

Statistics and Data Science 365 / 565

# Data Mining and Machine Learning

February 15

Yale

# Outline

- Recap: supervised learning goal
- Recap: Function classes
- Classification
- Surrogate Losses
- Population Risk
- Notebook

# Recap: Supervised learning goals

hat means estimate of sth

Create a function  $\hat{f} : \mathcal{X} \mapsto \mathcal{Y}$  by “learning” from data  $z_i = (x_i, y_i)$ .

Probably want  $\hat{f}(x_i) \approx y_i$  (**Notation:**  $\approx$  means approximately)

**Regression:**  $\mathcal{Y}$  is continuous/ordered (prices)

**Classification:**  $\mathcal{Y}$  is discrete/unordered (labels: hot dog v not hot dog)

Focus on classification today.

Recap: Function Classes 画出的类型

① **Linear functions:** Take  $x \in \mathbb{R}^d$

$$\mathcal{F} = \{f \mid f(x) = \theta^T x \text{ for some } \theta \in \mathbb{R}^d\}$$

$\nearrow$  weights  
 $\searrow$  features

**Learning:** Optimize loss for weights

Learning of how much a change in input will affect change in output

weight: sometimes weight and features are correlated colinearity

② **Nearest neighbors:** Make predictions based on neighbors

**Learning:** ① Training example proximity. ② Number of neighbors. ③ Choice of distance.

who were the  $k$  nearest neighbors

$k$

Often use Euclidean distance ( $L_2$  norm)

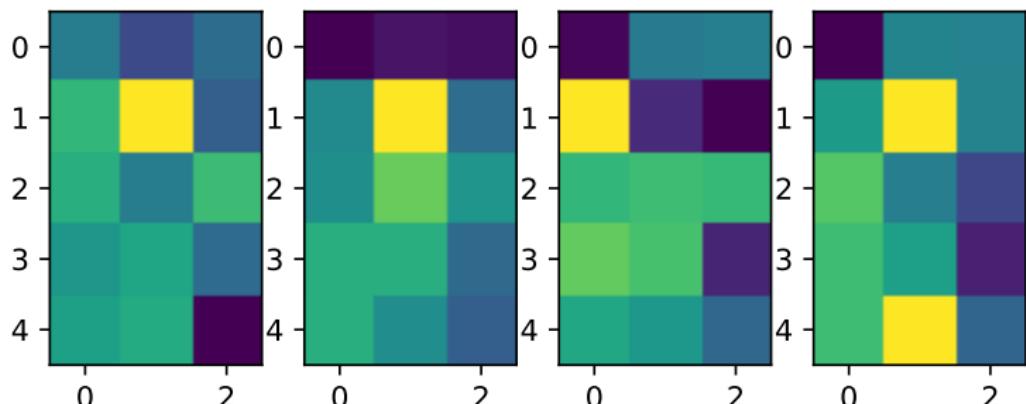
# Linear function weights: classification

$$y = -1$$

one label

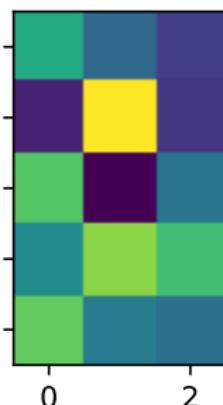
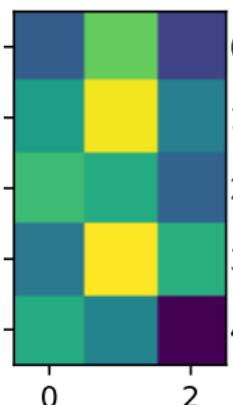
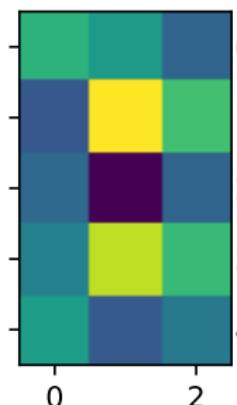
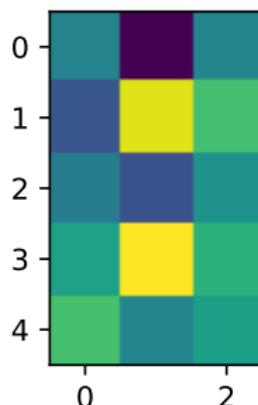
data set: images 7 and 8 with Gaussian noise

A Eg. of  $X$  with label 1



# Linear function weights: classification

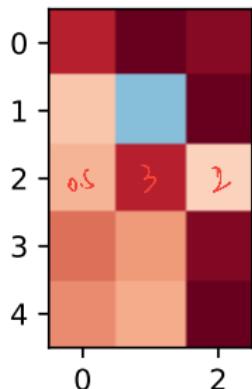
$y = 1$       another label



# Linear function weights

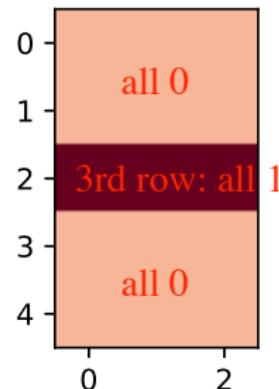
label  $y = \{-1, +1\}$        $\hat{y}_i = 5.5 > 1 \rightarrow +1$  label  
Discriminative learning       $R \rightarrow$  binary classification

$x$   
 $5 \times 3$  pixels = 15 dimensions



$\theta$  weight

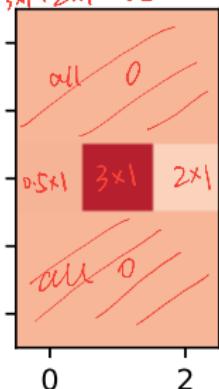
$\times$



scale  $x$  by weight  $\theta$   
point-wise multiplication  
 $x \cdot \theta$

$$x \cdot \theta = \text{a real number} \\ = 0.5x_1 + 3x_2 + 2x_3 = 5.5$$

=



# Linear function weights

This is a classification problem.

Above choice has mean Hamming error 0.128. Can we do better?

↓  
find optimal value of  
Hamming error

if  $f(x)$  is linear function  $\Rightarrow$  square loss  $\neq$  Hamming loss

error	y	$f(x)$	interpretation
81	+1	too positive ( $ x $ )	$f(x)$ very confident y is a positive thing
121	+1	too negative ( $- x $ )	bigger error

This is a classification problem. Directly optimizing Hamming loss is hard. NP hard computational intractable

Will just use squared loss. Add an intercept.

$$\text{sign}(f(x)) = \begin{cases} +1 & f(x) > 0 \\ 0 & f(x) = 0 \\ -1 & f(x) < 0 \end{cases}$$

if  $f(x)$  is discrete

error	y	$f(x)$
0	+1	1
4	+1	-1
4	-1	1

square loss + 1/4 = Hamming loss  
arg min don't care scaling

① Hamming loss  $\Rightarrow \mathbb{I}(y \neq f(x))$   
hard to have  $y = f(x)$   
 $\Rightarrow$  instead  $\mathbb{I}(y \neq \text{sign}(f(x)))$

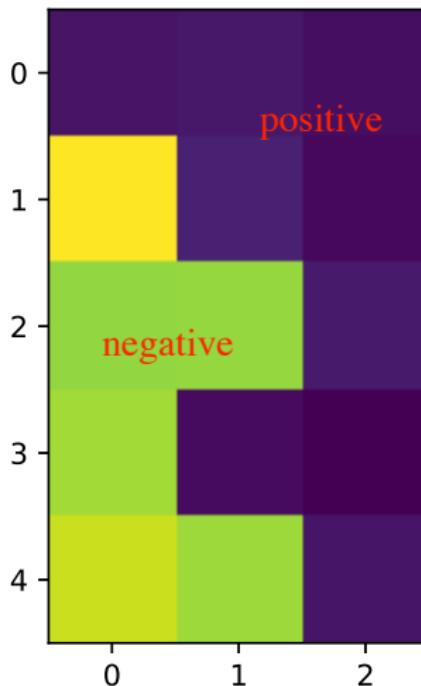
② Squared loss:  $(y - f(x))^2$  easier

$$y = 0, \pm 1$$

# Achieves 1% training error (Hamming loss)

Optimize square loss can also achieve very small (nice) Hamming loss

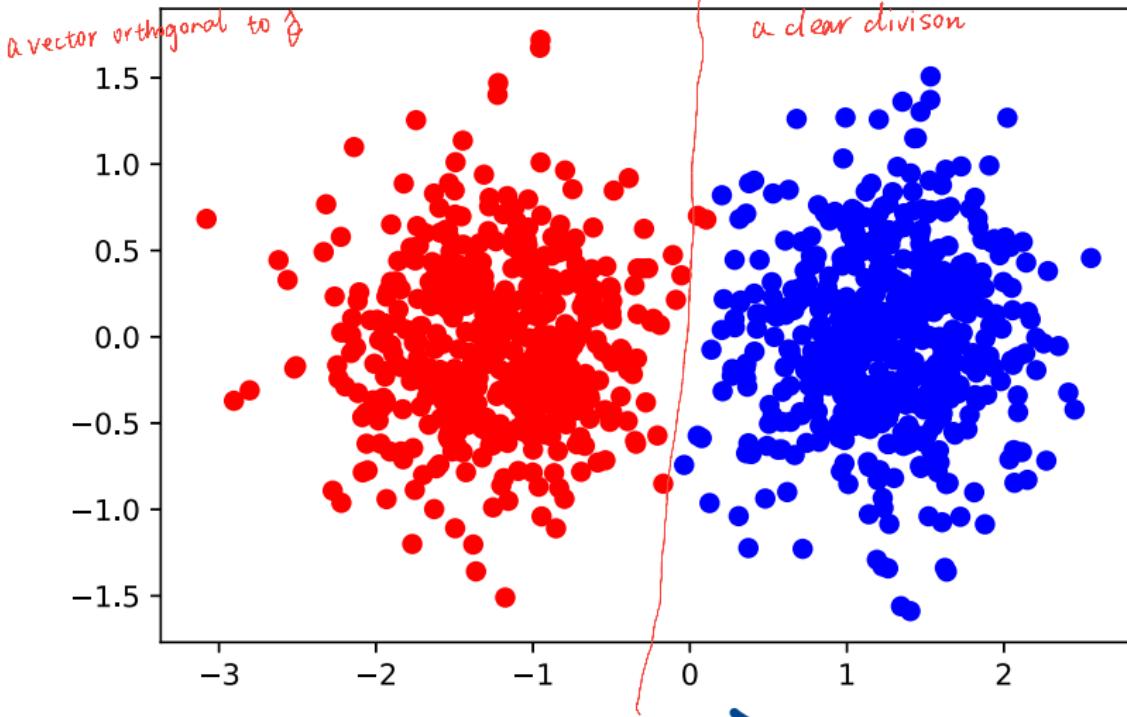
$$\hat{\theta}$$



Plot data points along  $\hat{\theta}$  direction

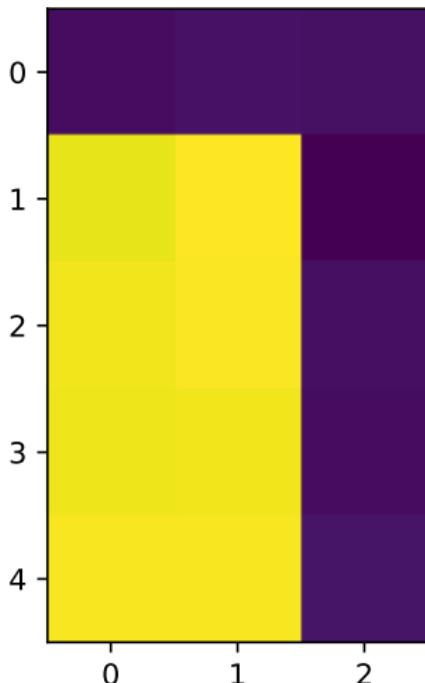
visualize 15D data to 2D

intercept term

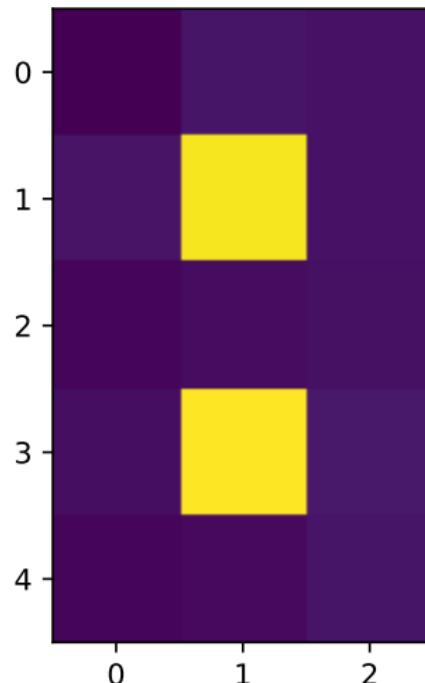


$$\hat{f}(x) = z \times \hat{\theta}$$

Labels == 0



Labels == 1



average of data point with label == 0 is image 7

# Classification

*true loss*

Directly optimizing Hamming loss is hard

Also not always a great idea

*don't have nice property*

Rely on surrogate losses  
*→ alternative*

easier to select parameters

- { ● squared loss above *(previous)*
- logistic loss
- hinge loss
- boosting loss

# Surrogate Losses

Let's take  $\mathcal{Y} = \{-1, +1\}$

$$= \mathbb{1}(y \neq \text{sign}(f(x)))$$

**Hamming loss:**  $\mathbb{1}(y \neq f(x)) = \mathbb{1}(f(x)y \leq 0)$ . Why? f(x)和y异号说明 f(x) ≠ y

optimal: often

**squared loss:**  $(y - f(x))^2 = y^2(1 - f(x)y)^2$

if  $f(x) = \{-1, +1\}$ , then sq loss = Ham loss

若y不是binary而是continuous

$$(y - f(x))^2 = [y(1 - \frac{f(x)}{y})]^2$$

**logistic loss:**  $\log(1 + \exp(-2f(x)y)) / \log(2)$

if  $f(x) = \{-1, +1\}$ , then logistic loss = hamming loss

**hinge loss:**  $(1 - yf(x))_+$

Notation: use to plot nice figure,

arg min don't care scaling 2,  $\log(2)$

if  $f(x) = \{-1, +1\}$ , then hin loss = Ham loss

$$(v)_+ = \begin{cases} v & v > 0 \\ 0 & v \leq 0 \end{cases}$$

positive parts of a number

**boosting loss:**  $\exp(-yf(x))$

# Surrogate Losses

All of the losses above function of  $f(x)y$

In classification could learn function  $f : \mathcal{X} \mapsto \{-1, +1\}$  Hamming loss

Hard. Usually learn  $f : \mathcal{X} \mapsto \mathbb{R}$  surrogate loss

Then take  $\text{sign}(f)$

margin:  $f(x)y$  - how correct or wrong

very positive: correct  
very negative: wrong

$f(x)y \geq 3$   
 $\geq 3$

If  $\mathcal{Y} = \{0, 1\}$  then take  $\mathbb{1}(f(x) > 0)$

# Linear logistic regression

to solve binary classification

most popular way

$$\hat{f}: \mathbb{X} \rightarrow \mathbb{R}$$

optimize over  $\theta$

$$\arg \min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T \theta))$$

Function class is linear functions

logistic loss

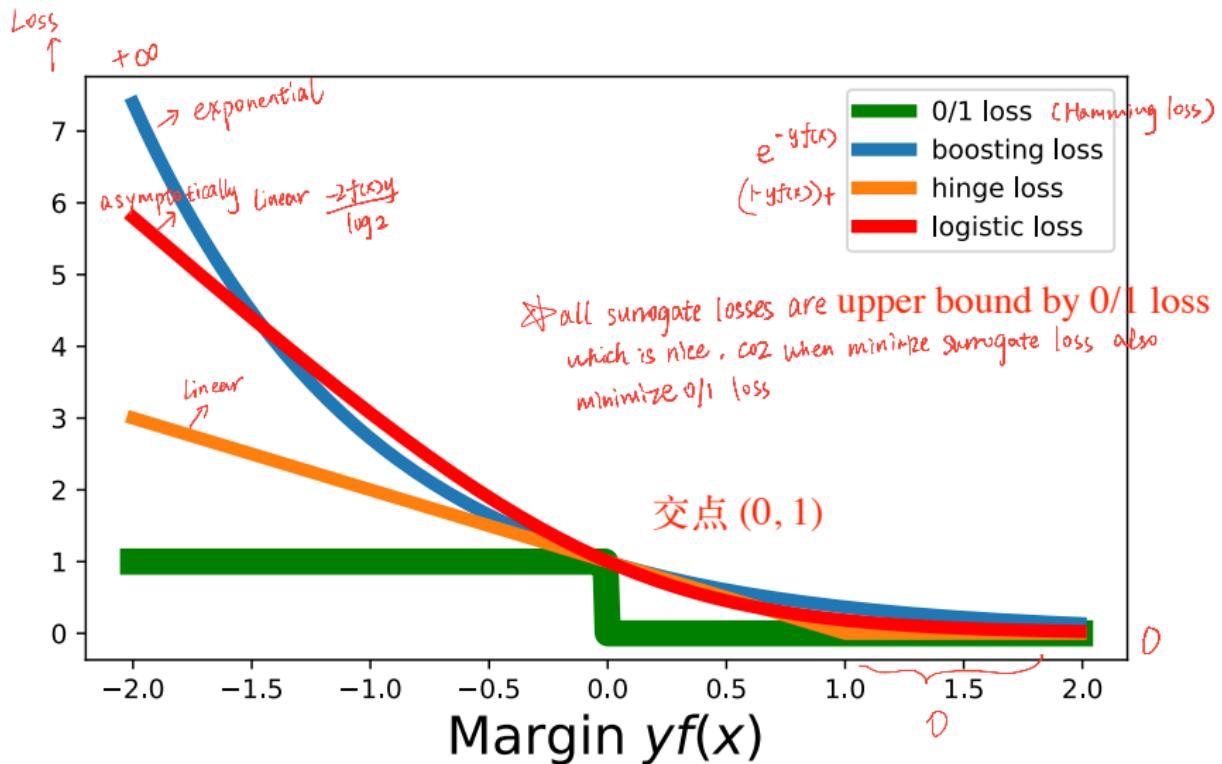
Easier to optimize over the above than to optimize

$$\arg \min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i \text{sign}(x_i^T \theta)))$$

non-linearity  
challenge with training  
neural network

no closed solution: but can use least squares solution  
and other methods

## Surrogate Losses



# Loss motivation

of both regression and classification

Have been optimizing

$$\arg \min_{f \in \mathcal{F}} \hat{R}(f)$$

$\mathcal{F}$ : function classes

linear function, KNN, neural networks, etc.

Motivation:

$$f^* = \arg \min_{f \in \mathcal{F}} \mathbb{E} \ell(f(x), y)$$

where  $x, y \sim \mathbb{P}(x, y)$

$x, y$  跟从某个分布

i.e. be good in expected value

Define population Risk or just Risk to be

$$R(f) = \mathbb{E} \ell(f(x), y)$$

Risk of function  $f$  is the expected loss  
of that function  
expectation is of pop, not sample

For fixed  $f$  if data identically distributed (not necessarily independent)

randomly uniformly pick a sample from a pop  $\xrightarrow{\text{coz expectation is linear}}$

$$\mathbb{E} \hat{R}(f) = R(f)$$

expected value of empirical risk = pop risk

Hope that with enough data  $\hat{R}(f)$  is close to  $R(f)$

Use central limit theorem, empirical process theory, Glivenko Cantelli,  
...  
*general*                           *specific*

## Caveats

- Data might not be independent (false sense of security for large  $n$ )  
*e.g. google flu*      *large n doesn't mean good estimate*  
*eq. time series*      *data correlate with adjacent data*
- Weighting data all the same (bias against minorities)  
*v e.g. Covid positive case*
- Choice of loss might be inappropriate for problem (robustness)  
*actually a paper proves different choice of loss is fine*

# Conclusion

Classification similar to regression, but with different losses

Use surrogates instead of Hamming loss

Empirical risk aims to approximate population risk

$$\hat{R}(f) \xrightarrow{\text{large sample}} R(f)$$

if data i.d.

$$E[\hat{R}(f)] = R(f)$$