$$

If  $p(x|y) \sim \mathcal{N}(\mu, \Sigma)$ .  $P(y|x)$  is a logistic function.

Naive Bayes.

• Dictionary

$$x = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 1 \end{bmatrix} \left. \begin{array}{l} a \\ \text{abandon} \\ \vdots \\ \text{Review} \\ \vdots \\ \text{Session} \end{array} \right\} \text{large!!}$$

$$x_1, x_2, \dots, x_n$$

$$P(x_1, \dots, x_n | y) = P(x_1 | y) \dots P(x_n | y) = \prod_{i=1}^n P(x_i | y)$$

$$y = \begin{cases} 0 & \text{no spam} \\ 1 & \text{spam!} \end{cases}$$

•  $y \sim \text{Bern}(\phi)$ .

$$x_i | y=b \sim \text{Bern}(\phi_{i|y=b}), \quad b=0,1.$$

$$\begin{aligned} & \mathcal{L}(\phi_y, \phi_{j|y=0}, \phi_{j|y=1}, j=1, \dots, n) \\ &= \sum_{i=1}^m \log(P(x^{(i)}, y^{(i)})) \\ &= \sum_{i=1}^m (\log(P(x^{(i)} | y^{(i)})) + \log(P(y^{(i)}))) \\ &= \sum_{i=1}^m \log \prod_{j=1}^n P(x_j^{(i)} | y^{(i)}) + \log P(y^{(i)}) \\ &= \sum_{i=1}^m \sum_{j=1}^n (\mathbb{1}\{y^{(i)}=0\} + \mathbb{1}\{y^{(i)}=1\}) \log P(x_j^{(i)} | y^{(i)}) + \log P(y^{(i)}) \end{aligned}$$

$$x^{(i)} = \begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$$

$$1 = \mathbb{1}\{y^{(i)}=0\} + \mathbb{1}\{y^{(i)}=1\}$$

$$\begin{aligned}
&= \sum_{i=1}^m \sum_{j=1}^n \sum_{b=0}^1 \mathbb{1}\{y^{(i)}=b\} \cdot \log P(x_j^{(i)} | y^{(i)}) + \sum_{i=1}^m y^{(i)} \log(\phi_y) + (1-y^{(i)}) \log(1-\phi_y) \\
&= \sum_{i=1}^m \sum_{j=1}^n \sum_{b=0}^1 \mathbb{1}\{y^{(i)}=b\} (x_j^{(i)} \log \phi_{j|y=b} + (1-x_j^{(i)}) \log(1-\phi_{j|y=b})) \\
&\quad + \sum_{i=1}^m y^{(i)} \log(\phi_y) + (1-y^{(i)}) \log(1-\phi_y)
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \phi_y} &= \frac{\partial}{\partial \phi_y} \left( \sum_{i=1}^m y^{(i)} \log(\phi_y) + (1-y^{(i)}) \log(1-\phi_y) \right) \\
&= \frac{\sum_{i=1}^m y^{(i)}}{\phi_y} - \frac{\sum_{i=1}^m (1-y^{(i)})}{1-\phi_y} = 0. \quad \frac{a}{b} = \frac{c}{d} = \frac{a+c}{b+d}
\end{aligned}$$

$$\Rightarrow \phi_y = \frac{1}{m} \cdot \mathbb{1}\{y^{(i)}=y\}, \quad y=0,1.$$

Support Vector Machine - (SVM)

(KKT)

• Given  $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$ , find  $(w^*, b^*)$  to.

$$\begin{aligned}
\text{prime problem} & \quad \min_{w,b} \frac{1}{2} \|w\|^2 \\
& \quad \text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq 1, \quad i=1, \dots, m. \\
\text{dual problem} & \quad \max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \langle \alpha_i y^{(i)} x^{(i)}, \alpha_j y^{(j)} x^{(j)} \rangle \\
& \quad \text{s.t. } \alpha_i \geq 0, \quad i=1, \dots, m. \\
& \quad \sum_{i=1}^m \alpha_i y^{(i)} = 0.
\end{aligned}$$

$$\text{Solution: } \begin{cases} w^* = \sum_i \alpha_i^* y^{(i)} x^{(i)} \\ b^* = -\frac{1}{2} \left( \max_{i: y^{(i)}=-1} w^{*T} x^{(i)} + \min_{i: y^{(i)}=1} w^{*T} x^{(i)} \right). \end{cases}$$

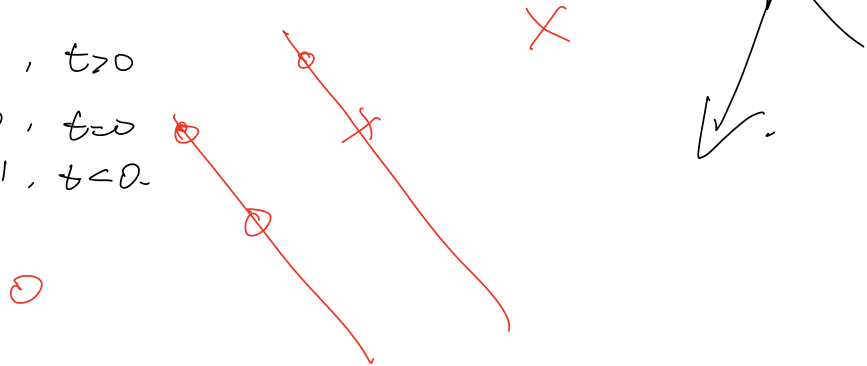
Given new sample  $z$ .

$$y = \text{sign}[\underbrace{w^{*T} z + b^*}]$$



$$= \text{sign} \left[ \sum_{i=1}^m \alpha_i y^{(i)} \langle x^{(i)}, z \rangle + b \right]$$

$$\text{sign}(t) = \begin{cases} 1, & t > 0 \\ 0, & t = 0 \\ -1, & t < 0 \end{cases}$$

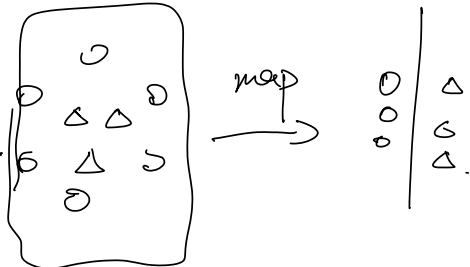


### \* Soft-SVM.

• kernel Trick.

$$x \in \mathbb{R}^d \mapsto \phi(x) \in \mathbb{R}^D, D \gg d.$$

$$K(x, x') \triangleq \phi^T(x) \phi(x') \in \mathbb{R}.$$



$$K \begin{matrix} \square \\ \square \\ \square \\ \square \end{matrix} \quad \# \{x\} \times \# \{x'\}.$$

$$\text{Dual : } \max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j K(x^{(i)}, x^{(j)}) \quad \left. \vphantom{\max_{\alpha}} \right\} \Rightarrow \alpha^*$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, \quad \sum \alpha_i y^{(i)} = 0$$

$$\text{Solution : } w^* = \sum \alpha_i^* y^{(i)} \phi(x^{(i)})$$

$$b^* = y^{(j)} - \sum_{i=1}^m \alpha_i^* y^{(i)} K(x^{(i)}, x^{(j)}), \quad \text{for some } 0 < \alpha_j < C.$$

$$f(w) = \text{sign}(w^T \phi(x) + b).$$



$$P_{xy}(x, y)$$

$$P_x(y) = \sum_x P_{xy}(x, y)$$

$$P_x(x) = \sum_y P_{xy}(x, y)$$

$$P_{x|y}$$

$$P_{y|x}$$