



# CS M146 Discussion: Week 8 Ensemble Method, Multi-Class Classification, ML Evaluation

Junheng Hao Friday, 02/26/2021







- Announcement
- Ensemble Method
- Multi-Class Classification
- ML Evaluation





- **5:00 pm PST, Feb 26 (Friday):** Weekly Quiz 8 released on Gradescope.
- **11:59 pm PST, Feb 28 (Sunday):** Weekly Quiz 8 closed on Gradescope!
  - Start the quiz before **11:00 pm PST, Feb 28** to have the full 60-minute time
- **Problem set 3** released on CCLE, submission on Gradescope.
  - Please assign pages of your submission with corresponding problem set outline items on GradeScope.
  - You need to submit code and the results required by the problem set
  - Due on Today 11:59pm PST, Feb 26 (Friday)!

Late Submission of PS will NOT be accepted!





- Quiz release date and time: Feb 26, 2021 (Friday) 05:00 PM PST
- Quiz due/close date and time: Feb 28, 2021 (Sunday) 11:59 PM PST
- You will have up to **60 minutes** to take this exam. → Start before **11:00 PM** Sunday
- You can find the exam entry named "Week 8 Quiz" on GradeScope.
- Topics: Ensemble Method, Multi-Class Classification, ML Evaluations
- Question Types
  - True/false, multiple choices
  - Some questions may include several subquestions.
- Some light calculations are expected. Some scratch paper and one scientific calculator (physical or online) are recommended for preparation.



# Quiz 7 Review: Kernel SVM



#### Q6 Kernel SVM

2 Points

We can introduce non-linearity to SVM using the kernel trick. Instead of searching for a hyperplane  $\mathbf{w}^T \mathbf{x} + b$  that maximizes the margin, we are looking for  $\mathbf{w}^T \phi(\mathbf{x}) + b$  where  $\phi$  is the non-linear basis function.

Which one of the following statements is wrong about kernel SVM?

 ${old O}$  We can learn the optimal value of the weights  ${f w}$  using only the kernel function.

O We can predict the label of a new sample using only the kernel function.

O If we apply kernel functions, non-separable data may be separable.

O A valid kernel function should have a positive-semidefinite kernel matrix.

	Linear SVM $oldsymbol{w}^{ op}oldsymbol{x}+b$	Kernel SVM $oldsymbol{w}^{\mathrm{T}}oldsymbol{\phi}(oldsymbol{x}_n)+b$
Weight parameter	$oldsymbol{w} = \sum_n lpha_n oldsymbol{y}_n oldsymbol{x}_n$	$oldsymbol{w} = \sum_n lpha_n y_n oldsymbol{\phi}(oldsymbol{x}_n)$
Predicting new data	$ ext{SIGN}ig(\sum_n y_n lpha_nig(oldsymbol{x}_n^Toldsymbol{x}ig)+big)$	$ ext{SIGN}(\sum_n y_n lpha_n k(oldsymbol{x}_n,oldsymbol{x})+b)$





## Q4 SVM on non-separable data

2 Points

# Which **one** of the following statements is **wrong** about SVM applied to non-separable data (soft-margin SVM)?

- O The use of slack variables permits samples to be misclassified.
- O As we increase the value of the hyperparameter C, the margin of the learned hyperplane decreases.
- O As we decrease the value of the hyperparameter C, the training accuracy of the learned SVM may decrease.
- During training, we do not only minimize  $\frac{1}{2} \|\omega\|_2^2$ , but also minimize  $\sum (1 \xi_n)$  so that the value of slack variable can be controlled.



- The C parameter tells the SVM optimization how much you want to avoid misclassifying each training example.
- For large values of C, the optimization will choose a smaller-margin hyperplane if that hyperplane does a better job of getting all the training points classified correctly.
- Conversely, a very small value of C will cause the optimizer to look for a larger-margin separating hyperplane, even if that hyperplane misclassified more points.

















- What if C is super large, i.e.  $C \rightarrow +inf$ ?
- What is C is super small, i.e.  $C \rightarrow 0$  (a very small positive number, such as 0.0001)?



# Additional: SVM on 1-dim data (code)



[1] from sklearn import svm import numpy as np

- [2] X = [[-3], [-2], [-1], [0], [1], [2], [3]] y = [0, 0, 0, 1, 1, 1, 1] X, y = np.array(X), np.array(y)

print\_svm\_decision(X, y, C=1)

□ Support Vectors are: [2 3] Predicting y=wx+b: [-2.36363017 -1.45454339 -0.54545661 0.36363017 1.27271694 2.18180372 3.0908905 ]

[5] print\_svm\_decision(X, y, C=0.5)

Support Vectors are: [2 3 4] Predicting y=wx+b: [-1.75000807 -1.08333882 -0.41666956 0.24999969 0.91666894 1.5833382 2.2500745]

[6] print\_svm\_decision(X, y, C=0.1)

Support Vectors are: [1 2 3 4 5] Predicting y=wx+b: [-1.09999783 -0.69999856 -0.299999929 0.09999998 0.49999925 0.89999852 1.29999779]

Case: C is suffciently small and every vector is within the margin, i.e. every data point is a support vector.

[7] print\_svm\_decision(X, y, C=0.01)

Support Vectors are: [0 1 2 3 4 5 6] Predicting y=wx+b: [-0.44399465 -0.29014848 -0.13630231 0.01754386 0.17139003 0.3252362 0.47908237]

Colab link:

https://colab.research.google.com/driv e/1Ru\_gN8UikD\_fGY3DfarHTfXQQDg4 b-D4?usp=sharing

# Ensemble: Bagging and Boosting

UCLA

**Engineer Change.** 







• "Multiple" dataset and multiple classifier



#### **Bagging Classifier Process Flow**





• **Single:** Decision Tree → **Bagging:** Random Forest







- Given: N samples  $\{x_n, y_n\}$ , where  $y_n \in \{+1, -1\}$ , and some way of constructing weak (or base) classifiers
- Initialize weights  $w_1(n) = \frac{1}{N}$  for every training sample
- For t = 1 to T

Engineer Change.

**1** Train a weak classifier  $h_t(x)$  using current weights  $w_t(n)$ , by minimizing the weighted classification error

$$\epsilon_t = \sum_n w_t(n) \mathbb{I}[y_n 
eq h_t(oldsymbol{x}_n)]$$

1. Calculating error

2 Compute contribution for this classifier: β<sub>t</sub> = <sup>1</sup>/<sub>2</sub> log <sup>1-ε<sub>t</sub></sup>/<sub>ε<sub>t</sub></sub>
 3 Update weights on training points

$$w_{t+1}(n) \propto w_t(n) e^{-\beta_t y_n h_t(\boldsymbol{x}_n)}$$

and normalize them such that  $\sum_{n} w_{t+1}(n) = 1$ • Output the final classifier

$$h[oldsymbol{x}] = \mathsf{sign}\left[\sum_{t=1}^T eta_t h_t(oldsymbol{x})
ight]$$

2. Calculating classifier weights of t

3. Reweighting training points

## **Bagging Classifier**





• Exponential loss, instead of 0/1 loss







Given C classes and N data points each class:

Comparison	One-vs-one	One-vs-rest (all)
# Binary Classifies		
# Training Data		
Pros		
Cons		











### Model

For each class  $C_k$ , we have a parameter vector  $\theta_k$  and model the probability of class  $C_k$  as

Model:

$$p(y = C_k | \boldsymbol{x}; \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K) = rac{e^{\boldsymbol{\theta}_k^{\mathrm{T}} \boldsymbol{x}}}{\sum_{k'} e^{\boldsymbol{\theta}_{k'}^{\mathrm{T}} \boldsymbol{x}}} \qquad \leftarrow$$

 $\leftarrow \quad \text{This is called } softmax$ 

**Decision boundary**: assign  $\boldsymbol{x}$  with the label that is the maximum of

$$\arg \max_k P(y = C_k | \boldsymbol{x}; \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K) \to \arg \max_k \boldsymbol{\theta}_k^{\mathrm{T}} \boldsymbol{x}_k$$

$$\sum_{n} \log P(y_n | \boldsymbol{x}_n; \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K) = \sum_{n} \log \prod_{k=1}^{K} P(y = C_k | \boldsymbol{x}_n; \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K)^{y_{nk}}$$
$$= \sum_{n} \sum_{k} y_{nk} \log P(y = C_k | \boldsymbol{x}_n; \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K)$$

Likelihood





Likelihood

**Cost Function** 

$$\sum_{n} \log P(y_{n} | \boldsymbol{x}_{n}; \boldsymbol{\theta}_{1}, \dots, \boldsymbol{\theta}_{K}) = \sum_{n} \log \prod_{k=1}^{K} P(y = C_{k} | \boldsymbol{x}_{n}; \boldsymbol{\theta}_{1}, \dots, \boldsymbol{\theta}_{K})^{y_{nk}}$$
$$= \sum_{n} \sum_{k} y_{nk} \log P(y = C_{k} | \boldsymbol{x}_{n}; \boldsymbol{\theta}_{1}, \dots, \boldsymbol{\theta}_{K})$$
$$J(\boldsymbol{\theta}_{1}, \boldsymbol{\theta}_{2}, \dots, \boldsymbol{\theta}_{K}) = -\sum_{n} \sum_{k} y_{nk} \log P(y = C_{k} | \boldsymbol{x}_{n}; \boldsymbol{\theta}_{1}, \boldsymbol{\theta}_{2}, \dots, \boldsymbol{\theta}_{K})$$
$$= -\sum_{n} \sum_{k} y_{nk} \log \left(\frac{e^{\boldsymbol{\theta}_{k}^{\mathrm{T}} \boldsymbol{x}}}{\sum_{k'} e^{\boldsymbol{\theta}_{k'}^{\mathrm{T}} \boldsymbol{x}}}\right)$$

 $= -\sum_n \sum_k y_{nk} \boldsymbol{\theta}_k^{\mathrm{T}} \boldsymbol{x} - y_{nk} \log \left( \sum_{k'} e^{\boldsymbol{\theta}_{k'}^{\mathrm{T}} \boldsymbol{x}} \right)$ 

Optimization

 $\mathsf{Convex} \to \mathsf{SGD}$ 



# Multiclass Classification in Neural Nets







5 separate **binary classifiers** Key: **sharing the same hidden layers** with **different weights at the end** 

https://developers.google.com/machine-learning/crash-course/multi-class-neural-networks/one-vs-all http://www.briandolhansky.com/blog/2013/9/23/artificial-neural-nets-linear-multiclass-part-3



# Softmax Layer in Neural Nets



#### Multi-Class Classification with NN and SoftMax Function



$$p(y=j|\mathbf{x}) = rac{e^{(\mathbf{w}_j^T\mathbf{x}+b_j)}}{\sum_{k\in K}e^{(\mathbf{w}_k^T\mathbf{x}+b_k)}}$$

→ Implementation in PyTorch?

Credit: https://developers.google.com/machine-learning/crash-cour

se/multi-class-neural-networks/softmax



# **Evaluation: Binary Classifier**



• Diagnostic testing table

		True condi	ition			
	Total population	Condition positive	Condition negative	$Prevalence = \frac{\Sigma \text{ Condition positive}}{\Sigma \text{ Total population}}$	$\frac{Accuracy (ACC) =}{\frac{\Sigma True \ positive + \Sigma \ True \ negative}{\Sigma \ Total \ population}}$	
condition	Predicted condition positive	True positive	False positive, Type I error	Positive predictive value (PPV), PrecisionFalse discovery rate ( $\Sigma$ False positive= $\frac{\Sigma}{\Sigma}$ True positive $\Sigma$ False positive $\Sigma$ Predicted condition positive $\Sigma$ Predicted condition		
Predicted	Predicted condition negative	False negative, Type II error	True negative	False omission rate (FOR) = $\Sigma$ False negative $\Sigma$ Predicted condition negative	Negative predictive value (NPV) = $\Sigma$ True negative $\Sigma$ Predicted condition negative	
		True positive rate (TPR), Recall, Sensitivity,probability of detection, Power = $\frac{\Sigma \text{ True positive}}{\Sigma \text{ Condition positive}}$	False positive rate (FPR), Fall-out,probability of false alarm= $\frac{\Sigma}{\Sigma}$ False positive $\Sigma$ Condition negative	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds $F_1$ score =	
		$\frac{\text{False negative rate (FNR)}}{\Sigma \text{ False negative}} \text{ Miss rate}$ $= \frac{\Sigma \text{ False negative}}{\Sigma \text{ Condition positive}}$	$\frac{\text{Specificity (SPC), Selectivity, True negative}}{\text{rate (TNR)} = \frac{\Sigma \text{ True negative}}{\Sigma \text{ Condition negative}}$	Negative likelihood ratio (LR–) = $\frac{FNR}{TNR}$	2 · Precision + Recall	









## **Evaluation: Binary Classifier**



• Precision and recall



Credit: https://en.wikipedia.org/wiki/F-score



## Evaluation: Example



• Calculation from confusion matrix



	Predicted class POSITIVE (spam 🖂 )	Predicted class NEGATIVE (normal 鹵 )
Actual class POSITIVE (spam 🖂 )	TRUE POSITIVE (TP)	FALSE NEGATIVE (FN) 전 원 43
Actual class NEGATIVE (normal ⊠ )	FALSE POSITIVE (FP)	TRUE NEGATIVE (TN) 回 函 538

Credit: https://towardsdatascience.com/confusion-matrix-and-class-statistics-68b79f4f510b





UCLA

**Engineer Change.** 

# Evaluation: Multi-class Confusion Matrix







- ROC [Wiki]
- Area under the ROC Curve: AUC/AUROC
- A random classifier has AUROC =  $\frac{1}{2}$ .
- A perfect classifier has AUROC = 1.
- Others: Precision-Recall Curve and AUPR

In-class Exercise: Building ROC Curves







A		В		С		_	C'				
TP=63	FP=28	91	TP=77	FP=77	154	TP=24	FP=88	112	TP=76	FP=12	88
FN=37	TN=72	109	FN=23	TN=23	46	FN=76	TN=12	88	FN=24	TN=88	112
100	100	200	100	100	200	100	100	200	100	100	200
TPR = 0.6	3		TPR = 0.7	7		TPR = 0.2	24		TPR = 0.7	76	
FPR = 0.28 FPR = 0.77			FPR = 0.88			FPR = 0.12					
PPV = 0.69 PPV = 0.50			PPV = 0.21			PPV = 0.86					
F1 = 0.66		F1 = 0.61		F1 = 0.23			F1 = 0.81				
ACC = 0.68		ACC = 0.50			ACC = 0.18			ACC = 0.82			

---- Random guess Pefect Classification 0.9 0.8 ٠ C 0.7 • TPR or sensitivity A Better 0.3 ċ 0.2 Worse 0.1 0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1 FPR or (1 - specificity)

**ROC Space** 









- Official Doc (stable version):
   <u>https://scikit-learn.org/stable/modules/model\_evaluation.html</u>
- Some helpful demos:
  - Precision-Recall:

https://scikit-learn.org/stable/auto\_examples/model\_selection/plot\_precision\_recall.html#sp hx-glr-auto-examples-model-selection-plot-precision-recall-py

• Confusion matrix:

https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion\_matrix.html#skl earn.metrics.confusion\_matrix

 Receiver Operating Characteristic (ROC): <u>https://scikit-learn.org/stable/auto\_examples/model\_selection/plot\_roc.html#sphx-glr-auto-examples-model-selection-plot-roc-py</u>



## Whiteboard







# Thank you!