

Lec 07:

- ① Finish up regression
 - ② Classification
-

$$\|w\|_1 = \sum_{i=1}^d |w_i|$$

Regularization: Ridge: $\|\Phi w - y\|^2 + \lambda \|w\|_2^2 \rightarrow \text{argmin}$

LASSO: $\|\Phi w - y\|^2 + \lambda \|w\|_1$

OPT:

Ridge

$$\begin{aligned} \min & \|\Phi w - y\|^2 \\ \text{s.t.} & \|w\|_2 \leq c_1 \end{aligned}$$

LASSO:

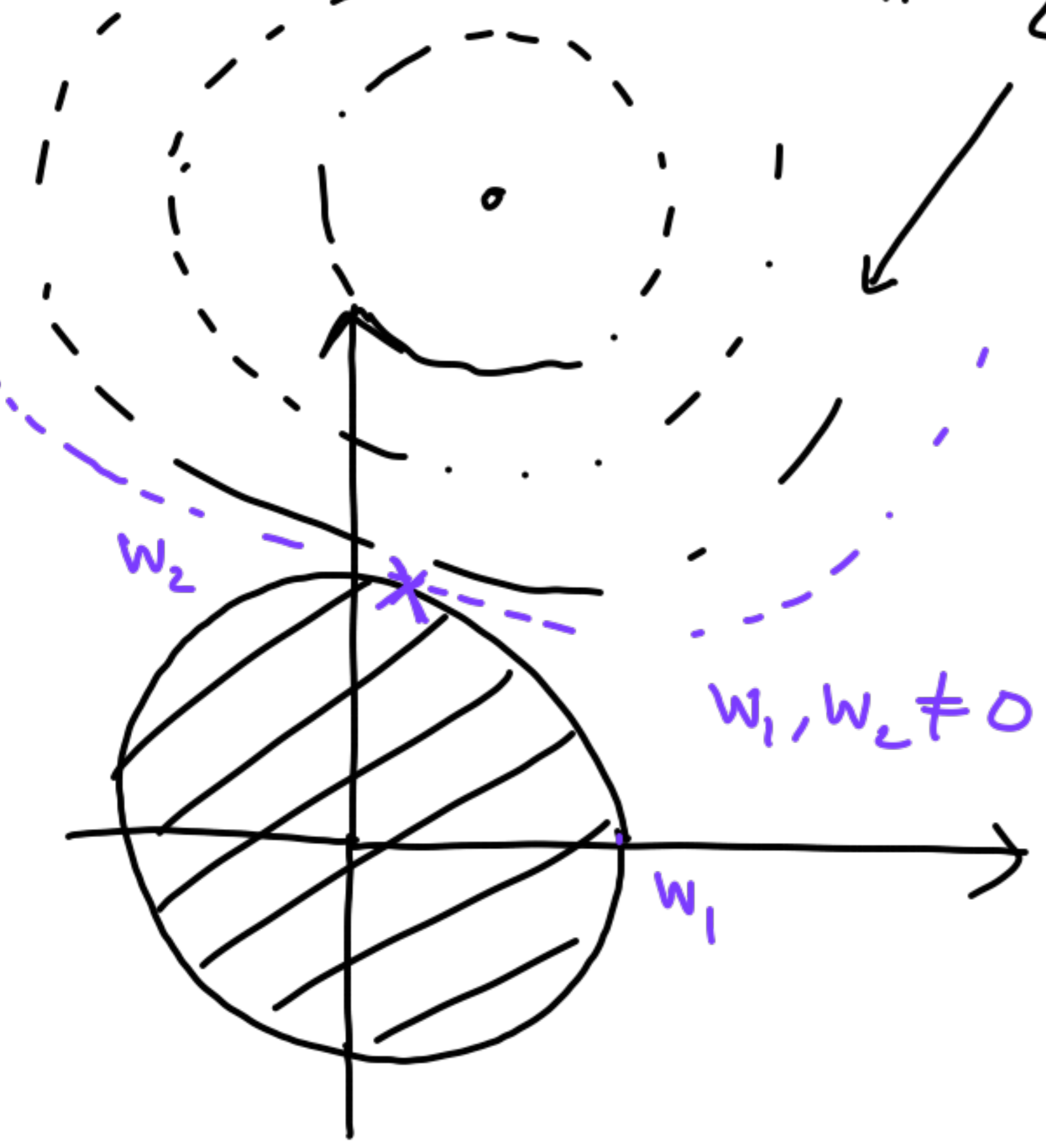
$$\begin{aligned} \min & \|\Phi w - y\|^2 \\ \text{s.t.} & \|w\|_1 \leq c_2 \end{aligned}$$

Constraints:

$$\|W\|_2 \leq C_1$$

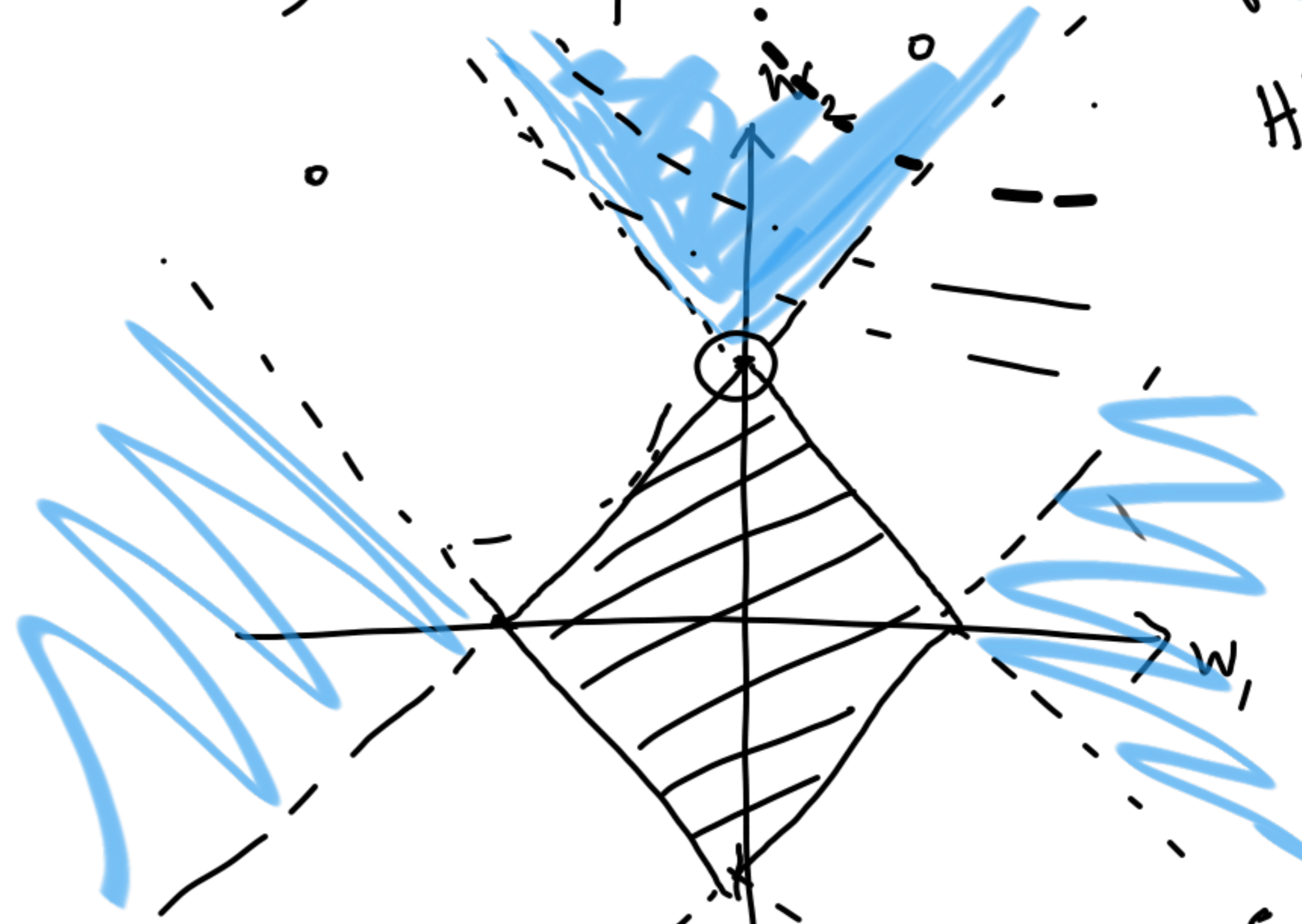
$$\|W\|_1 \leq C_2$$

Elements of Stat learning
HTF



$$w_1^2 + w_2^2 \leq C_1^2$$

Applications: all components are imp.



$$|w_1| + |w_2| \leq C_2$$

Some components are imp.

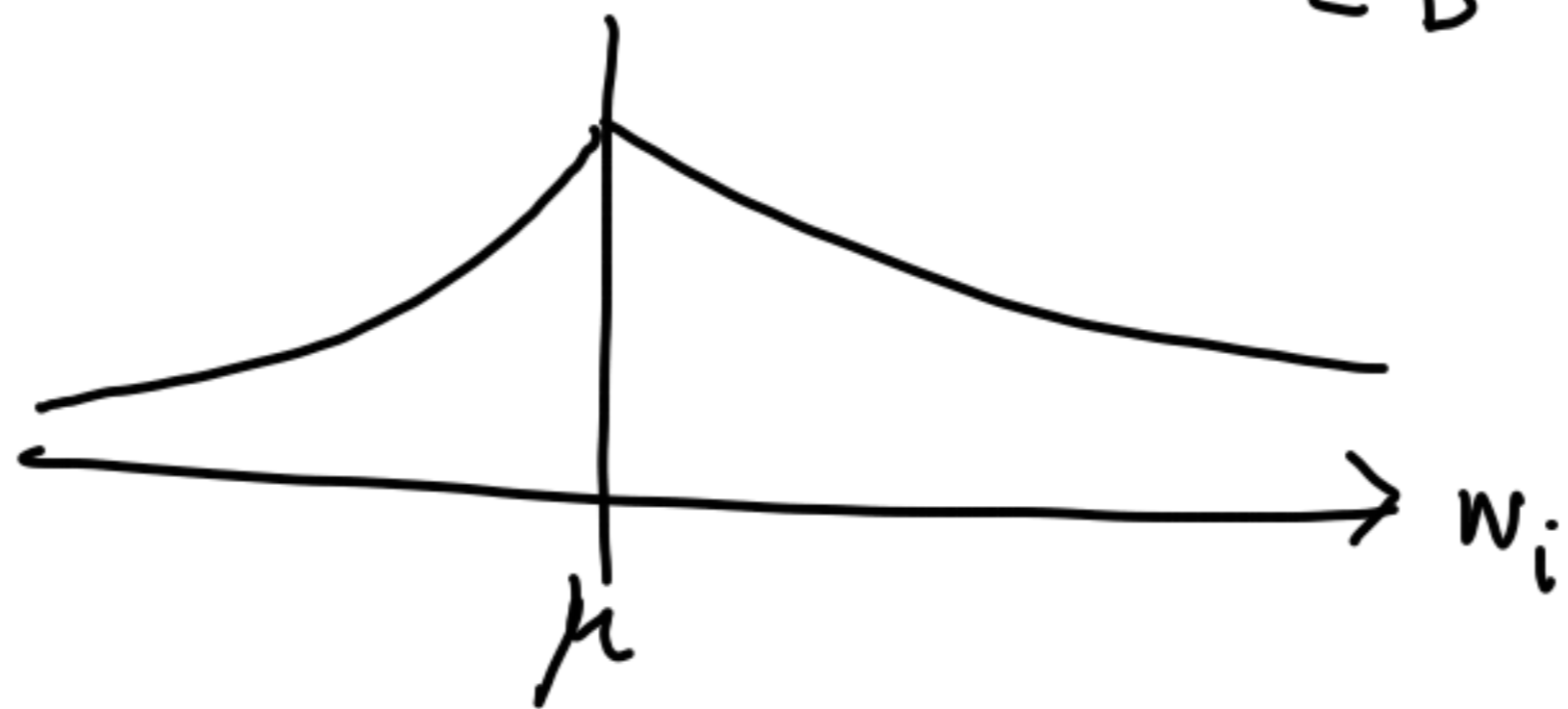
MAP : prior on w

Gaussian prior \rightarrow Ridge

Laplace prior \rightarrow LASSO

$$\text{Laplace}(w_i | \mu, b) = \frac{1}{2b} \exp\left\{-\frac{|w_i - \mu|}{b}\right\}$$

HW



Classification

Regression: $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$

Y $y_i \in S \leftarrow$ finite set, simplest $S = \{0, 1\}$

Binary classification

Probabilistic approach:

$$\hat{y} = \underset{y}{\operatorname{argmax}} P(Y=y | \hat{x})$$
$$P(Y=y | \hat{x}) = \frac{P(\hat{x} | Y=y) P(Y=y)}{P(\hat{x})}$$

$\hat{x} \rightarrow \hat{y}$
 $\hat{x} = (\hat{x}_1, \dots, \hat{x}_d)^T$

Naive Bayes

x_1, x_2, \dots, x_d \nearrow $(2^d - 1) \times (k-1)$

Assumption:

$$P(x | Y=y) = \prod_{i=1}^d P(x_i | Y=y)$$

$$\begin{aligned} \underset{y}{\operatorname{argmax}} P(Y=y | x) &= \underset{y}{\operatorname{argmax}} \left(\prod_{i=1}^d P(x_i | Y=y) \right) P(Y=y) \\ &= \underset{y}{\operatorname{argmax}} \left\{ \sum_{i=1}^d \log P(x_i | Y=y) + \log P(Y=y) \right\} \end{aligned}$$

$(d \times 1 \times (k-1))$

$$= \frac{P(x_i=1 | Y=y)}{\# \text{ of times } x_i \text{ and } y \text{ happen together}}$$

$$= \frac{\# \text{ of times } x_i \text{ and } y \text{ happen together}}{\# \text{ of times } Y=y}$$

$$P(Y=y) = \frac{\# \text{ of times } Y=y}{n}$$

Uses of NB classifier: Topic classification

Given a document \rightarrow find its topic

Treat the document as a "bag of words"

given sentences \rightarrow create a vocab of words

(tokenization) \rightarrow index them arbitrarily

Article: $x = \{x_1, x_2, \dots, x_N\}$ x_i is the i^{th} word

$$P(x_i | Y=y) = \frac{\# \text{ of times } x_i \text{ appears in the doc of } y + 1}{\sum_{w \in V} \# \text{ of times } w \text{ appears in } y + |V|}$$

for unseen words.

$$P(x | Y=y) = \prod_{i=1}^N P(x_i | Y=y)$$

$P(Y=y)$ = relative fraction of y in all topics.

extending x_i for a group of words \rightarrow n-grams

Disadvantages of NB

True distribution of binary variable X_1 and Y

$Y=0$	$Y=1$
0.8	0.2

$P(Y)$

	$X_1=0$	$X_1=1$
$Y=0$	0.7	0.3
$Y=1$	0.3	0.7

$P(X_1|Y)$

$$P(Y=0) = 0.8$$
$$P(Y=1) = 0.2$$

NB classifier

$$\hat{y} = \underset{y}{\operatorname{argmax}} P(Y=y | X_1) = \underset{y}{\operatorname{argmax}} P(X_1 | Y=y) P(Y=y)$$

	$P(X_1, Y=0)$	$P(X_1, Y=1)$	\hat{y}
$X_1=0$	0.7×0.8	0.3×0.2	0
$X_1=1$	0.3×0.8	0.7×0.2	0

X_2 identical to X_1 → copy of

	$P(X_1, X_2, Y=0)$	$P(\dots, Y=1)$	\hat{y}
$X_1=X_2=0$	$0.7^2 \times 0.8$	$0.3^2 \times 0.2$	0
$X_1=X_2=1$	$0.3^2 \times 0.8$	$(0.7)^2 \times 0.2$	1
$X_1=0, X_2=1$	✓	✓	0
$X_1=1, X_2=0$	✓	✓	0

$$\begin{aligned} \text{Expected error} &= P(Y \neq \hat{y}) \\ &= P(Y \neq \hat{y}, X_1=0) + P(Y \neq \hat{y}, X_1=1) \\ &= P(Y=1, X_1=0) + P(Y=0, X_1=1) \\ &= 0.06 + 0.14 = \underline{0.2} \end{aligned}$$

$$\begin{aligned} P(Y \neq \hat{y}) &= P(Y=1, X_1=0) + P(Y=0, X_1=1) \\ &= 0.3 \end{aligned}$$

Logistic Regression (Classification)

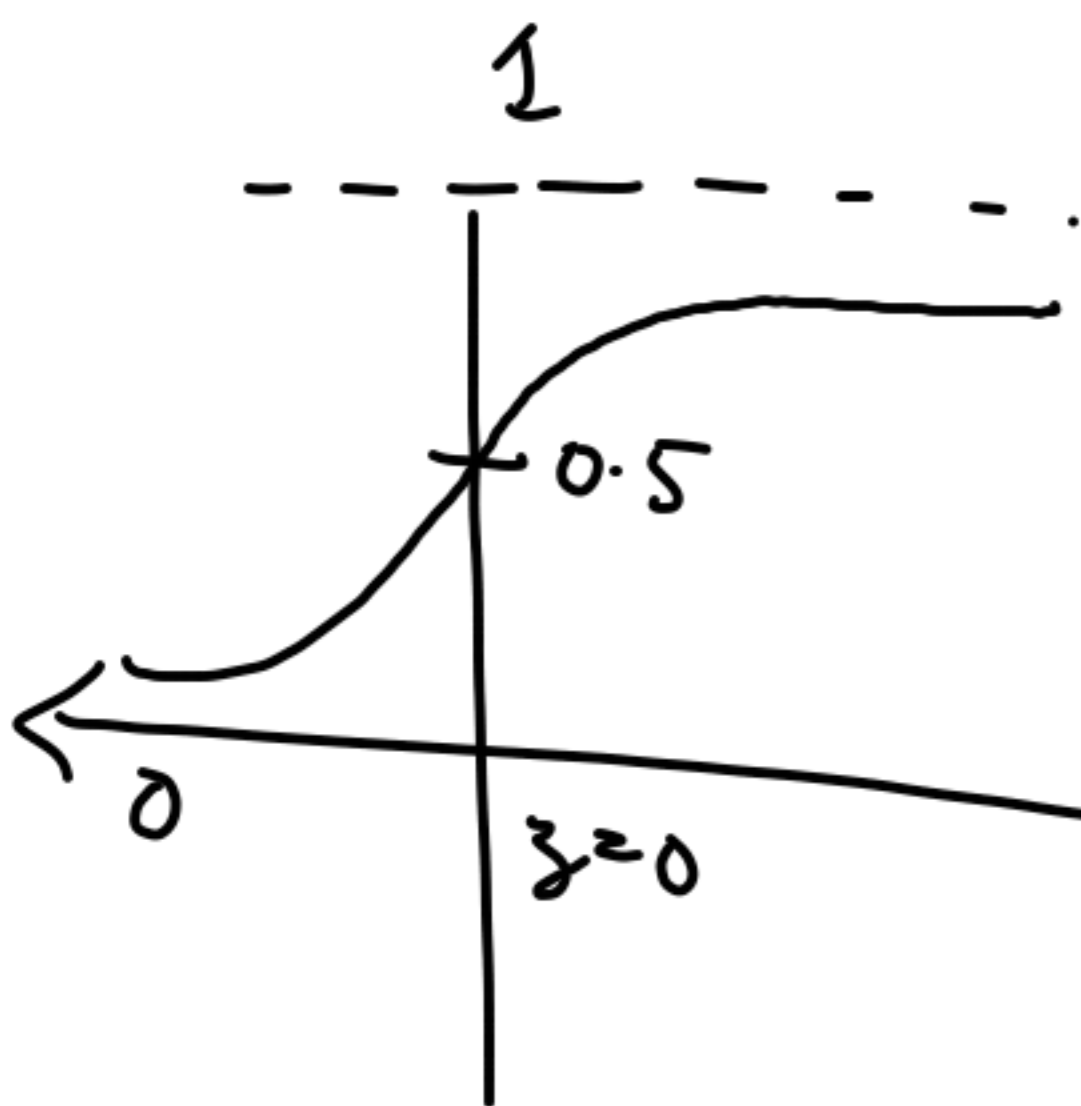
$$y \approx W^T x$$

Binary classification

$$W_1^T x$$

$$x \in \mathbb{R}^d$$

$$W_2^T x$$



$$\frac{e^{W_1^T x}}{e^{W_1^T x} + e^{W_2^T x}}$$
$$\frac{e^{W_2^T x}}{e^{W_1^T x} + e^{W_2^T x}}$$

$$\frac{1}{1 + e^{\underbrace{(W_2 - W_1)^T x}_{-W}}} = \frac{1}{1 + e^{-W^T x}} = \sigma(W^T x)$$

Assumption: $P(Y=1 | x, W)$
 $= \sigma(W^T x)$

$$\frac{P(y=1 | x, w)}{P(y=0 | x, w)} > 1 \Rightarrow \hat{y} = 1$$

$$0 < \frac{P(y=1 | x, w)}{P(y=0 | x, w)} < 1 \Rightarrow \hat{y} = 0$$

$$\exp(w^T x) > 1 \Rightarrow w^T x > 0$$

