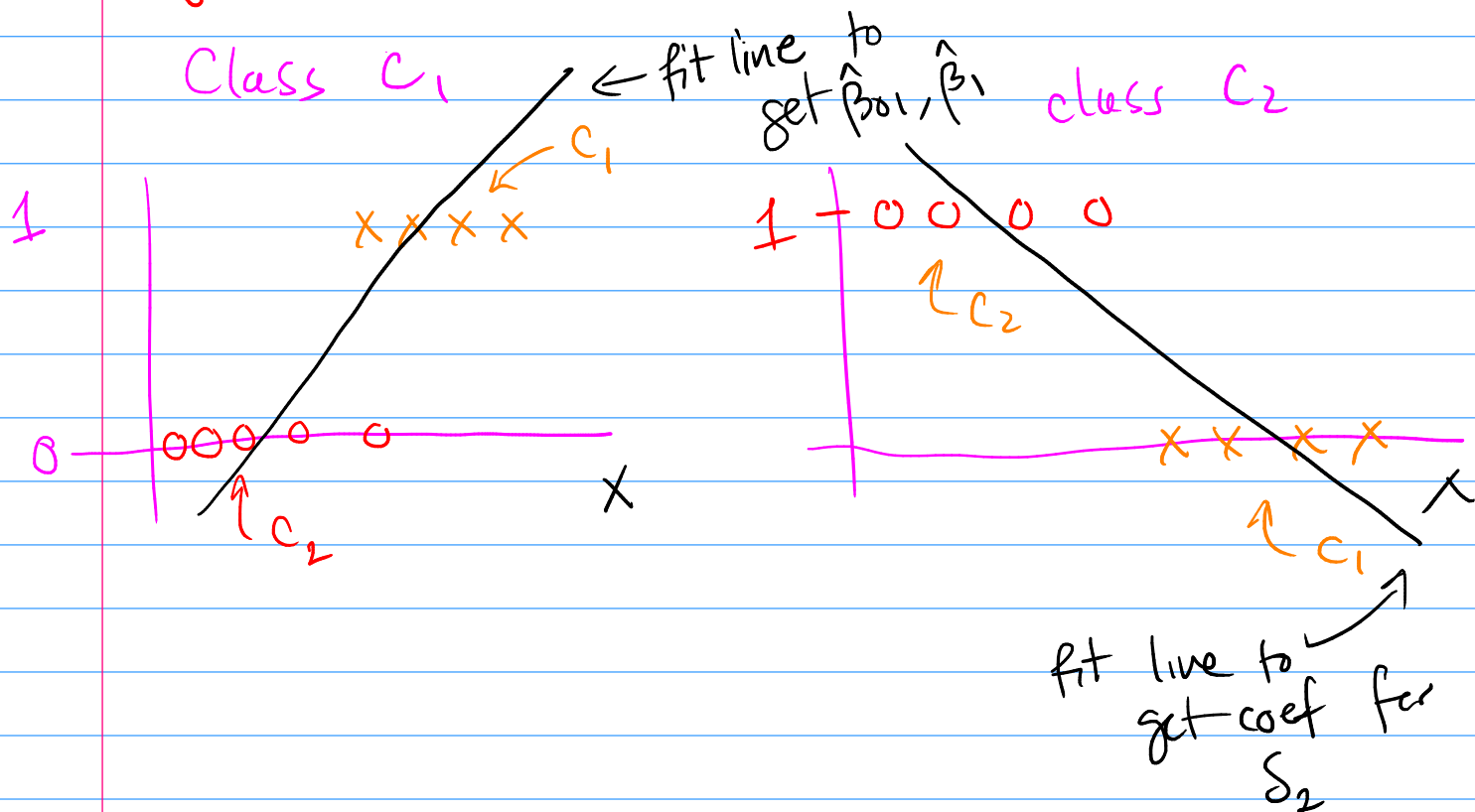# Lecture 11: Logistic Regression

LDA is a linear classifier

$$\delta_c(\underline{x}) = \hat{\beta}_{0c} + \hat{\beta}_c^{\top} \underline{x}$$

Why not just fit $\delta_c$ using linear regression?
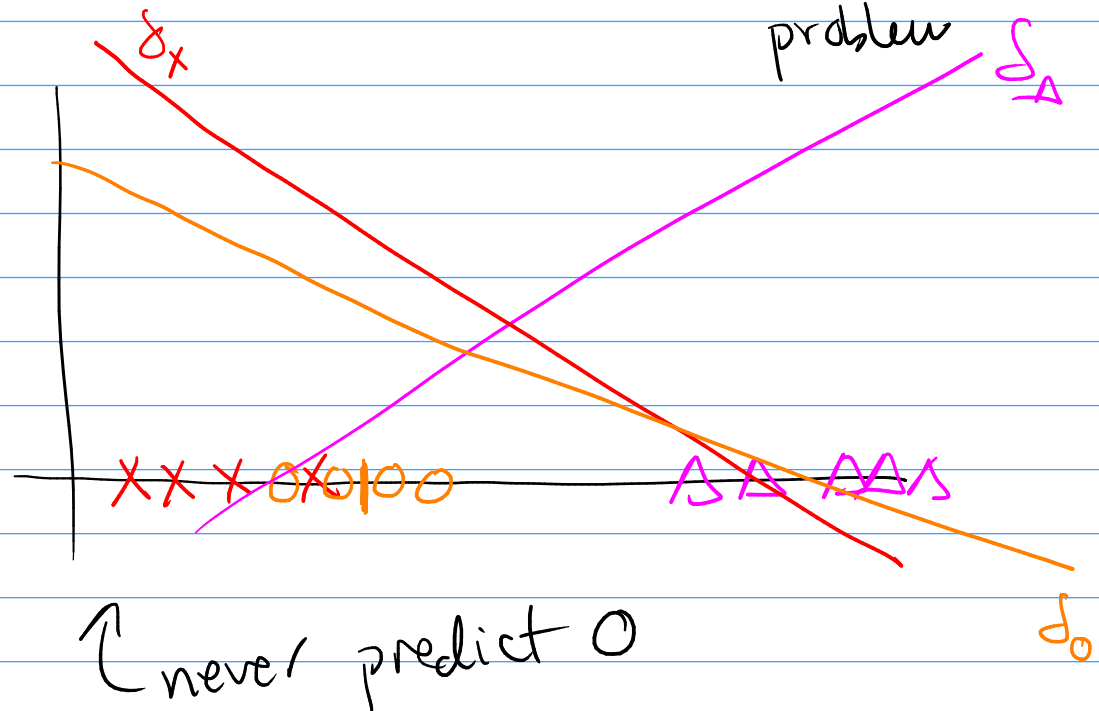
Binary example $Y = c_1$ or $Y = c_2$

Class $c_1$ ← fit line to get $\hat{\beta}_{01}, \hat{\beta}_1$  class $c_2$

$C_1$

1    × × × ×          1 ┬ o o  o  o

0 ─ o o o o  o    o              × × × ×  ×

$C_2$                                                        $C_1$

$C_2$                                        fit line to get coef for $\delta_2$

Punchline: — reasonable when $K = $ # classes is small

$(K = 2$ exactly gives LDA$)$

When K is "large" can have a <u>masking</u>
<u>problem</u>

$K=3$

$\delta_+$

$\delta_\triangle$

X X X OXOIOO    $\triangle \triangle \quad \triangle\triangle\triangle \triangle$

$\delta_0$

$\Big\{$ never predict O

---

## <u>Logistic Regression</u>

<u>LDA:</u>  $\delta_c(\underline{x}) = \mathbb{P}(Y=c \mid \underline{X}=\underline{x})$

$\qquad\qquad \propto \mathbb{P}(\underline{X}=\underline{x} \mid Y=c)\,\mathbb{P}(Y=c)$

## <u>Logistic Reg.</u>

directly model $\delta_c(\underline{x}) = \mathbb{P}(Y=c \mid \underline{X}=\underline{x})$

---

<u>Binary Classification ($K=2$)</u>

So $Y=0$ or $Y=1$    $(Y=\pm 1)$

$\delta_0(\underline{x}) = \mathbb{P}(Y=0 \mid \underline{X}=\underline{x})$

$\delta_1(\underline{x}) = \mathbb{P}(Y=1 \mid \underline{X}=\underline{x}) = 1 - \mathbb{P}(Y=0 \mid \underline{X}=\underline{x}) = 1 - \delta_0(\underline{x})$

So I really only need $\delta_1$

$$\hat{f}(\underline{x}) = \arg\max_c \delta_c(\underline{x})$$

$$\hat{f}(\underline{x}) = 1 \quad \text{when} \quad \delta_1(\underline{x}) > \delta_0(\underline{x})$$
$$> 1 - \delta_1(\underline{x})$$

so
$$\hat{f}(\underline{x}) = 1 \quad \text{when} \quad \delta_1(\underline{x}) > 1/2$$

Traditionally $\quad p(\underline{x}) = \delta_1(\underline{x}) = P(Y=1 \mid \underline{X}=\underline{x})$

Given $\underline{X} = \underline{x}$, $Y = 0$ or $Y = 1$

$$Y \mid \underline{X} = \underline{x} \sim \text{Bern}(p(\underline{x}))$$
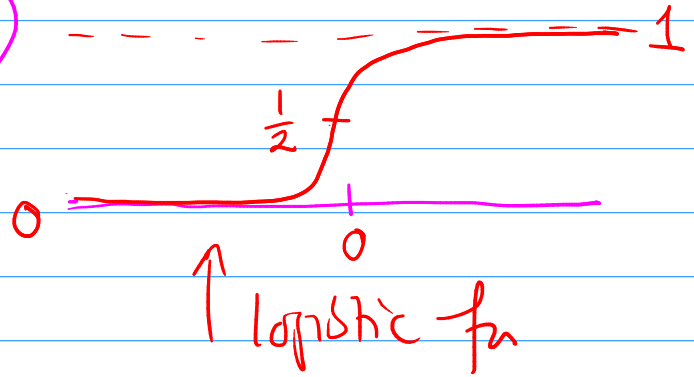
<u>Game</u>: reasonable model for $p(\underline{x})$



linear? $p(\underline{x}) = \underline{x}^T \beta$

S-shaped curve

# Logistic regression Says!

$$p(\underset{\sim}{x}) = \text{logistic}(\underset{\sim}{x}^T\hat{\beta})$$

$$\text{logistic}(x) = 1/(1+e^{-x})$$

**Inverse:**

$$\text{logit}(x) = \text{logistic}^{-1}(x)$$

$$= \log(x/(1-x))$$


↑ logistic fn

$$p(\underset{\sim}{x}) = \text{logistic}(\underset{\sim}{x}^T\hat{\beta})$$

$$= \frac{1}{1+\exp(-\underset{\sim}{x}^T\hat{\beta})}$$

$$= \frac{1}{1+\exp\left(-\left(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \cdots + \hat{\beta}_p x_p\right)\right)}$$

**Notice!** $\quad \delta_1(\underset{\sim}{x}) = P(\underset{\sim}{x}) = \text{logistic}(\underset{\sim}{x}^T\hat{\beta})$

So

$$\text{logit}(\delta_1(\underset{\sim}{x})) = \underset{\sim}{x}^T\hat{\beta}$$

← linear
a linear classifier

How do we get $\hat{\beta}$?

Training data: $Y_n | \underset{\sim}{X}_n = \underset{\sim}{x}_n \overset{indep}{\sim} Bern(P_\beta(\underset{\sim}{x}_n))$

$$P_\beta(\underset{\sim}{x}) = \frac{1}{1 + \exp(-\underset{\sim}{x}^T \beta)}$$

Maximum Likelihood Est. (MLEs)

parameter $\theta$ in a model $P_\theta(\underset{\sim}{x})$

$$\hat{\theta}_{MLE} = \underset{\theta}{argmax} \; P_\theta(\underset{\sim}{x})$$

Logistic regression sets $\hat{\beta}$ as MLE of $\beta$ under this model.

$$\hat{\beta} = \underset{\beta}{argmax} \; \underbrace{P(Y_1, Y_2, Y_2, ..., Y_N | \underset{\sim}{X}_1, \underset{\sim}{X}_2, ..., \underset{\sim}{X}_N)}_{L(\beta)}$$

$$L(\beta) = \prod_{n=1}^{N} P_\beta(Y_n | \underset{\sim}{X}_n = \underset{\sim}{x}_n)$$

$$= \prod_{n=1}^{N} P_\beta(\underset{\sim}{x}_n)^{y_n} (1 - P_\beta(\underset{\sim}{x}_n))^{1-y_n}$$

$$= \prod_{n=1}^{N} \left( \frac{1}{1+\exp(-\underset{\sim}{x}_n^T \beta)} \right)^{y_n} \left( 1 - \frac{1}{1+\exp(-\underset{\sim}{x}_n^T \beta)} \right)^{1-y_n}$$

Bernoulli (p)

$$p(x) = p^x (1-p)^{1-x}$$

for $x = 0$ or $1$

No closed form soln, need numerical opt to get $\hat{\beta}$.

---

## Multi-nomial Logistic Regression

$K-1$ probs

When $K > 2$

$$Y_n \mid X_n = x_n \overset{indep}{\sim} Categorical(P_1(x_n), \ldots, P_{K-1}(x_n))$$

$$\delta_k(x) = \mathbb{P}(Y = k \mid x = x) = P_k(x)$$

$$= \text{multi-variate logistic} (x^T \hat{\beta}_k)$$

soft-max function

$$= \frac{\exp(\hat{\beta}_k^T x)}{1 + \sum_{k=1}^{K-1} \exp(\hat{\beta}_k^T x)}$$

Similarly, fit each/all $\hat{\beta}_k s$ by MLE.

---

## LDA v. Logistic Regression

| LDA | Logistic Regression |
|---|---|
| ① models $X/Y$ and $Y$ using normality assumption on $X$ | ① models $Y/X$ no model of $X$ |
| ② easier to fit | ② harder to fit |
| ③ both linear | |