

# Hongyi Zhang

## Lecture 1: Structural Risk Minimization (SRM)

Objective: minimizing a combination of **loss function** and **model complexity penalty**.

1.1 Machine Learning Problem: dataset  $D = \{(\vec{x}_i, y_i)\}_{i=1,\dots,n}$

goal to learn  $y = h(\vec{x})$  s.t.  $\{h(\vec{x}_i) = y_i, i=1,\dots,n\}$   
 $h(\vec{x}^*) = y^*$ , for unseen test data  
 $y = h(f(\vec{x}))$  with  $f: \vec{x} \rightarrow R$   $\begin{cases} \text{classification } f \in \mathcal{F}, h(a) = \text{sign}(a) \\ \text{regression } h(a) = I(a) \end{cases}$

Empirical risk minimization (ERM) (minimize training error)

$$\min_{\vec{w}} \text{R}_{\text{emp}}(\vec{w}) = \min_{\vec{w}} \frac{1}{n} \sum_{i=1}^n l(h(\vec{x}_i), y_i) \quad \begin{cases} \text{tends to overfit training data} \\ \text{a simpler model is preferred} \end{cases}$$

### 1.2 Structural Risk Minimization (SRM)

balance fitting the training data against model complexity  
 $\min_{\vec{w}} L(\vec{w}) = \min_{\vec{w}} \frac{1}{n} \sum_{i=1}^n l(h_{\vec{w}}(\vec{x}_i), y_i) + \lambda r(\vec{w})$   $\begin{cases} \text{hyperparameter} \\ \text{training loss} \\ \text{regularizer} \end{cases}$

## 2. Loss Function

$$l(f_{\vec{w}}(\vec{x}_i), y_i)$$

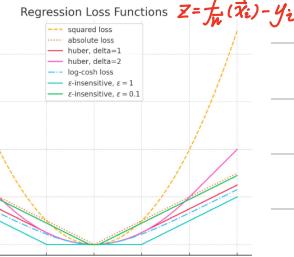
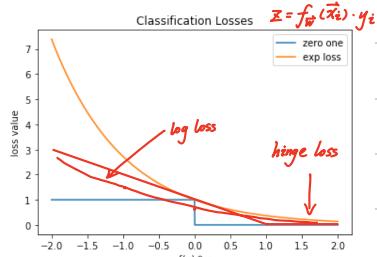
### ✓ 2.1 Binary Classification Loss Functions ( $y \in \{-1, +1\}$ )

zero-one loss:  $\delta(h_{\vec{w}}(\vec{x}_i) = y_i)$  non-continuous, impractical

hinge loss:  $\max[1 - f_{\vec{w}}(\vec{x}_i) \cdot y_i, 0]^p$  standard SVM ( $p=1$ )

log loss:  $\log(1 + e^{-f_{\vec{w}}(\vec{x}_i) \cdot y_i})$  logistic regression

exponential loss:  $e^{-f_{\vec{w}}(\vec{x}_i) \cdot y_i}$  AdaBoost, aggressive, sensitive



### ✓ 2.2 Regression Loss Functions ( $h_{\vec{w}}(\vec{x}) = f_{\vec{w}}(\vec{x})$ )

Squared loss:  $(f_{\vec{w}}(\vec{x}_i) - y_i)^2$  Ordinary Least Squares (OLS)

Absolute loss:  $|f_{\vec{w}}(\vec{x}_i) - y_i|$  not differentiable at 0

(Smooth AL) Huber loss:  $\begin{cases} \frac{1}{2} \varepsilon_i^2, & \text{if } |\varepsilon_i| < \delta \\ \delta(|\varepsilon_i| - \frac{\delta}{2}), & \text{otherwise} \end{cases}$  where  $\varepsilon_i = f_{\vec{w}}(\vec{x}_i) - y_i$ ,  $\delta$ -constant

Log-Cosh loss:  $\log(\cosh(f_{\vec{w}}(\vec{x}_i) - y_i))$  where  $\cosh = \frac{e^x + e^{-x}}{2}$

$\epsilon$ -Insensitive loss:  $\begin{cases} 0, & \text{if } |\varepsilon_i| < \epsilon \\ |\varepsilon_i| - \epsilon, & \text{otherwise} \end{cases}$  where  $\varepsilon_i = f_{\vec{w}}(\vec{x}_i) - y_i$

✓ 3. Regularizer  $\min_{\vec{w}} \sum_{i=1}^n l(\vec{w}^T \vec{x}_i + b, y_i) + \lambda r(\vec{w}) \Leftrightarrow \min_{\vec{w}} \sum_{i=1}^n l(\vec{w}^T \vec{x}_i + b, y_i)$  subject to  $r(\vec{w}) \leq B$

$L_1$ -regularizer:  $r(\vec{w}) = \vec{w}^T \vec{w} = \|\vec{w}\|_2^2 = \sum_i w_i^2$

$L_p$ -regularizer:  $r(\vec{w}) = \|\vec{w}\|_p = (\sum_i |w_i|^p)^{1/p}$

Elastic Net:  $\alpha \|\vec{w}\|_1 + (1-\alpha) \|\vec{w}\|_2^2, \alpha \in [0, 1]$   $(X^T \vec{w} - \vec{y})^T (X^T \vec{w} - \vec{y})$

4 Famous SRM model (Notation:  $X \in \mathbb{R}^{d \times n} = [\vec{x}_1; \dots; \vec{x}_n], \vec{y} = [y_1; \dots; y_n]$ )  
(DLS)

Ordinary Least Square:  $\min_{\vec{w}} \sum_{i=1}^n (\vec{w}^T \vec{x}_i - y_i)^2 \xrightarrow{\frac{\partial L}{\partial \vec{w}} = 0} \vec{w} = (X^T X)^{-1} X^T \vec{y}$

Ridge Regression:  $\min_{\vec{w}} \sum_{i=1}^n (\vec{w}^T \vec{x}_i - y_i)^2 + \lambda \|\vec{w}\|_2^2 \rightarrow \vec{w} = (X^T X + \lambda I)^{-1} X^T \vec{y}$

Lasso:  $\min_{\vec{w}} \frac{1}{n} \sum_{i=1}^n (\vec{w}^T \vec{x}_i - y_i)^2 + \lambda \|\vec{w}\|_1$

Logistic Regression:  $\min_{\vec{w}} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i(\vec{w}^T \vec{x}_i + b)})$  AD or Newton

## Lecture 2: Optimization for Machine Learning

2. Gradient Descent AD Algorithm: Initialize  $\vec{w}_0$

repeat  $\vec{w}_{t+1} \leftarrow \vec{w}_t + s$

until  $\|\vec{w}_{t+1} - \vec{w}_t\| \leq \epsilon$

+ ..

If  $\|s\|$  is small, Taylor expansion:  $L(\vec{w}+s) = L(\vec{w}) + \nabla L(\vec{w})^T s + \frac{1}{2} s^T V^2 L(\vec{w}) s$

Gradient descent use  $\nabla L(\vec{w})^T$   $L(\vec{w}+s) \approx L(\vec{w}_t) + \nabla L(\vec{w})^T s \leq L(\vec{w}_t)$

$\Rightarrow \nabla L(\vec{w})^T s \leq 0$  ( $g(\vec{w}) = \nabla L(\vec{w})^T$ )  $\Rightarrow s = -\alpha g(\vec{w}_t)$  learning rate or stepsize

3. Newton's Method

sufficient small - gd converge  
too large - gd diverge

$L(\vec{w}+s) = L(\vec{w}) + \nabla L(\vec{w})^T s + \frac{1}{2} s^T V^2 L(\vec{w}) s$   $L(\vec{w})$  is twice differentiable

$\nabla^2 L(\vec{w}) = H(\vec{w})$   $H(\vec{w}) = \frac{\partial^2 L}{\partial \vec{w}_i \partial \vec{w}_j}$  symmetric, positive semidefinite

choose  $s$ , we want  $L(\vec{w}+s)$  to be small

(\*)

$s = \arg \min_s (L(\vec{w}) + \nabla L(\vec{w})^T s + \frac{1}{2} s^T V^2 L(\vec{w}) s) = \arg \min_s L(\vec{w}) + g(\vec{w})^T s + \frac{1}{2} s^T H(\vec{w}) s$

$\frac{\partial (*)}{\partial s} = 0 \Rightarrow g(\vec{w}) + H(\vec{w})s = 0 \Rightarrow s = -H(\vec{w})^{-1}g(\vec{w})$

## Lecture 3: Estimating Probabilities from Data

✓ Discriminative learning, estimate  $p(y|\vec{x})$  directly

Generative learning, estimate  $p(\vec{x}, y) = p(\vec{x}|y)p(y)$

可用 Bayes rule 计算  $p(y|\vec{x}) = \frac{p(\vec{x}|y)p(y)}{p(\vec{x})}$

## 2. Maximum Likelihood Estimation (MLE)

(the goal is to choose  $\theta$  s.t. the observed data  $D$  is most likely)

MLE principle: Find  $\theta$  max likelihood  $P(D|\theta)$ :  $\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} P(D|\theta)$

(disadv: MLE overfit data if  $n$  small; x correct model, MLE terrible way)

## 3. The Bayesian Way and Maximum-a-posterior Estimation (MAP)

Bayes rules:  $p(\theta|D) = \frac{P(D|\theta)p(\theta)}{P(D)}$

MAP principle: Find  $\theta$  max  $p(\theta|D)$

$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|D) = \arg \max_{\theta} \frac{P(D|\theta)p(\theta)}{P(D)} = \arg \max_{\theta} \log P(D|\theta) + \log p(\theta)$

## 4. True Bayesian Approach

Posterior Mean:  $\hat{\theta}_{\text{mean}} = E(\theta|D) = \int_{\theta} \theta p(\theta|D) d\theta$

## Lecture 4: MLE and MAP for Discriminative Supervised Learning

(Discriminative learning, estimate  $p(y|\vec{x})$  directly)

✓ Assumption for Discriminative Supervised Learning:

(1)  $\vec{x}_i$  are known  $\rightarrow \vec{x}_i$  independent of the model parameters  $\vec{w}$   
 $\rightarrow p(x|\vec{w}) = p(x)$ ,  $p(\vec{w}|x) = p(\vec{w})$

(2)  $y_i$  are independent given the input feature  $\vec{x}_i$  and  $\vec{w}$   
 $p(y|\vec{x}, \vec{w})$

Goal: estimate  $\vec{w}$  from  $D = \{(\vec{x}_i, y_i)\}$  using joint conditional likelihood

Lemma: Max likelihood  $p(D|\vec{w}) = p(\vec{y}, X|\vec{w})$  is equivalent to max the joint conditional likelihood  $p(\vec{y}|X, \vec{w})$ .

Since  $p(\vec{y}, X|\vec{w}) = p(\vec{y}|X, \vec{w}) \cdot p(X|\vec{w}) = p(\vec{y}|X, \vec{w}) \cdot p(X)$   $\leftarrow$  not depends on  $\vec{w}$

✓ MLE: choose  $\vec{w}$  to max  $p(\vec{y}|X, \vec{w})$  log-likelihood

$$\hat{\vec{w}}_{MLE} = \underset{\vec{w}}{\operatorname{argmax}} p(\vec{y}|X, \vec{w}) \stackrel{(1)}{=} \underset{\vec{w}}{\operatorname{argmax}} \prod_{i=1}^n p(y_i|\vec{x}_i, \vec{w}) = \underset{\vec{w}}{\operatorname{argmax}} \sum_{i=1}^n \log p(y_i|\vec{x}_i, \vec{w})$$

✓ MAP: choose  $\vec{w}$  to max  $p(\vec{w}|X, \vec{y})$

$$\hat{\vec{w}}_{MAP} = \underset{\vec{w}}{\operatorname{argmax}} p(\vec{w}|X, \vec{y}) = \underset{\vec{w}}{\operatorname{argmax}} P(X, \vec{y}|\vec{w}) p(\vec{w})$$

$$= \underset{\vec{w}}{\operatorname{argmax}} P(\vec{y}|X, \vec{w}) p(\vec{w}) \underset{\text{Likelihood prior}}{=} \underset{\vec{w}}{\operatorname{argmax}} \sum_{i=1}^n \log p(y_i|\vec{x}_i, \vec{w}) + \log p(\vec{w})$$

Lecture 5: Naive Bayes Classifier  $P(\vec{x}|y, D)$

$$\text{Bayes rule } P(y|\vec{x}) = \frac{P(\vec{x}|y)P(y)}{P(\vec{x})} \stackrel{(1)}{=} P(y|D)$$

## 2. Naive Bayes

### 2.1 Estimate $P(y)$ counting

binary classification  $\hat{\pi}_{y=+1} = P(y=+1)$ ,  $\hat{\pi}_{y=-1} = P(y=-1)$

multi-class classification  $\hat{\pi}_c = P(y=c) = \frac{1}{n} \sum_{i=1}^n I(y_i=c)$

### 2.2 Estimate $P(\vec{x}|y)$

✓ Naive Bayes Assumption: features are independent given the label

$$P(\vec{x}|y) = \prod_{\alpha=1}^d P(\vec{x}_\alpha|y)$$

$$y = h(\vec{x}) = \underset{y}{\operatorname{argmax}} P(y|\vec{x}) = \underset{y}{\operatorname{argmax}} \frac{P(\vec{x}|y)P(y)}{P(\vec{x})}$$

$$= \underset{y}{\operatorname{argmax}} P(\vec{x}|y)P(y) = \underset{y}{\operatorname{argmax}} \prod_{\alpha=1}^d P(\vec{x}_\alpha|y)P(y)$$

Case 1: Categorical Features; Case 2: Multinomial Features

Case 3: Continuous Features (Gaussian Naive Bayes)

$$P([\vec{x}]_\alpha | y=c) = N(\mu_{yc}, \sigma_{yc}^2) = \frac{1}{\sqrt{2\pi\sigma_{yc}^2}} \exp\left(-\frac{(x - \mu_{yc})^2}{2\sigma_{yc}^2}\right)$$

$$\hat{\mu}_{yc} \leftarrow \frac{1}{n_c} \sum_{i=1}^n I(y_i=c) \vec{x}_i ; \sigma_{yc}^2 \leftarrow \frac{1}{n_c} \sum_{i=1}^n I(y_i=c) (\vec{x}_i - \hat{\mu}_{yc}) (\vec{x}_i - \hat{\mu}_{yc})^\top$$

$$\text{where } n_c = \sum_{i=1}^n I(y_i=c)$$

$$\hat{\mu}_c \leftarrow \left( \begin{array}{c} \hat{\mu}_{1c} \\ \vdots \\ \hat{\mu}_{dc} \end{array} \right)$$

$$\text{without naive assumption } \hat{\Sigma}_c \leftarrow \left( \begin{array}{cc} \hat{\sigma}_{1c}^2 & \hat{\sigma}_{1c} \hat{\sigma}_{2c} \\ \hat{\sigma}_{2c} & \hat{\sigma}_{2c}^2 \end{array} \right)$$

$$\text{with naive assumption } \hat{\mu}_c \leftarrow \left( \begin{array}{c} \hat{\mu}_{1c} \\ \vdots \\ \hat{\mu}_{dc} \end{array} \right) ; \hat{\Sigma}_c = \left( \begin{array}{cc} \hat{\sigma}_{1c}^2 & 0 \\ 0 & \hat{\sigma}_{2c}^2 \end{array} \right)$$

✓ NB is a linear classifier  $\rightarrow h(\vec{x}) = +1 \Leftrightarrow \vec{w}^\top \vec{x} + b > 0$

Suppose features - binary  $h(\vec{x}) = \underset{y \in \{-1, +1\}}{\operatorname{argmax}} P(y=c) \stackrel{(1)}{=} P(\vec{x}_1|y=c) = \text{sign}(\vec{w}_1^\top \vec{x} + b_{1,y})$

features - multinomial  $h(\vec{x}) = \underset{y \in \{-1, +1\}}{\operatorname{argmax}} P(y=c) P(\vec{x}|m, y=c) = \text{sign}(\vec{w}_1^\top \vec{x} + b_{1,y})$

feature-continuous  $P(y=c|\vec{x}) = \frac{1}{1 + \exp(-y(\vec{w}_1^\top \vec{x} + b_{1,y}))}$

logodds:  $\log \frac{P(y=1|\vec{x})}{P(y=-1|\vec{x})} = \vec{w}^\top \vec{x}$  (logistic regression)

$$\text{Multinomial: } P(X_1=x_1, \dots, X_k=x_k) = \frac{n!}{x_1! \dots x_k!} P_1^{x_1} \dots P_k^{x_k}$$

$$\text{Gaussian: } f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$\text{Bernoulli: } P(X=x|p) = \begin{cases} p, & x=1 \\ 1-p, & x=0 \end{cases}$$

$$P(X=x) = p^x (1-p)^{1-x}, x \in \{0, 1\} \rightarrow \frac{n!}{k!(n-k)!} \text{ 次方次数}$$

$$\text{Binomial Distribution: } P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}, k=0, 1, \dots, n$$

$(P(w|D) \propto P(D|w)P(w)$ , Bayes rule)

## Lecture A1: Unsupervised Learning - Clustering

1.2 Clustering Given:  $D = \{\vec{x}_1, \dots, \vec{x}_n\} \subset \mathbb{R}^d$

### 2. K-means clustering

1) number of clustering  $K$  is known

Assumption: 2) each input can only belong to a single cluster

Notation: 1)  $\vec{z}_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \in \mathbb{R}^K$  one-hot coding, 每个点对应一个  $\vec{z}_i$

如果数据点  $x_i$  归属于簇  $k$ , 则  $[z_i]_k = 1$ ,

### 2) cluster center (also cluster mean)

$$\hat{\mu}_\alpha = \frac{\sum_{i=1}^n [z_i]_\alpha \vec{x}_i}{\sum_{i=1}^n [z_i]_\alpha} \leftarrow \text{属于簇的所有点的坐标和} / \text{属于簇的所有点的个数}$$

$$\text{distance} \quad d_{ik} = \|\vec{x}_i - \hat{\mu}_k\|^2$$

### 2.1 Objective and Algorithm

$$\text{objective function: } D(\vec{\mu}_1, \dots, \vec{\mu}_K, z_1, \dots, z_n) = \sum_{i=1}^n \sum_{\alpha=1}^K [z_i]_\alpha d_{ik}$$

### K-means clustering Algorithm

Initialize  $\vec{\mu}_1, \dots, \vec{\mu}_K$  to distinct data points in  $D$  计算距离  
判断属于哪簇

repeat 所有数据点  
for all  $i=1, \dots, n$  do  $[z_i]_\alpha = \begin{cases} 1, & \text{if } \alpha = \operatorname{argmin}_\alpha \|\vec{x}_i - \hat{\mu}_\alpha\|^2 \\ 0, & \text{otherwise} \end{cases}$   
end for  
for all  $\alpha=1, \dots, K$  do  $\hat{\mu}_\alpha = \frac{\sum_{i=1}^n [z_i]_\alpha \vec{x}_i}{\sum_{i=1}^n [z_i]_\alpha}$   
end for  
until converge ( $z_i^{(t)} = z_i^{(t+1)}$  or  $\hat{\mu}_\alpha^{(t)} = \hat{\mu}_\alpha^{(t+1)}$ )

### 2.3 choosing the number of clusters $K$ — kink

### 3. Gaussian Mixture Model Clustering

Define  $[z_i]_\alpha \in [0, 1]$  as the probability that  $\vec{x}_i$  belongs to cluster  $\alpha$ .

Assumption: data set  $D$  comes from a mixture of Gaussian

$$\Rightarrow \vec{x}_i \sim N(\hat{\mu}_\alpha, \hat{\Sigma}_\alpha) \leftarrow \text{each cluster } \alpha$$

### 3.2 Generative Model and GMM Algorithm

Initialize  $\vec{\mu}_1, \dots, \vec{\mu}_K$  to be distinct random data points in  $D$

Initialize  $\Sigma_1, \dots, \Sigma_K$  to  $\sigma^2 I$  for any  $\sigma > 0$

Initialize  $\pi_1, \dots, \pi_K$  to  $\frac{1}{K}$  (uniform prior)

repeat

$$\text{E-step: } [\vec{z}_i]_\alpha = \frac{P(\vec{x}_i|\hat{\mu}_\alpha, \hat{\Sigma}_\alpha) \pi_\alpha}{\sum_{\alpha=1}^K P(\vec{x}_i|\hat{\mu}_\alpha, \hat{\Sigma}_\alpha) \pi_\alpha}$$

end for

$$\text{M-step: } \hat{\mu}_\alpha = \frac{\sum_{i=1}^n [\vec{z}_i]_\alpha \vec{x}_i}{\sum_{i=1}^n [\vec{z}_i]_\alpha} ; \hat{\Sigma}_\alpha = \frac{\sum_{i=1}^n [\vec{z}_i]_\alpha (\vec{x}_i - \hat{\mu}_\alpha)(\vec{x}_i - \hat{\mu}_\alpha)^\top}{\sum_{i=1}^n [\vec{z}_i]_\alpha}$$

end for

$$\pi_\alpha = \frac{1}{n} \sum_{i=1}^n [\vec{z}_i]_\alpha$$

missing features  $y = \underset{c}{\operatorname{argmax}} \prod_{\alpha=1}^K P(\vec{x}_\alpha|y=c) = \underset{c}{\operatorname{argmax}} \prod_{\alpha=1}^K \prod_{\alpha \in \text{OBS}} P(\vec{x}_\alpha|y=c)$

$$\frac{\partial(\vec{w}^\top \vec{w})}{\partial \vec{w}} = 2A\vec{w} \quad \frac{\partial(\vec{w}^\top B)}{\partial \vec{w}} = B^\top \quad \text{sign}(x) = \begin{cases} +1, & x > 0 \\ 0, & x = 0 \\ -1, & x < 0 \end{cases}$$

$P(\vec{w}|X, \vec{y})$  posterior distribution (over parameters)

$P(y|x, w)$  model  $P(y^*|x^*, w)$  (posterior) predictive distribution

$P(\vec{y}|X, \vec{w})$  conditional likelihood (i.i.d.  $= \prod_{i=1}^n P(y_i|\vec{x}_i, \vec{w})$ )

$P(D|\vec{w})$  likelihood, MLE  $\vec{w} = \operatorname{argmax}_{\vec{w}} P(D|\vec{w})$

$P(w)$  prior distribution (over parameters  $w$ )