



Statistics and Data Science for Engineers

E178 / ME276DS

Logistic regression

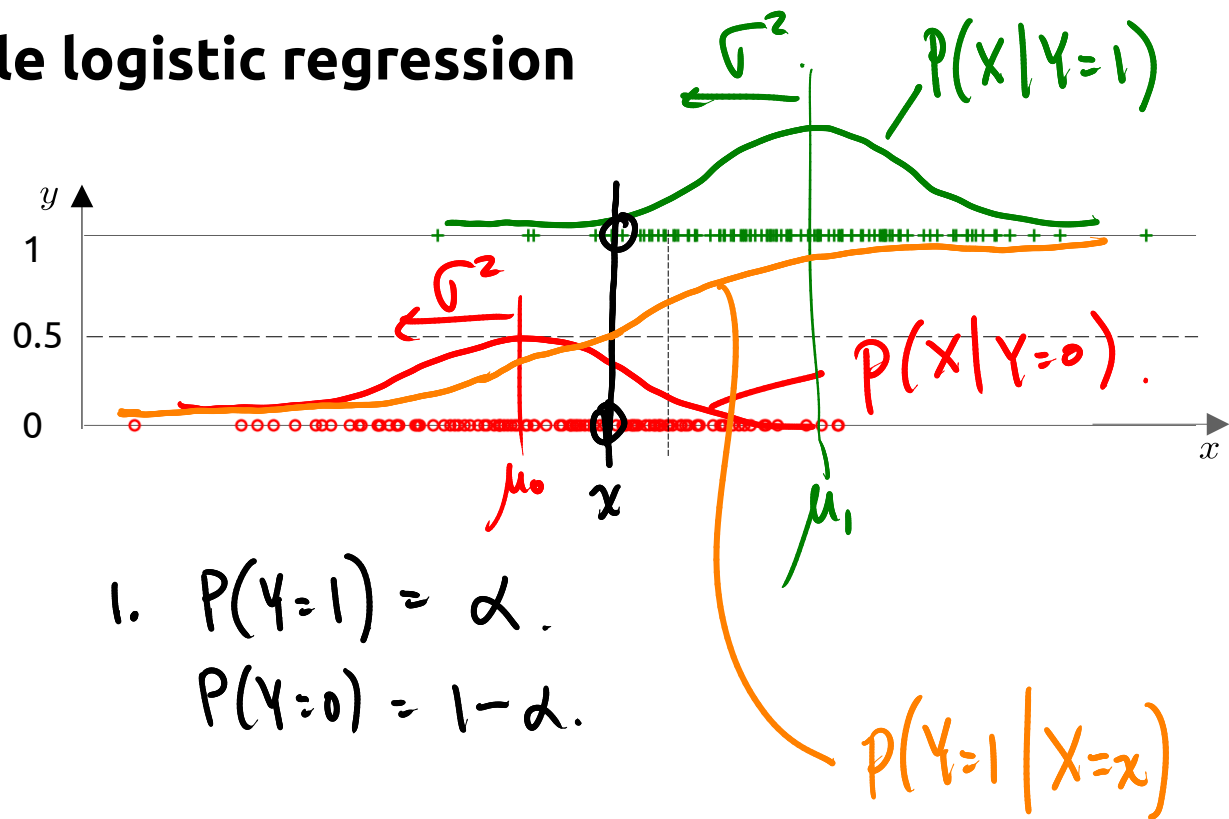
Naïve Bayes: Used assumptions to approximate $p(Y=c) \prod_{d=1}^D p(X^d=x^d | Y=c)$

Simple logistic regression

Assumptions:

$K=2$.

1. $Y \sim \mathcal{B}(\alpha)$, α is known
2. $D=1$, single input.
3. $X | Y=0 \sim \mathcal{N}(\mu_0, \underline{\sigma^2})$
4. $X | Y=1 \sim \mathcal{N}(\mu_1, \underline{\sigma^2})$
5. μ_0, μ_1, σ^2 are known



$$\hat{y} = h(x) = \operatorname{argmax}_{c \in \{0,1\}} p(Y=c \mid X=x)$$

$$h(x) = \begin{cases} 1 & p(Y=1 \mid X=x) > p(Y=0 \mid X=x) \\ 0 & \text{otherwise} \end{cases}$$

$$= \begin{cases} 1 & \frac{p(Y=1 \mid X=x)}{p(Y=0 \mid X=x)} > 1 \\ 0 & \text{otherwise} \end{cases} \quad \dots \text{ odds ratio } \checkmark$$

$$= \begin{cases} 1 & \log \left(\frac{p(Y=1 \mid X=x)}{p(Y=0 \mid X=x)} \right) > 0 \\ 0 & \text{otherwise} \end{cases} \quad \dots \text{ logs odds ratio } \checkmark$$

Bayes' rule

$$\log \left(\frac{p(Y=1|X=x)}{p(Y=0|X=x)} \right) = \log \left(\frac{\frac{p(Y=1)p(X=x|Y=1)}{\cancel{p(X=x)}}}{\frac{p(Y=0)p(X=x|Y=0)}{\cancel{p(X=x)}}} \right)$$

$$= \log \left(\frac{p(Y=1)}{p(Y=0)} \frac{p(X=x|Y=1)}{p(X=x|Y=0)} \right)$$

$$= \log \left(\frac{\alpha}{1-\alpha} \frac{\cancel{\frac{1}{\sqrt{2\pi\sigma^2}}}}{\cancel{\frac{1}{\sqrt{2\pi\sigma^2}}}} \frac{\exp \left(-\frac{1}{2} \frac{(x-\mu_1)^2}{\sigma^2} \right)}{\exp \left(-\frac{1}{2} \frac{(x-\mu_0)^2}{\sigma^2} \right)} \right)$$

⋮

$$\begin{aligned}
& \vdots \\
&= \log \left(\frac{\alpha}{1-\alpha} \right) - \frac{1}{2\sigma^2} \left((x - \mu_1)^2 - (x - \mu_0)^2 \right) \\
&= \log \left(\frac{\alpha}{1-\alpha} \right) - \frac{1}{2\sigma^2} \left(-2(\mu_1 - \mu_0)x + (\mu_1^2 - \mu_0^2) \right) \\
&= \underbrace{\log \left(\frac{\alpha}{1-\alpha} \right) + \frac{\mu_0^2 - \mu_1^2}{2\sigma^2}}_{\hat{\theta}_0} + x \underbrace{\frac{\mu_1 - \mu_0}{\sigma^2}}_{\hat{\theta}_1} \\
&= \hat{\theta}_0 + x \hat{\theta}_1
\end{aligned}$$

$$\therefore h(x) = \begin{cases} 1 & \hat{\theta}_0 + x \hat{\theta}_1 > 0 \\ 0 & \text{otherwise} \end{cases} \quad \tau: -\hat{\theta}_0/\hat{\theta}_1$$




$\tau = -\hat{\theta}_0 / \hat{\theta}_1$

$h(x) = 0 \leftarrow$ $h(x) = 1 \rightarrow$

Compute the conditional probabilities in terms of $\hat{\theta}_0$ and $\hat{\theta}_1$

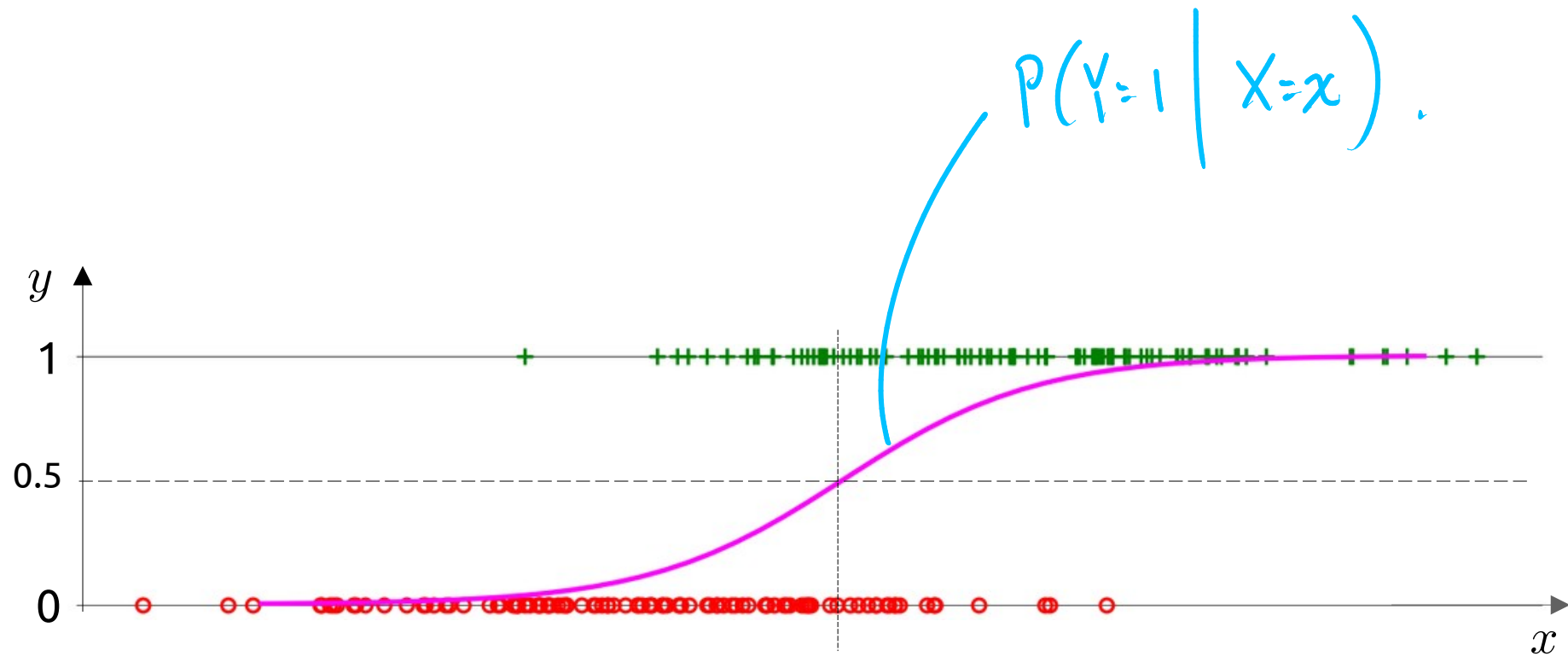
$$\ln \left(\frac{p(Y=1|X=x)}{p(Y=0|X=x)} \right) = \hat{\theta}_0 + \hat{\theta}_1 x$$

$$\Leftrightarrow \frac{p(Y=1|X=x)}{1 - p(Y=1|X=x)} = \exp(\hat{\theta}_0 + \hat{\theta}_1 x)$$

$$\Leftrightarrow p(Y=1|X=x) = \frac{1}{1 + \exp(-(\hat{\theta}_0 + \hat{\theta}_1 x))}$$


Define the **sigmoid** function: $\sigma(z) = \frac{1}{1 + e^{-z}}$ $z = \hat{\theta}_0 + \hat{\theta}_1 x.$

Then $p(Y=1|X=x) = \sigma(\hat{\theta}_0 + \hat{\theta}_1 x)$



Full logistic regression

$$x = \begin{bmatrix} x' \\ \vdots \\ x^D \end{bmatrix}$$

$$\underline{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_D \end{bmatrix}$$

Assumptions:

$K=2$

1. $Y \sim \mathcal{B}(\alpha)$ ~~α is known~~

✓ 2. ~~$D=1$~~ ,

3. ~~$X | Y=0 \sim \mathcal{N}(\mu_0, \sigma^2)$~~

4. ~~$X | Y=1 \sim \mathcal{N}(\mu_1, \sigma^2)$~~

5. ~~μ_0, μ_1, σ^2 are known~~

New assumption:

$$\begin{aligned} p(Y=1|X=x) &= \sigma(\theta_0 + x^1\theta_1 + \dots + x^D\theta_D) \\ &= \sigma(\theta_0 + x^T \underline{\theta}_1) \end{aligned}$$

for some $\theta_0, \theta_1, \dots, \theta_D$

$$P(Y=0|X=x) = 1 - \sigma(\theta_0 + x^T \underline{\theta}_1).$$

MLE solution to full logistic regression

$$\begin{aligned}\mathcal{L}(\theta_0, \underline{\theta}_1) &= \prod_{i=1}^N p(y_i | x_i; \theta_0, \underline{\theta}_1) \\ &= \prod_{\{i: y_i=1\}} \sigma(\theta_0 + x_i^T \underline{\theta}_1) \prod_{\{i: y_i=0\}} (1 - \sigma(\theta_0 + x_i^T \underline{\theta}_1))\end{aligned}$$

$$(\hat{\theta}_0, \hat{\theta}_1) = \underset{\theta_0, \theta_1}{\operatorname{argmax}} \left(\prod_{\{i: y_i=1\}} \sigma(\theta_0 + x_i^T \theta_1) \prod_{\{i: y_i=0\}} (1 - \sigma(\theta_0 + x_i^T \theta_1)) \right)$$

$$= \underset{\theta_0, \theta_1}{\operatorname{argmax}} \prod_{i=1}^N \sigma(\theta_0 + x_i^T \theta_1)^{y_i} (1 - \sigma(\theta_0 + x_i^T \theta_1))^{(1-y_i)}$$

$$= \underset{\theta_0, \theta_1}{\operatorname{argmax}} \sum_{i=1}^N \left(y_i \log \sigma(\theta_0 + x_i^T \theta_1) + (1 - y_i) \log (1 - \sigma(\theta_0 + x_i^T \theta_1)) \right)$$

$\times (-1)$

$$= \underset{\theta_0, \theta_1}{\operatorname{argmin}} \sum_{i=1}^N \text{CE}(y_i, \hat{p}_i)$$

cross entropy.

True output
estimated $P(Y=1 | X=x)$.

where $\hat{p}_i = \sigma(\theta_0 + x_i^T \theta_1) \approx P(Y=1 | X=x)$.

→ $\text{CE}(y, p) = -y \log(p) - (1 - y) \log(1 - p)$

clever
trick

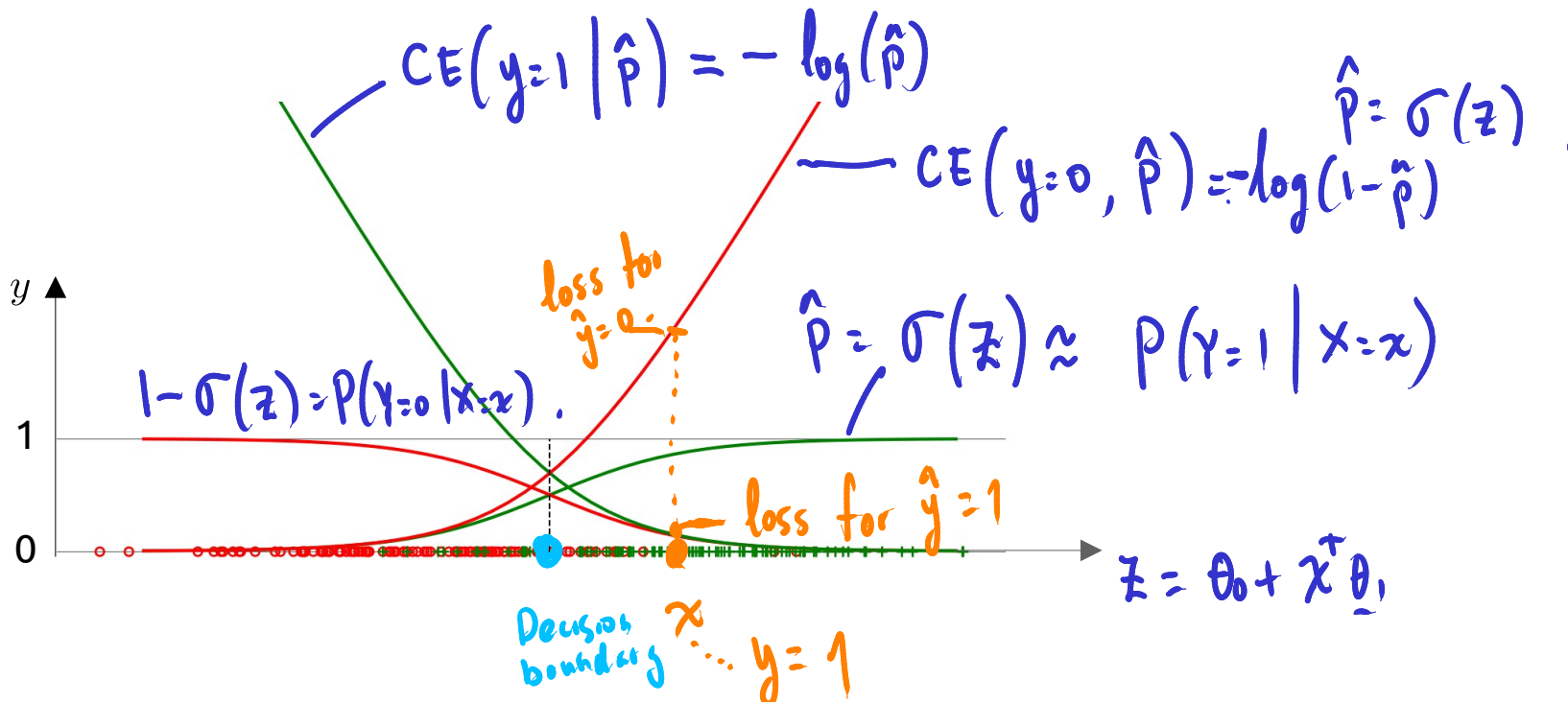
log.

Cross entropy loss

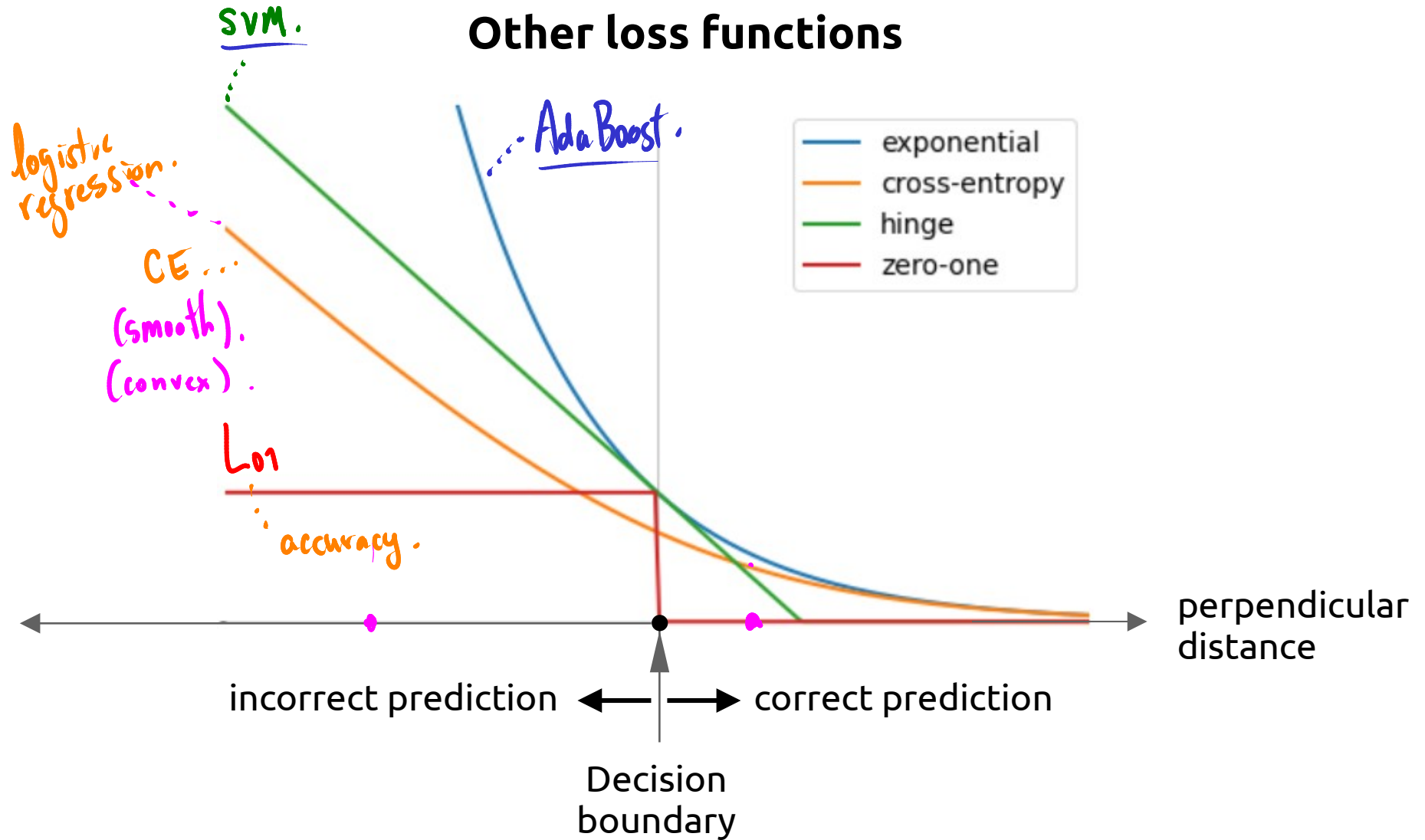
$\{0,1\}$ $(0,1)$

$$\text{CE}(y, p) = -y \log(p) - (1 - y) \log(1 - p)$$

$$\text{CE}(y, p) = -y \log(p) - (1 - y) \log(1 - p) = \begin{cases} -\log(p) & y = 1 \text{ — green} \\ -\log(1 - p) & y = 0 \text{ — red} \end{cases}$$



Other loss functions



Multi-class logistic regression

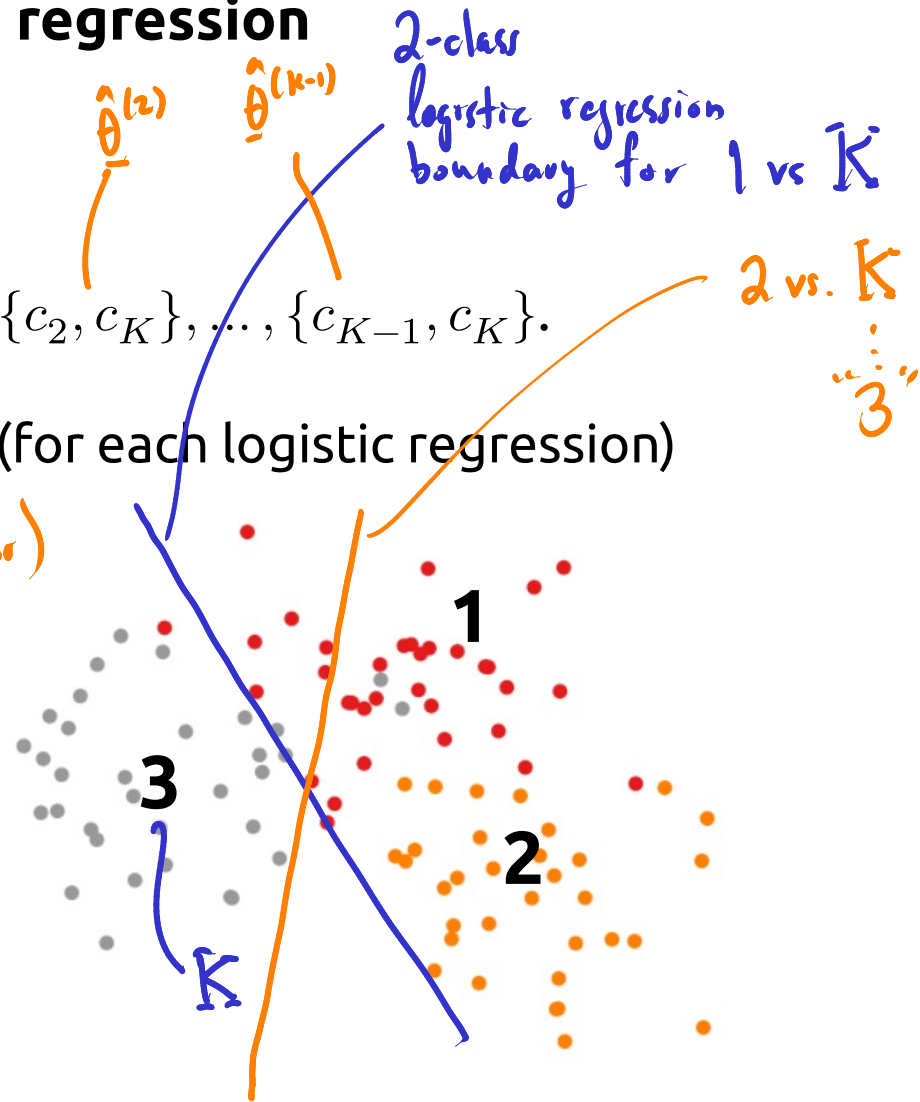
- $\Omega_Y = \{c_1, \dots, c_K\}$
- c_K is a "reference class".
- Solve $K - 1$ logistic regressions: $\{c_1, c_K\}, \{c_2, c_K\}, \dots, \{c_{K-1}, c_K\}$.
- Obtain parameter **vectors** $\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(K-1)}$ (for each logistic regression)

e.g. $K = 3$

(added a "1" to x vector)

— $\log \frac{p(Y=1 | X=x)}{p(Y=3 | X=x)} = x^T \hat{\theta}^{(1)}$

— $\log \frac{p(Y=2 | X=x)}{p(Y=3 | X=x)} = x^T \hat{\theta}^{(2)}$



In general:

$$\log \frac{p(Y=c_k | X=x)}{p(Y=c_K | X=x)} = x^T \hat{\theta}^{(k)} \quad k = 1 \dots K-1$$

\exp (

$$\therefore p(Y=c_k | X=x) = p(Y=c_K | X=x) \exp(x^T \hat{\theta}^{(k)}) \quad k = 1 \dots K-1 \quad (\text{I})$$

Also:

$$p(Y=c_K | X=x) = 1 - \sum_{\kappa=1}^{K-1} p(Y=c_{\kappa} | X=x)$$

$$= 1 - p(Y=c_K | X=x) \sum_{\kappa=1}^{K-1} \exp(x^T \hat{\theta}^{(\kappa)})$$

$$\therefore p(Y=c_K | X=x) = \frac{1}{1 + \sum_{\kappa=1}^{K-1} \exp(x^T \hat{\theta}^{(\kappa)})} \quad (\text{II})$$

$$\sum_{k=1}^K p(Y=c_k | X=x) = 1$$

Softmax

$$p(Y = c_k | X = x) = \frac{\exp(x^T \hat{\underline{\theta}}^{(k)})}{1 + \sum_{\kappa=1}^{K-1} \exp(x^T \hat{\underline{\theta}}^{(\kappa)})}$$

