LDA: $\delta_c(x) = \mathbb{P}(Y = c \mid X = x)$

$$= \frac{\mathbb{P}(X = x \mid Y = c)\,\mathbb{P}(Y = c)}{\mathbb{P}(X = x)}$$

$$\propto \underbrace{\mathbb{P}(X = x \mid Y = c)}_{N(\mu_c,\, \sigma^2)}\, \underbrace{\mathbb{P}(Y = c)}_{\pi_c}$$

## Logistic Regression

$$\delta_c(x) = \mathbb{P}(Y = c \mid X = x) \quad \leftarrow \text{model this directly.}$$

## Binary Logistic Regression $(K = 2)$

So here $Y = 0$ or $Y = 1$

discr. fns

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

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

So in the binary case only need one discr function (we can always get the other)

Call $p(x) = \delta_1(x) = P(Y=1|X=x)$

Rule: classify as class 1 if $p(x) > \frac{1}{2}$

since iff $p(x) = \delta_1(x) > \frac{1}{2}$ since $\delta_1$ and $\delta_0$
sum to 1 then $\delta_0(x) < \frac{1}{2}$,
hence $\delta_1(x) > \delta_0(x)$

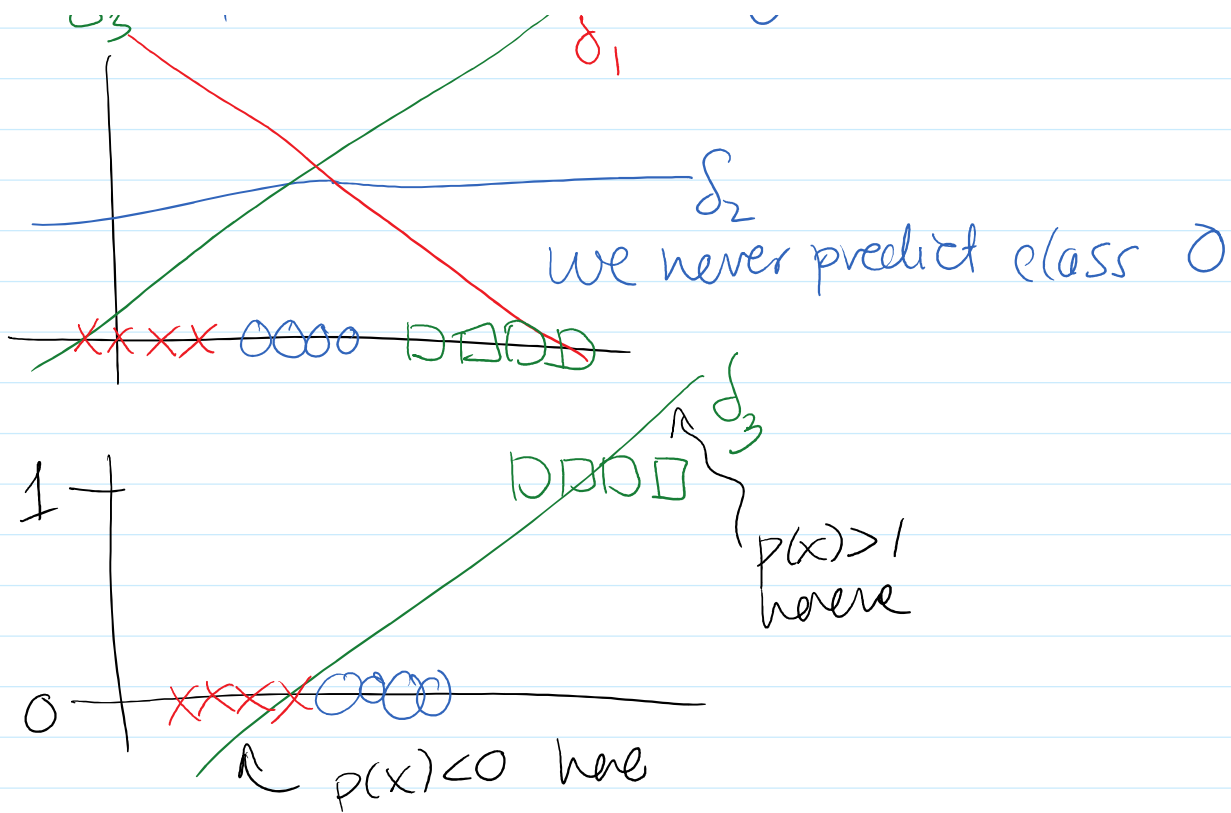Given $X=x$ notice that $Y=0$ or $Y=1$.

We call $Y/X=x$ a Bernoulli R.V.

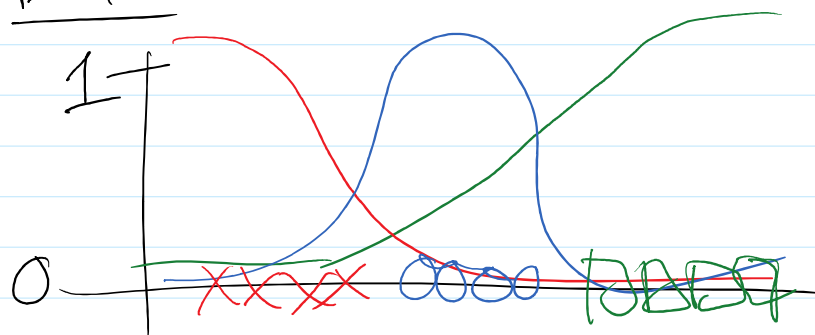We could write our setup simply as

$$Y/X=x \sim Bernoulli(p(x))$$

Game: How do we model $p(x)$?

Way 1: $p(x) = \hat{\beta}^T x$ could learn $\hat{\beta}$s from
regression.

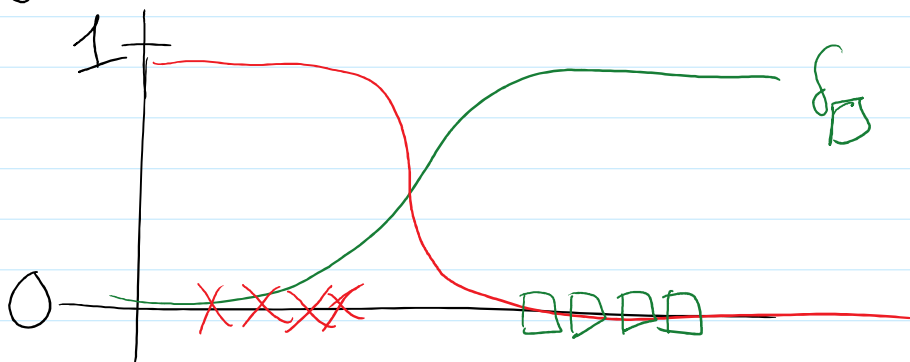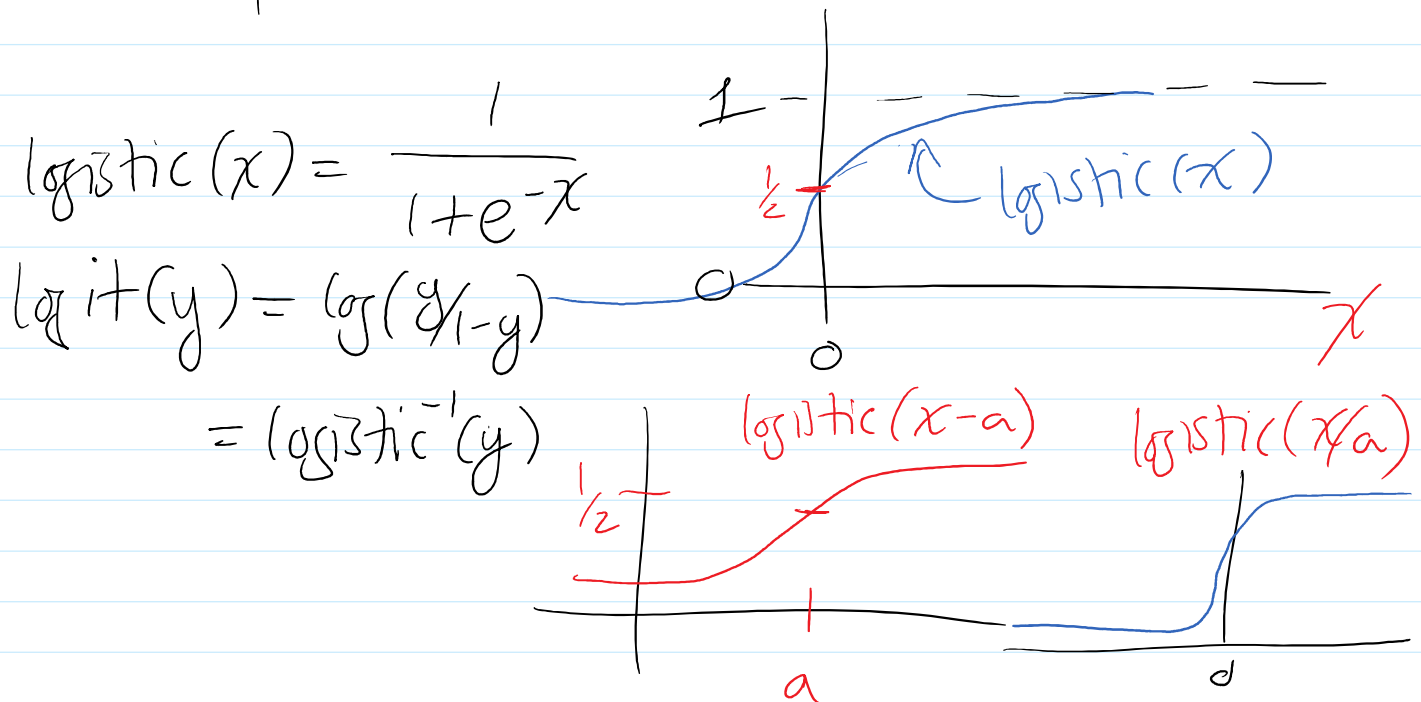problems: masking

$\delta_3$                     $\delta_1$

$\delta_3$      $\delta_1$

$\delta_2$

we never predict class 0

xxxx oooo □□□□

$\delta_3$

□□□□

$p(x) \gg 1$ here

1

0

xxxx ooooo

↑ $p(x) < 0$ here

## Better

1

0

xxxx ooo □□□□

## Way 2: Logistic Regression (Binary)

1

0

$\delta_{□}$

xxxx □□□□

Logistic regression says let

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

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

$$\text{logit}(y) = \log(y/(1-y))$$

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

logistic(x)

logistic(x-a)

logistic(x/a)

---

logistic regression (in Binary Case)

$$Y | X = x \sim \text{Bernoulli}(p(x))$$

$$\text{and} \quad p(x) = \frac{1}{1 + e^{-\beta^T x}}$$

How do we learn $\hat{\beta}$?  Maximum Likelihood Estimation

$$y_n | x_n \overset{\text{indep}}{\sim} \text{Bernoulli}(p_\beta(x))$$

$$p_\beta(x) = \text{logistic}(\beta^T x).$$

Joint density of $y_n$s conditioned on $x_n$s

liklihood

$$L(\beta) = \mathbb{P}(y_1, y_2, y_3, \dots, y_N \mid x_1, x_2, x_3, \dots, x_N)$$

$$= \prod_{n=1}^{N} \mathbb{P}(y_n \mid x_n)$$

$$= \prod_{n=1}^{N} p_\beta(x)^y (1 - p_\beta(x))^{1-y}$$

$$= \prod_{n=1}^{N} \left(\frac{1}{1+e^{-\beta^T x_n}}\right)^{y_n} \left(1 - \frac{1}{1+e^{-\beta^T x_n}}\right)^{1-y_n}$$

$$Y \sim \text{Bernoulli}(p)$$
$$f(x) = p^y (1-p)^{1-y}$$
$$= \begin{cases} p, & 1 \\ 1-p, & 0 \end{cases}$$

## MLE!

$$\hat{\beta} = \arg\max_\beta L(\beta)$$
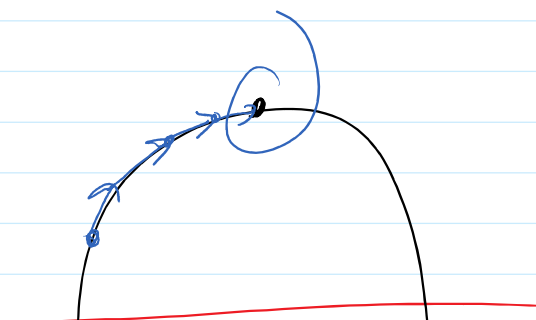
→ No analytical solution — give to friends in OR.

→ Solve w/ gradient descent (iterative)
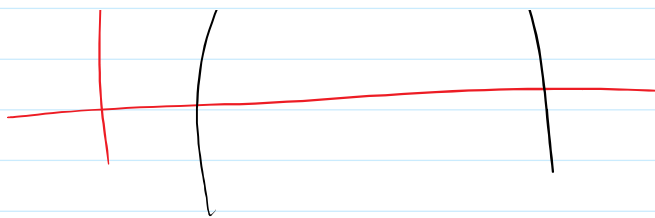   Can cast GD as
   a series of (weighted)
   regression.
   (IRLS)
   ↳ iteratively

⌐ iteratnenely
reweighted
least-squares

---

Multinomial Logistic Regression

## Multi-Class $(K > 2)$

generalizata of Binomial
when I have $K$
discrete outcomes

$$y_n \mid x_n \overset{indep}{\sim} \text{Multinomial}(P_1(x), P_2(x), \cdots, P_{K-1}(x))$$

Sum to 1 so
only need $K-1$

For $k = 1, \cdots, K-1$

$$\delta_k(x) = P_k(x) = P(Y = k \mid X = x)$$

$$= MV\text{Logistic}_k(\cdots)$$

$$= \frac{e^{\beta_k^T x}}{1 + \sum_{\ell=1}^{K} e^{\beta_\ell^T x}} \longleftarrow \beta \text{ for each class.}$$

Now we have $\beta_1, \beta_2, \cdots, \beta_K$
and we estimate these as MLEs

$$\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \cdots, \hat{\beta}_K = \arg\max_{\beta_1 \cdots \beta_K} L(\beta_1, \cdots, \beta_K)$$
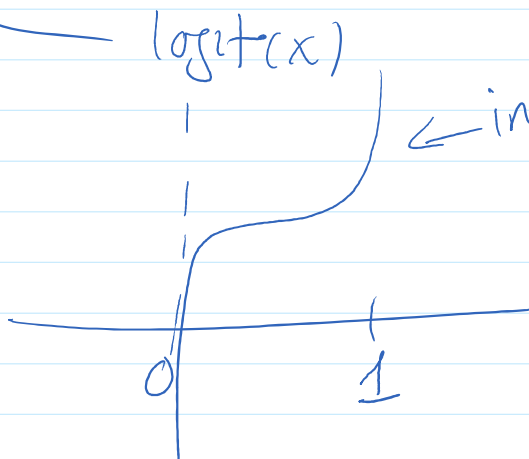
---

## Back to Binary $(K = 2)$

$$\delta(x) = P(x) = \frac{1}{\phantom{---}} = \text{logistic}(\beta^T x)$$

$$1 + e^{-\beta^T x}$$

$$\text{logit } \delta_1(x) = \log\left(\frac{P(x)}{1-P(x)}\right) = \beta^T x$$

$$1 - P(x) = \frac{1 + e^{-\beta^T x} - 1}{1 + e^{-\beta^T x}} = \frac{e^{-\beta^T x}}{1 + e^{-\beta^T x}}$$

$$P(x)\Big/1-P(x) = \frac{1 + e^{-\beta^T x}}{e^{-\beta^T x}} \cdot \frac{1}{1 + e^{-\beta^T x}} = e^{\beta^T x}$$

logit(x)

← increasy

0          1

$\Rightarrow$ increesy transf. of my $\delta_1$ is linear

$\Rightarrow$ logistic Regression is a linear classifier.

## LDA v. Logistic Regression

| LDA | Logistic Reg |
|---|---|
| ① models $P(X\mid Y)$ and $P(Y)$ using a normality | ① models $Y\mid X$ directly — no normality assumpta about X (more cannon  ) |

models "(MLE) vs "LE")
Using a **normality**
assumption about X|Y
(X continous)

② normality is
make estimates
easier to get
(more efficient if
Xs truly normal)

no normality assumption about
X (more general)
(can have categorical covariates)

② estimating βs is
comp. more expensive.

Both $\delta_1(x) = \hat{\beta}^T x$