

Announcement

- 1) HW 1 due today
- 2) HW 2 release today, due 5/5

Logistic Regression

Process

$$\text{Data } \{(x_i, y_i)\}_{i=1}^n \quad y_i \in \{0, 1\}, x \in \mathbb{R}^d$$

Decide on a **model**

$$f: \mathbb{R}^d \rightarrow \{0, 1\}, f \in \mathcal{F}$$
$$f(x) = \underset{y}{\operatorname{argmax}} p(y|x)$$

Find the function which fits the data best

Choose a loss function $l(f(x), y)$

Pick the function which minimizes loss

on data

$$\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i)$$

Use function to make prediction on new examples

$$x_{\text{new}} \in \mathbb{R}^d, \quad \hat{f}(x_{\text{new}})$$

Logistic Regression

$$\begin{aligned} x &\in \mathbb{R}^d \\ w &\in \mathbb{R}^d \\ y &\in \{0, 1\} \end{aligned}$$

Actually classification, not regression :)

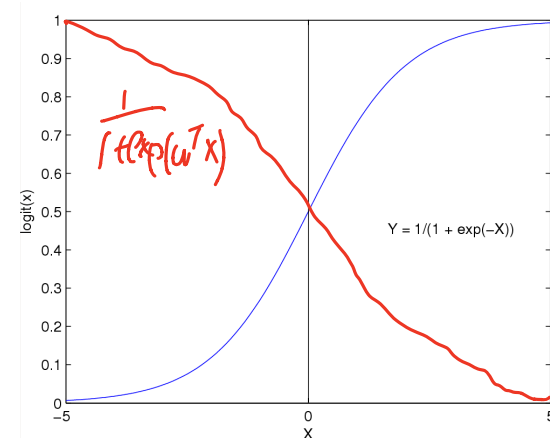
Learn $\mathbb{P}(Y = 1|X = x)$ using $\sigma(w^T x)$, for link function $\sigma =$

Logistic function(or Sigmoid):

$$\frac{1}{1 + \exp(-z)}$$

$$\mathbb{P}[Y = 1|X = x, w] = \sigma(w^T x) = \frac{1}{1 + \exp(-w^T x)}$$

$$\begin{aligned} \mathbb{P}[Y = 0|X = x, w] &= 1 - \sigma(w^T x) = \frac{\exp(-w^T x)}{1 + \exp(-w^T x)} \\ &= \frac{1}{1 + \exp(w^T x)} \end{aligned}$$



Features can be discrete or continuous!

Sigmoid for binary classes

$w_0, w_1, \dots, w_d \in \mathbb{R}$
 w_0 : offset
 $X_k \in \mathbb{R}, X \in \mathbb{R}^d$

$$\mathbb{P}(Y = 0|w, X) = \frac{1}{1 + \exp(w_0 + \sum_k w_k X_k)}$$

$$\mathbb{P}(Y = 1|w, X) = 1 - \mathbb{P}(Y = 0|w, X) = \frac{\exp(w_0 + \sum_k w_k X_k)}{1 + \exp(w_0 + \sum_k w_k X_k)}$$

$$\frac{\mathbb{P}(Y = 1|w, X)}{\mathbb{P}(Y = 0|w, X)} = \exp\left(w_0 + \sum_{k=1}^d w_k X_k\right)$$

exp in X, w
if magnitude is large (or small)
ratio is extremely large
($\rightarrow 0$)

Sigmoid for binary classes

$$\mathbb{P}(Y = 0|w, X) = \frac{1}{1 + \exp(w_0 + \sum_k w_k X_k)}$$

$$\mathbb{P}(Y = 1|w, X) = 1 - \mathbb{P}(Y = 0|w, X) = \frac{\exp(w_0 + \sum_k w_k X_k)}{1 + \exp(w_0 + \sum_k w_k X_k)}$$

$f(x) = \underset{y}{\operatorname{argmax}} \mathbb{P}(y|X)$

$$\frac{\mathbb{P}(Y = 1|w, X)}{\mathbb{P}(Y = 0|w, X)} = \exp(w_0 + \sum_k w_k X_k)$$

\Rightarrow

$> 1 \Rightarrow \text{predict } 1$
 $< 1 \Rightarrow \text{predict } 0$
 $= 1 \Rightarrow \text{both are OK}$

\Leftrightarrow

$$\log \frac{\mathbb{P}(Y = 1|w, X)}{\mathbb{P}(Y = 0|w, X)} = w_0 + \sum_k w_k X_k$$

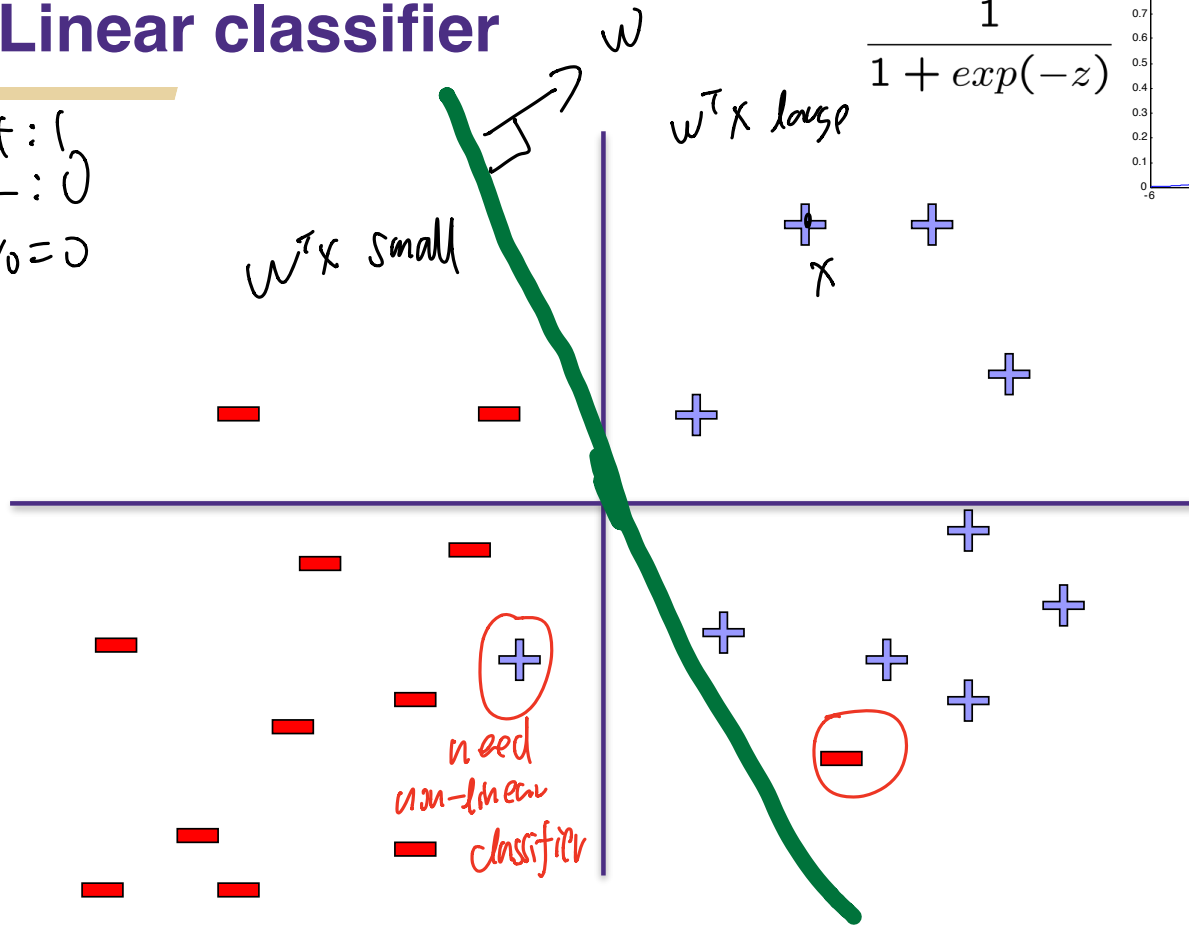
Linear Decision Rule!

\Rightarrow

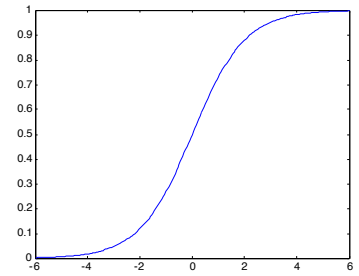
$> 0 \Rightarrow \text{predict } 1$
 $< 0 \Rightarrow \text{predict } 0$
 $= 0 \Rightarrow \text{both are OK}$

Logistic Regression – a Linear classifier

$f: \{ \}$
 $-: 0$
 $w_0 = 0$



$$\frac{1}{1 + \exp(-z)}$$



$$\log \frac{\mathbb{P}(Y = 1|w, X)}{\mathbb{P}(Y = 0|w, X)} = w_0 + \sum_k w_k X_k$$

Process

Decide on a **model**

$$f(x) = \begin{cases} 1 & w^T x > 0 \\ 0 & \text{otherwise} \end{cases}$$

Find the function which fits the data best

Choose a loss function

**Pick the function which minimizes loss
on data**

Use function to make prediction on new
examples

Loss function: Conditional Likelihood

encoding y_i
only simplicity
if $y_i \in \{-1, 1\}$
 $\Rightarrow 2y_i - 1 \in \{-1, 1\}$

- Have a bunch of iid data: $\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d, \quad y_i \in \{-1, 1\}$

$$P(Y = -1|x, w) = \frac{1}{1 + \exp(w^T x)}$$

$$P(Y = 1|x, w) = \frac{\exp(w^T x)}{1 + \exp(w^T x)}$$

- This is equivalent to:

$$P(Y = y|x, w) = \frac{1}{1 + \exp(-y w^T x)}$$

- So we can compute the maximum likelihood estimator:

$$\hat{w}_{MLE} = \arg \max_w \prod_{i=1}^n P(y_i|x_i, w)$$

w : parameter
want to learn

Loss function: Conditional Likelihood

- Have a bunch of iid data: $\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d, \quad y_i \in \{-1, 1\}$

$$P(Y = y|x, w) = \frac{1}{1 + \exp(-y w^T x)}$$

$$\hat{w}_{MLE} = \arg \max_w \prod_{i=1}^n P(y_i|x_i, w)$$

log is monotonic

$$= \arg \min_w \sum_{i=1}^n \underbrace{\log(1 + \exp(-y_i x_i^T w))}_{\ell(x_i, y_i)}$$

Logistic Loss: $\ell_i(w) = \log(1 + \exp(-y_i x_i^T w))$

for classification

Squared error Loss: $\ell_i(w) = (y_i - x_i^T w)^2$

for regression

(MLE for Gaussian noise)

Process

- what we really care
- o/r: $\{ f(x) \neq y \}$
 - $\log(1 + \exp(-y w^T x))$
for training

Decide on a **model**

- MLE principle
- o/r it's hard to optimize

Find the function which fits the data best

Choose a loss function

Pick the function which minimizes loss

on data $\arg\min_w \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$

Use function to make prediction on new examples

Loss function: Conditional Likelihood

- Have a bunch of iid data: $\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d, \quad y_i \in \{-1, 1\}$

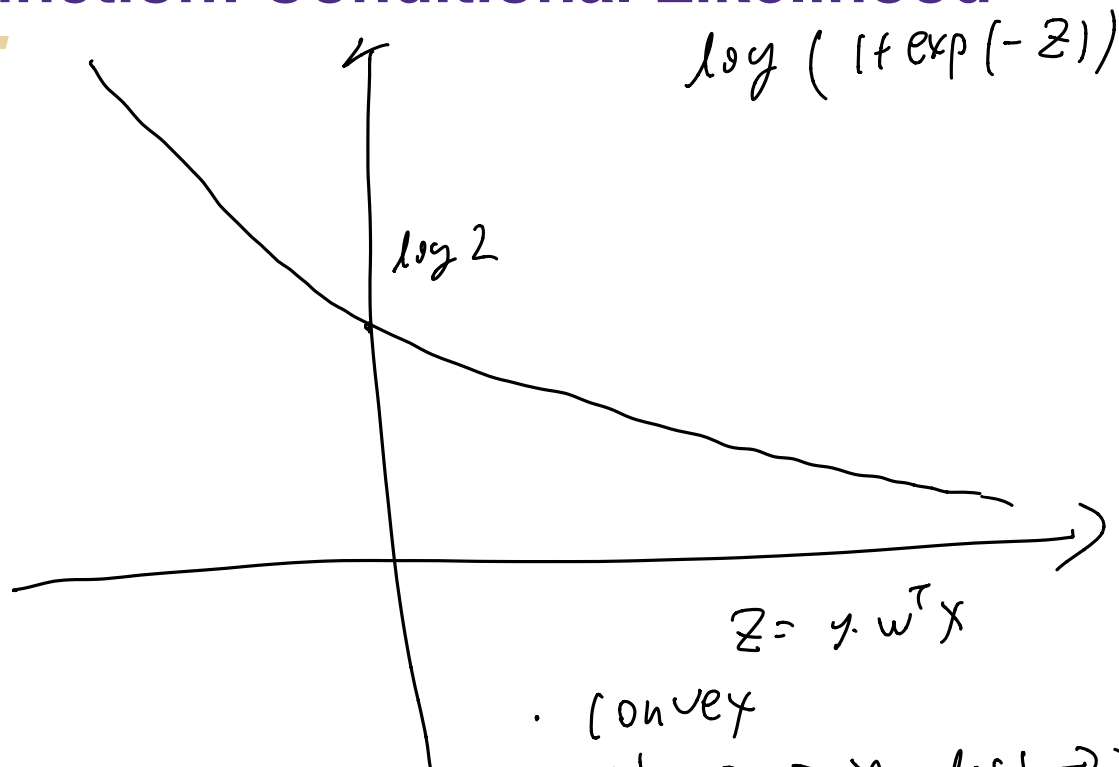
$$P(Y = y|x, w) = \frac{1}{1 + \exp(-y w^T x)}$$

$$\begin{aligned} \hat{w}_{MLE} &= \arg \max_w \prod_{i=1}^n P(y_i|x_i, w) & \mathcal{J}(w) &= \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w)) \\ &= \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w)) = J(w) \end{aligned}$$

What does $J(w)$ look like? Is it convex?

$$\begin{aligned} \mathcal{J}''(w) &\succ 0 \\ w &\in \mathcal{P} \end{aligned}$$

Loss function: Conditional Likelihood



- convex

- if $z \rightarrow \infty$, loss $\rightarrow 0$
 $z \rightarrow -\infty$, loss $\rightarrow \infty$

intuition: if y & $w^T x$ have same sign \rightarrow loss ^{small}
if y & $w^T x$ have different sign \rightarrow loss ^{big}

Loss function: Conditional Likelihood

- **Have a bunch of iid data:** $\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d, \quad y_i \in \{-1, 1\}$

$$P(Y = y|x, w) = \frac{1}{1 + \exp(-y w^T x)}$$

$$\begin{aligned} \hat{w}_{MLE} &= \arg \max_w \prod_{i=1}^n P(y_i | x_i, w) \\ &= \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w)) = J(w) \end{aligned}$$

Good news: $J(\mathbf{w})$ is convex function of \mathbf{w} , no local optima problems

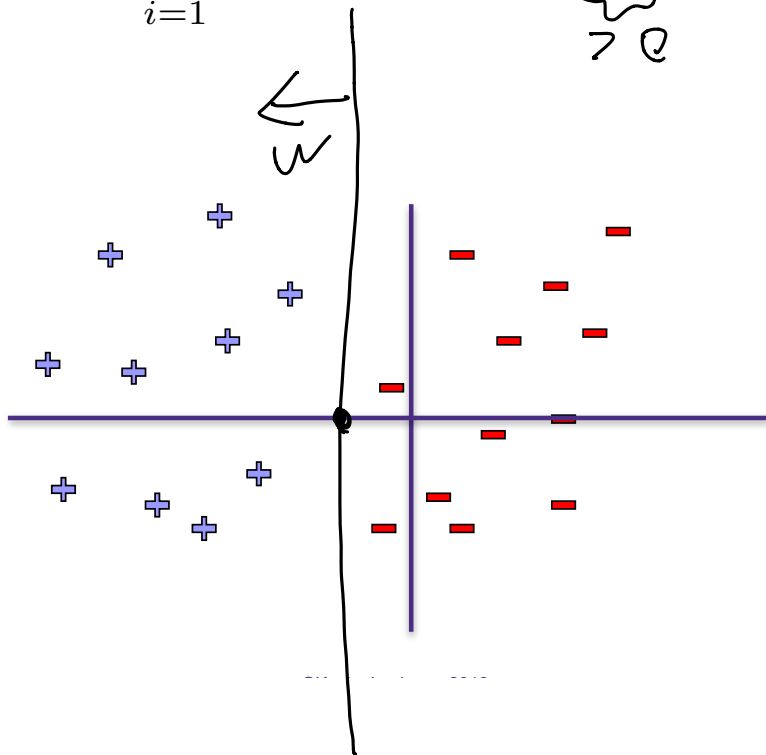
Bad news: no closed-form solution to maximize $J(\mathbf{w})$

Good news: convex functions easy to optimize

(gradient descent)

Overfitting and Linear Separability

$$\arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i \underbrace{x_i^T w}_{> 0}))$$



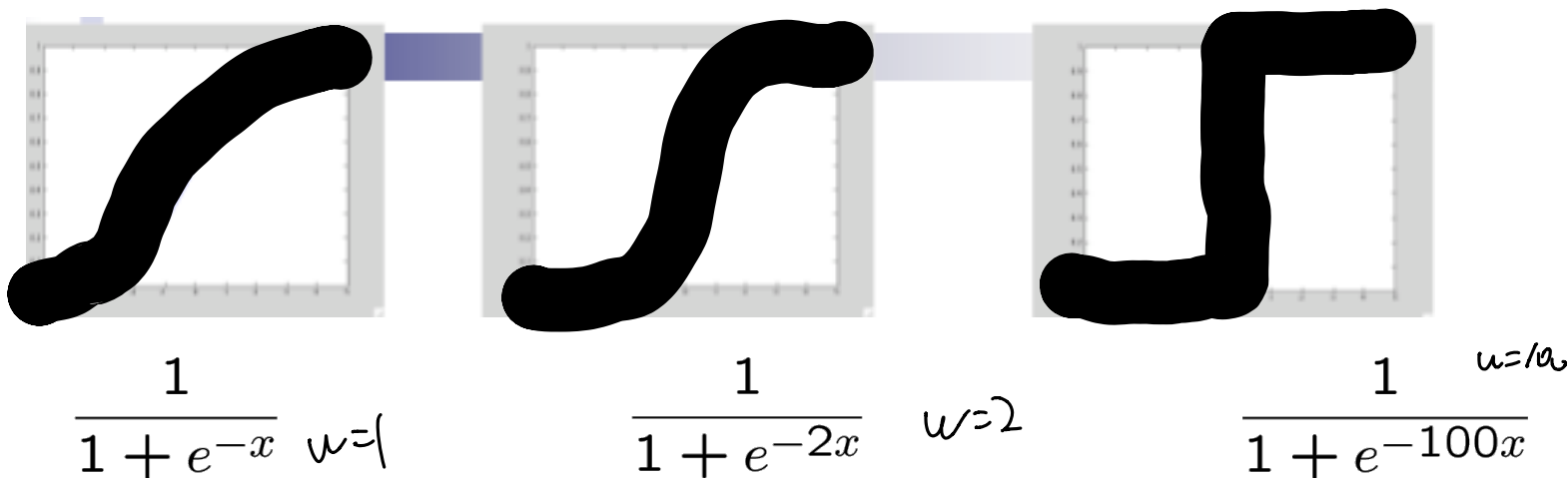
When is this loss small?

(1) divergence, offset
 $\text{sign}(x_i^T w) = \text{sign}(y_i)$

(2) magnitude of w
 $w \rightarrow 2w \rightarrow 4w \dots$
 $\|w\|_2 \rightarrow \infty$
 $J(w) \rightarrow 0$

Large parameters \rightarrow Overfitting

When data is linearly separable, weights $\Rightarrow \infty$



Overfitting

large weight $p(y|x) \rightarrow 0$ or 1
 often not accurate \neq want $p(y|x) \in (0,1)$
 (i) constant

Penalize high weights to prevent overfitting?

Regularized Conditional Log Likelihood

$w_1(x, y) \rightarrow y > 0 \quad b < 0 \quad x^T w + b \neq y$
 $x^T w > 0 \rightarrow \text{if } \|w\| \text{ small, } x^T w + b < 0, \text{ if } \|w\| \text{ large, } w^T x + b > 0$

Add a penalty to avoid high weights/overfitting?:

$$\arg \min_{w, b} \underbrace{\sum_{i=1}^n \log(1 + \exp(-y_i (x_i^T w + b)))}_{\text{fit data}} + \underbrace{\lambda \|w\|_2^2}_{\text{regularization}}$$

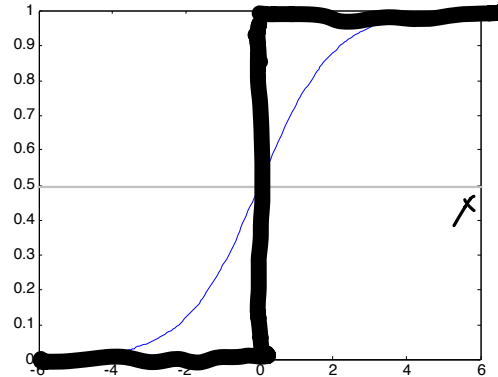
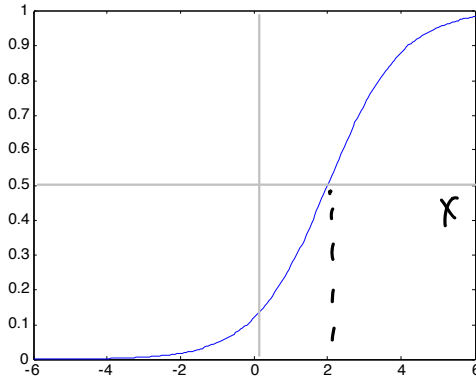
can also use l_1

Be sure to not regularize the offset b !

$$y = b(w_0 + w_1 \cdot x)$$

$$w_1 = 1$$

$$w_0 = -2$$



$$w_1 = 1$$

$$w_0 = 0$$