PROPERTY AND AND AND AND

TREO PROCESSION

AT 424815 1 miles that a baseline of the second s

Machine Learning Basics

Prof. Kuan-Ting Lai 2021/3/11 THE COURT OF ANY ANY THE TOUS ANY CONTRACTOR

FERRE STATE IN

A CONTRACTOR OF THE OWNER

PUBLIC CLAIR CLARNES

Mare state coolars

Chief in all the

Constant and the second second

Cambrid and a straight a branches a

The Grand of Station of Station

Machine Learning



Francois Chollet, "Deep Learning with Python," Manning, 2017

Machine Learning Flow











scikit-learn.org



Home Installation Documentation - Examples



scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- · Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- · Open source, commercially usable BSD license

Classification

Identifying to which category an object belongs to.

Applications: Spam detection, Image recognition. Algorithms: SVM, nearest neighbors, random forest, ... — Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency Algorithms: PCA, feature selection, nonnegative matrix factorization. — Examples

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices. Algorithms: SVR, ridge regression, Lasso, — Examples

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning Modules: grid search, cross validation, metrics. — Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes Algorithms: k-Means, spectral clustering, mean-shift, ... — Examples

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms. Modules: preprocessing, feature extraction.

Types of Data

Data Types (Measurement Scales)



https://towardsdatascience.com/data-types-in-statistics-347e152e8bee

10

Nominal Data (Labels)

- Nominal data are labeling variables without any quantitative value
- Encoded by one-hot encoding for machine learning
- Examples:

What is your Gender?	What languages do you speak?
O Female	O Englisch
O Male	O French
	O German
	O Spanish

Ordinal Data

- Ordinal values represent discrete and ordered units
- The order is meaningful and important

What Is Your Educational Background?

- 1 Elementary
- 🔘 2 High School
- 3 Undegraduate
- 🔵 4 Graduate

Interval Data

- Interval values represent ordered units that have the same difference
- Problem of Interval: Don't have a true zero
- Example: Temperature Celsius (°C) vs. Fahrenheit (°F)

Temperature? -10
-5
0
+5
+10
+15



- Same as interval data but have absolute zero
- Can be applied to both descriptive and inferential statistics
- Example: weight & height





Machine Learning vs. Statistics

<u>https://www.r-bloggers.com/whats-the-difference-between-machine-learning-statistics-and-data-mining/</u>

Machine learning	Statistics
network, graphs	model
weights	parameters
learning	fitting
generalization	test set performance
supervised learning	regression/classification
unsupervised learning	density estimation, clustering
large grant = \$1,000,000	large grant = \$50,000
nice place to have a meeting: Snowbird, Utah, French Alps	nice place to have a meeting: Las Vegas in August

Supervised and Unsupervised Learning



Iris Flower Classification (鳶尾花分類)



Iris Versicolor

Iris Setosa

Iris Virginica

Extracting Features of Iris (抽取特徵值)

• Width and Length of Petal (花瓣) and Sepal (花萼)



Iris Flower Dataset



Classify Iris Species via Petals and Sepals

 Iris versicolor and virginica are not linearly separable



https://www.tensorflow.org/tutorials/customization/custom_training_walkthrough



Evaluation (Loss Function)

 $y' = w^T x - b < 0$ x_2 $\Delta = \mathcal{Y} - \mathcal{Y}'$ Learning Rule (Perceptron) $W \leftarrow W + \alpha(y - y') X$ Learning Rate



22,

Support Vector Machine (SVM)

• Choose the hyperplanes that have the largest separation (margin)



Loss Function of SVM

• Calculate prediction errors

$$y'_{i} = W^{T}x_{i} - b \ge 1$$

$$y'_{i} = W^{T}x_{i} - b \le -1$$

$$y_{i}(W^{T}x_{i} - b) \ge 1$$

$$Loss = max[0, 1 - y_{i}(W^{T}x_{i} - b)]$$

Hinge Loss



SVM Optimization

- Maximize the margin while reduce hinge loss
- Hinge loss: $\max(0, 1 y_i(\vec{w} \cdot \vec{x}_i b)) x_2 \uparrow$

mīn. || W || s.b.t. max(o, | - y:(wx:-b))



Nonlinear Problem?

• How to separate Versicolor and Virginica?



SVM Kernel Trick

• Project data into higher dimension and calculate the inner products



https://datascience.stackexchange.com/questions/17536/kernel-trick-explanation

Nonlinear SVM for Iris Classification



Using Neural Network



https://www.tensorflow.org/tutorials/customization/custom_training_walkthrough

Supervised and Unsupervised Learning



Linear Regression (Least squares)

• Find a "line of best fit" that minimizes the total of the square of the errors



Supervised and Unsupervised Learning



Logistic Regression

Sigmoid function

$$S(x) = \frac{e^x}{e^x + 1} = \frac{1}{1 + e^{-x}}$$

• Derivative of Sigmoid

S(x) = S(x)(1 - S(x))



https://en.wikipedia.org/wiki/Sigmoid_function

Decision Boundary

• Binary classification with decision boundary t

$$y' = P(x, w) = P_{\theta}(x) = \frac{1}{1 + e^{-(w^T x + b)}}$$

$$y' = \begin{cases} 0, & x < t \\ 1, & x \ge t \end{cases}$$



Cross Entropy Loss

Loss function: cross entropy

$$loss = \begin{cases} -\log(1 - P_{\theta}(x)), & \text{if } y = 0\\ -\log(P_{\theta}(x)), & \text{if } y = 1 \end{cases}$$

$$\Rightarrow L_{\theta}(\mathbf{x}) = -y \log(P_{\theta}(\mathbf{x})) + -(1-y)\log(1-P_{\theta}(\mathbf{x}))$$

$$\nabla L_W(\mathbf{x}) = -(\underline{y - P_\theta(x)})x$$

Machine Learning Workflow



https://towardsdatascience.com/workflow-of-a-machine-learning-project-ec1dba419b94

Į.

Overfitting and Underfitting

Overfitting

Underfitting



https://en.wikipedia.org/wiki/Overfitting

Overfitting (以偏概全)

• Overfitting is common, especially for neural networks



Neural Network Urban Legend: Detecting Tanks

• Detector learned the illumination of photos





Bias and Variance Trade-off

• Model with high variance overfits to training data and does not generalize on unseen test data



Model Selection



42

Training, Validation, Testing

- Never leak test data information into our model
- Tuning the *hyperparameters* of our model on validation dataset



K-Fold Cross Validation

• Lower the variance of validation set



Regularization

 <u>https://developers.google.com/machine-learning/crash-</u> course/regularization-for-sparsity/l1-regularization

5

п

ÞI



Metrics: Accuracy vs. Precision in Binary Classification

Confusion Matrix

		True co	ondition
	Total population	Condition positive	Condition negative
Predicted	Predicted condition positive	True positive , Power	False positive, Type I error
condition	Predicted condition negative	False negative, Type II error	True negative

https://en.wikipedia.org/wiki/Confusion_matrix

Confusion Matrix

https://en.wikipedia.org/wiki/Confusion_matrix

		True cond	dition			
	Total population	Condition positive	Condition negative	$\frac{\text{Prevalence}}{\sum \text{ Condition positive}} = \frac{\Sigma \text{ Condition positive}}{\Sigma \text{ Total population}}$	Accu <u>Σ True positi</u> Σ Το	iracy (ACC) = ive + Σ True negative tal population
Predicted	Predicted condition positive	True positive	False positive, Type I error	Positive predictive value (PPV), Precision = Σ True positive $\overline{\Sigma}$ Predicted condition positive	False disc Σ False disc Σ Predicte	overy rate (FDR) = alse positive d condition positive
condition	Predicted condition negative	False negative , Type II error	True negative	False omission rate (FOR) = Σ False negative Σ Predicted condition negative	Negative pre Σ Τ Σ Predicted	edictive value (NPV) = rue negative d 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 = $\frac{\Sigma \text{ False positive}}{\Sigma \text{ Condition negative}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio	F ₁ score =
		False negative rate (FNR), Miss rate = $\frac{\Sigma \text{ False negative}}{\Sigma \text{ Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) = $\frac{\Sigma \text{ True negative}}{\Sigma \text{ Condition negative}}$	Negative likelihood ratio (LR-) = <u>FNR</u> TNR	$(DOR) = \frac{LR+}{LR-}$	2 · Precision · Recall Precision + Recall

Popular Metrics

Notations

–P: positive samples, N: negative samples, P': predicted positive samples, TP: true positives, TN: true negatives

• Recall = $\frac{TP}{P}$ • Precision = $\frac{TP}{P'}$ • Accuracy = $\frac{TP}{P'+TN}$

• Accuracy =
$$\frac{P+N}{2}$$

• F1 score =
$$\frac{1}{\frac{1}{recall} + \frac{1}{precision}}$$

• Miss rate = false negative rate = 1 – recall

Coronavirus Example

- Precision = 8 / 18 = 44%
- Accuracy = (8 + 90) / 110 = 89%



一種錯誤率, 各自表述!?

這幾天有新聞報導,捷克使用中國製的快篩劑來檢驗新冠病毒 結果錯誤率高達 80%!





Accuracy

• Example: Classifying tumors as malignant or benign (Google Machine Learning Crash Course)

True Positive (TP):	False Positive (FP):
Reality: Malignant	Reality: Benign
ML model predicted: Malignant	ML model predicted: Malignant
Number of TP results: 1	Number of FP results: 1
False Negative (FN):	True Negative (TN):
False Negative (FN): • Reality: Malignant	True Negative (TN): Reality: Benign
False Negative (FN): • Reality: Malignant • ML model predicted: Benign	True Negative (TN): Reality: Benign ML model predicted: Benign

 $\label{eq:accuracy} {\rm Accuracy} = \frac{TP+TN}{TP+TN+FP+FN} = \frac{1+90}{1+90+1+8} = 0.91$

https://developers.google.com/machine-learning/crash-course/classification/accuracy

Precision and Recall

• Classifying tumors as malignant or benign



https://developers.google.com/machine-learning/crash-course/classification/precision-and-recall

ROC Curve

• ROC (Receiver Operating Characteristic) Curve

True Positive Rate (TPR) is a synonym for recall and is therefore defined as follows:

$$TPR = \frac{TP}{TP + FN}$$

False Positive Rate (FPR) is defined as follows:







https://bwfinsight.blog/2018/05/10/the-link-betweenworld-war-ii-and-your-predictive-models/ 53

AUC (Area under the ROC Curve)

• AUC ranges in value from 0 to 1.

- AUC of 0.0 = predictions are 100% wrong (X)

– AUC of 1.0 = predictions are 100% correct (O)

Advantages

- Scale-invariant
- Classification-threshold-invariant
 - Not suitable for email spam detection



Evaluate Decision Boundary t

 ROC (Receiver Operating Characteristic) Curve



• Precision-Recall (PR) Curve



Summary of ML Training Flow

- 1. Defining the problem and assembling a dataset
- 2. Choosing a measure of success
- 3. Deciding on an evaluation protocol
- 4. Preparing your data
- 5. Developing a model that does better than a baseline
- 6. Scaling up: developing a model that overfits
- 7. Regularizing your model and tuning your hyperparameters

Pedro Domingos – Things to Know about Machine Learning



Useful Things to Know about Machine Learning

- 1. It's generalization that counts
- 2. Data alone is not enough
- 3. Overfitting has many faces
- 4. Intuition fails in high dimensions
- 5. Theoretical guarantees are not what they seem
- 6. More data beats a cleverer algorithm
- 7. Learn many models, not just one

Pedro Domingos, "A Few Useful Things to Know about Machine Learning," Commun. ACM, 2012

It's Generalization that Counts

• The goal of machine learning is to *generalize* beyond the examples in the training set

• Don't use test data for training

• Use cross validation to verify your model

Data Alone Is Not Enough

• No free lunch theorem (Wolpert)

-Every learner must embody some knowledge or assumptions beyond the data

Learners combine knowledge with data to grow programs

Overfitting Has Many Faces

- Ex: when your model accuracy is 100% on training data but only 50% on test data, when in fact it could have 75% on both, it has overfit.
- Overfitting has many forms. Example: bias & variance
- Combat overfitting
 - Cross validation
 - Add regularization term



Intuition Fails in High Dimensions (Number of Features)

- Curse of Dimensionality
- Algorithms that work fine in low dimensions fail when the input is high-dimensional
- Generalizing correctly becomes exponentially harder as the dimensionality of the examples grows
- Our intuition only comes from 3-dimension

Theoretical Guarantees Are Not What They Seem

- Theoretical bounds are usually very loose
- The main role of theoretical guarantees in machine learning is to help understand and drive force for algorithm design

More Data Beats a Cleverer Algorithm

• Try simplest algorithm first



Learn Many Models, Not Just One

- Ensembling methods: Random Forest ,XGBoost, Late Fusion
- Combining different models can get better results



References

- Francois Chollet, "Deep Learning with Python", Chapter 4
- Pedro Domingos, "A Few Useful Things to Know about Machine Learning," Commun. ACM, 2012
- <u>https://ml-cheatsheet.readthedocs.io/en/latest/index.html</u>
- <u>https://machinelearningmastery.com/roc-curves-and-precision-recall-curves-for-classification-in-python/</u>
- <u>https://towardsdatascience.com/data-types-in-statistics-347e152e8bee</u>
- <u>https://en.wikipedia.org/wiki/Naive_Bayes_classifier</u>