CS 6501: Deep Learning for Computer Graphics

Basics of Machine Learning

Connelly Barnes

Overview

- Supervised, unsupervised, and reinforcement learning
- Simple learning models
 - Clustering
 - Linear regression
 - Linear Support Vector Machines (SVM)
 - k-Nearest Neighbors
- Overfitting and generalization
- Training, testing, validation
- Balanced datasets
- Measuring performance of a classifier

Supervised, Unsupervised, Reinforcement

- 3 broad categories:
 - **Supervised learning**: computer presented with example inputs and desired outputs by a "teacher", goal is to learn general rule that maps inputs to outputs.
 - **Unsupervised learning**: No output labels are given to the algorithm, leaving it on its own to find structure in the inputs.
 - **Reinforcement learning**: An agent determines what actions to best take in an environment to maximize some notion of cumulative reward.

What Kind of Learning is This?

• Learn given input image, whether it is truck or car? Training data:



Images are Creative Commons, sources: [1], [2], [3], [4], [5], [6]

What Kind of Learning is This?

 \mathbf{X}_1

• We have a dataset of customers, each with 2 associated attributes (x₁ and x₂). We want to discover groups of similar customers.



What Kind of Learning is This?

Outtakes

[Peng et al., Terrain-Adaptive Locomotion..., SIGGRAPH 2016]

Overview

- Supervised, unsupervised, and reinforcement learning
- Simple learning models
 - Clustering
 - Linear regression
 - Linear Support Vector Machines (SVM)
 - k-Nearest Neighbors
- Overfitting and generalization
- Training, testing, validation
- Balanced datasets
- Measuring performance of a classifier

- Unsupervised learning
- Requires input data, but no labels
- Detects patterns, e.g.
 - Groups of similar emails, similar webpages in search results
 - Similar customer shopping patterns
 - Regions of images



- Idea: group together similar instances
- Example: 2D point patterns



- Idea: group together similar instances
- **Example**: 2D point patterns



- Idea: group together similar instances
- **Example**: 2D point patterns



- Idea: group together similar instances
- **Problem**: How to define "similar"?
- **Problem**: How many clusters?



- Similarity: in Euclidean space **R**^{*n*}, could be a distance function.
- For example: $D(x, y) = ||x y||_2^2$
- Clustering results will depend on measure of similarity / dissimilarity



Clustering Algorithms

- Partitioning algorithms (flat)
 - K-means

- Hierarchical algorithms
 - Bottom-up: agglomerative
 - Top-down: divisive





Clustering Examples: Image Segmentation

• Divide an image into regions that are perceptually similar to humans.



Clustering Examples: Biology

Α

• Cluster species based on e.g. genetic or phenotype similarity.



- Iterative clustering algorithm based on partitioning (flat).
- Initialize: Pick k random points as cluster centers.
- Iterate until convergence:
 - Assign each point based on the closest cluster center.
 - Update each cluster center based on the mean of the points assigned to it.

- Iterative clustering algorithm based on partitioning (flat).
- Initialize: Pick *k* random points as cluster centers.
- Iterate until convergence:
 - Assign each point based on the closest cluster center.
 - Update each cluster center based on the mean of the points assigned to it.

Cluster centers

- Iterative clustering algorithm based on partitioning (flat).
- Initialize: Pick *k* random points as cluster centers.
- Iterate until convergence:
 - Assign each point based on the closest cluster center.
 - Update each cluster center based on the mean of the points assigned to it.

What color should this point be?

- Iterative clustering algorithm based on partitioning (flat).
- Initialize: Pick *k* random points as cluster centers.
- Iterate until convergence:
 - Assign each point based on the closest cluster center.
 - Update each cluster center based on the mean of the points assigned to it.

What color should this point be?

- Iterative clustering algorithm based on partitioning (flat).
- Initialize: Pick *k* random points as cluster centers.
- Iterate until convergence:
 - Assign each point based on the closest cluster center.
 - Update each cluster center based on the mean of the points assigned to it.



- Iterative clustering algorithm based on partitioning (flat).
- Initialize: Pick *k* random points as cluster centers.
- Iterate until convergence:
 - Assign each point based on the closest cluster center.
 - Update each cluster center based on the mean of the points assigned to it.

What color should this point be?

- Iterative clustering algorithm based on partitioning (flat).
- Initialize: Pick *k* random points as cluster centers.
- Iterate until convergence:
 - Assign each point based on the closest cluster center.
 - Update each cluster center based on the mean of the points assigned to it.



- Iterative clustering algorithm based on partitioning (flat).
- Initialize: Pick *k* random points as cluster centers.
- Iterate until convergence:
 - Assign each point based on the closest cluster center.
 - Update each cluster center based on the mean of the points assigned to it.



- Iterative clustering algorithm based on partitioning (flat).
- Initialize: Pick *k* randon Any changes?
- Iterate until convergence:
 - Assign each point based on the closest cluster center.
 - Update each cluster center based on the mean of the points assigned to it.





- Iterative clustering algorithm based on partitioning (flat).
- Initialize: Pick *k* random points as cluster centers.
- Iterate until convergence:
 - Assign each point based on the closest cluster center.
 - Update each cluster center based on the mean of the points assigned to it.



- Iterative clustering algorithm based on partitioning (flat).
- Initialize: Pick *k* random points as cluster centers.
- Iterate until convergence:
 - Assign each point based on the closest cluster center.
 - Update each cluster center based on the mean of the points assigned to it.

Result of k-Means:

Minimizes within-cluster sum of squares distance:

$$rgmin_{\mathbf{S}} \sum_{i=1}^{\kappa} \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2$$
 (1)

- Here μ_i is the mean of the points belonging to cluster S_i.
- No guarantee algorithm will converge to global minimum.
- Can run several times and take best result according to (1).

Overview

- Supervised, unsupervised, and reinforcement learning
- Simple learning models
 - Clustering
 - Linear regression
 - Linear Support Vector Machines (SVM)
 - k-Nearest Neighbors
- Overfitting and generalization
- Training, testing, validation
- Balanced datasets

Linear Regression

- Uses a linear model to model relationship between dependent variable $y \in \mathbb{R}$, and input (independent) variables $x_1, ..., x_n \in \mathbb{R}^n$
- Is this supervised or unsupervised learning?



Linear Regression

• Uses a linear model to model relationship between dependent variable $y \in \mathbb{R}$, and input (independent) variables $x_1, ..., x_n \in \mathbb{R}^n$

For each observation (data point)
$$i = 1, ..., m$$
:
 $y_i = \mathbf{w} \cdot \mathbf{x}_i + b$
 $= w_1 x_{i,1} + \dots + w_n x_{i,n} + \dots + b$
Here $x_{i,j}$ is observation i of input variable j .
Parameters of model: \mathbf{w}, b .

Linear Regression

- Can simply the model by adding additional input that is always one: $x_{i,n+1} = 1$ i = 1, ..., m
- The corresponding parameter in w is called the **intercept**.

For each observation (data point)
$$i = 1, ..., m$$
:
 $y_i = \mathbf{w} \cdot \mathbf{x}_i$
 $= w_1 x_{i,1} + \dots + w_{n+1} x_{i,n+1}$
Parameters of model: **w**.

Linear Least Squares Regression

- Define an objective function or loss function to optimize the model
- One loss function: least squares ("least squared error").

$$E = \sum_{i=1}^{m} (\mathbf{w} \cdot \mathbf{x}_i - y_i)^2$$

Parameters of model: w.

- What is *E* for the 2D line fitting case at right? (blackboard)
- How to minimize *E*?



Linear Least Squares Regression

Set derivatives of objective function with respect to parameters equal to zero.

$$\frac{\partial E}{\partial w_j} = \frac{\partial}{\partial w_j} \sum_{i=1}^m (\mathbf{w} \cdot \mathbf{x}_i - y_i)^2 = 0$$

$$2\sum_{i=1}^{m} (\mathbf{w} \cdot \mathbf{x}_i - y_i) x_{ij} = 0$$

$$\sum_{i=1}^{m} \left(\sum_{k=1}^{n} x_{ik} w_k - y_i \right) x_{ij} = 0$$

Normal equations:



Linear Least Squares Regression

• Normal equations in matrix form:

 $(\mathbf{X}^{\mathrm{T}}\mathbf{X})\mathbf{w} = \mathbf{X}^{\mathrm{T}}\mathbf{y}$

$$\mathbf{w} = \left(\mathbf{X}^{\mathrm{T}}\mathbf{X}\right)^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}$$

X is the matrix with x_{ij} being observation *i* of input variable *j*.

y is the vector of dependent variable (output) observations.

Linear Least Squares Example

• Suppose we have three observations (m=3) of one input variable x_1 :


Linear Least Squares Example

- Suppose we have three observations (m=3) of one input variable x_1 .
- Add additional constant variable x₂:

Xa

•

$$\mathbf{X}_{1} \qquad \mathbf{X}_{2} \qquad \mathbf{y}$$

$$0 \qquad 1 \qquad 0$$

$$0.5 \qquad 1 \qquad 0.6$$

$$1 \qquad 1 \qquad 0.9$$

$$\mathbf{X} = \begin{bmatrix} 0 & 1\\ 0.5 & 1\\ 1 & 1 \end{bmatrix}, \mathbf{y} = \begin{bmatrix} 0\\ 0.6\\ 0.9 \end{bmatrix}, \text{ so } \mathbf{w} = (\mathbf{X}^{T}\mathbf{X})^{-1}\mathbf{X}^{T}\mathbf{y} = \begin{bmatrix} 0.9\\ 0.05 \end{bmatrix}$$

Linear Least Squares Example



Overview

- Supervised, unsupervised, and reinforcement learning
- Simple learning models
 - Clustering
 - Linear regression
 - Linear Support Vector Machines (SVM)
 - k-Nearest Neighbors
- Overfitting and generalization
- Training, testing, validation
- Balanced datasets
- Measuring performance of a classifier

- In linear regression, we had input variables $x_1, ..., x_n$ and we regressed them against a dependent variable $y \in \mathbb{R}$
- But what if we want to make a classifier?
- For example, a binary classifier could predict either y = -1, y = 1
- One simple option: use linear regression to find a linear model that best fits the data
- But this will not necessarily generalize well to new inputs.

• Idea: if data are separable by a linear hyperplane, then maximize separation distance (margin) between points.



• Two hyperplanes:

$$ec w \cdot ec x + b = 1$$

and

 $ec{w} \cdot ec{x} + b = -1.$

- Distance between hyperplanes is:
 - $rac{2}{\|ec{w}\|}$
- So to maximize distance, minimize $\|ec{w}\|$



From Wikipedia

• For each point *i*, either:

$$ec{w}\cdotec{x}_i+b\geq 1,$$
 if $y_i=1$

or

$$ec{w}\cdotec{x}_i+b\leq -1,$$
 if $y_i=-1.$

• This can be rewritten as, for each *i*:

 $y_i(ec{w}\cdotec{x}_i+b)\geq 1$



- So our minimization problem becomes:
- Minimize $\|ec{w}\|$ subject to the constraint:

$$y_i(ec{w}\cdotec{x}_i+b)\geq 1$$
 $i=1,\ldots,m$

- Can be solved with quadratic programming
- Maximizes distance (margin) between two classes of data



Support Vectors

• If data are not linearly separable, can use a **soft margin classifier**, which has an objective function that sums for all data points *i*, a penalty of zero if the data point is correctly classified, otherwise, the distance to the margin.



How to Use Linear SVMs in Deep Learning?

- Linear SVMs tend to perform well with small amounts of training data
- Deep learning tends to perform well with large amounts of data
- What to do if we have a new problem with only a small dataset?
- One solution: use linear SVM
- Another solution: transfer learning.
 - Use a deep learning model trained on different problem
 - Train a linear SVM using features extracted from neural network.

Overview

- Supervised, unsupervised, and reinforcement learning
- Simple learning models
 - Clustering
 - Linear regression
 - Linear Support Vector Machines (SVM)

k-Nearest Neighbors

- Overfitting and generalization
- Training, testing, validation
- Balanced datasets
- Measuring performance of a classifier

k-Nearest Neighbors

- Suppose we can measure distance between input features.
- For example, Euclidean distance: $D(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} \mathbf{y}\|_2^2$
- *k*-Nearest Neighbors simply uses the distance to the nearest *k* points to determine the classification or regression.
 - Classifier: take most common class within the *k* nearest points
 - Regression: take mean of *k* nearest points
- No parameters, so no need to "train" the algorithm

k-Nearest Neighbors Example, k=3



k-Nearest Neighbors Example, k=5



Overview

- Supervised, unsupervised, and reinforcement learning
- Simple learning models
 - Clustering
 - Linear regression
 - Linear Support Vector Machines (SVM)
 - k-Nearest Neighbors
- Overfitting and generalization
- Training, testing, validation
- Balanced datasets
- Measuring performance of a classifier

• Will this model have decent prediction for new inputs? (i.e. inputs similar to the training exemplars in blue)





• How about the model here, shown as the blue curve?



• **Overfitting**: the model describes random noise or errors instead of the underlying relationship.



- **Overfitting**: the model describes random noise or errors instead of the underlying relationship.
- Frequently occurs when model is overly complex (e.g. has too many parameters) relative to the number of observations.
- Has poor predictive performance.



- **Overfitting**: the model describes random noise or errors instead of the underlying relationship.
- Frequently occurs when model is overly complex (e.g. has too many parameters) relative to the number of observations.
- Has poor predictive performance.



- A rule of thumb for linear regression: <u>one in ten rule</u>
- One predictive variable can be studied for every ten events.
- In general, want number of data points >> number of parameters.
- But models with more parameters often perform better!
- One solution: gradually increase number of parameters in model until it starts to overfit, and then stop.

Overfitting Example with 2D Classifier

From <u>ConvnetJS Demo: 2D Classification with Neural Networks</u>



23 data points, 17 parameters (5 neurons)

Overfitting Example with 2D Classifier

From <u>ConvnetJS Demo: 2D Classification with Neural Networks</u>



23 data points, 32 parameters (10 neurons)

Overfitting Example with 2D Classifier

• From <u>ConvnetJS Demo: 2D Classification with Neural Networks</u>



23 data points, 102 parameters (22 neurons total)

Generalization

- Generalization error is a measure of how accurately an algorithm is able to predict outcome values for previously unseen data.
- A model that is **overfit** will have poor generalization.





Overview

- Supervised, unsupervised, and reinforcement learning
- Simple learning models
 - Clustering
 - Linear regression
 - Linear Support Vector Machines (SVM)
 - k-Nearest Neighbors
- Overfitting and generalization
- Training, testing, validation
- Balanced datasets
- Measuring performance of a classifier

Training, testing, validation

- Break dataset into three parts by random sampling:
 - Training dataset: model is fit directly to this data
 - **Testing dataset**: model sees this data only once; used to measure the final performance of the classifier.
 - Validation dataset: model is repeatedly "tested" on this data during the training process to gain insight into overfitting
- Common percentages:
 - Training (80%), testing (15%), validation (5%).

Training, testing, validation

• For neural networks, typically keep running training until validation error increases, then stop.



Cross-validation

- Repeatedly partition data into subsets: training and test.
- Take mean performance of classifier over all such partitions.
- Leave one out: train on n-1 samples, test on 1 sample.
 - Requires training n times.
- **k-fold cross-validation**: randomly partition data into k subsets (folds), at each iteration, train on k-1 folds, test on the other fold.
 - Requires training k times.
 - Common: 10-fold cross-validation
- Less common for deep learning (why?)

Overview

- Supervised, unsupervised, and reinforcement learning
- Simple learning models
 - Clustering
 - Linear regression
 - Linear Support Vector Machines (SVM)
 - k-Nearest Neighbors
- Overfitting and generalization
- Training, testing, validation
- Balanced datasets
- Measuring performance of a classifier

Balanced datasets

- Unbalanced dataset:
 - Suppose we have a binary classification problem (labels: 0, 1)
 - Suppose 99% of our observations are class 0.
 - We might learn the model "everything is zero."
 - This model would be 99% accurate, but not model class 1 at all.
- Balanced dataset:
 - Equal numbers of observations of each class

Overview

- Supervised, unsupervised, and reinforcement learning
- Simple learning models
 - Clustering
 - Linear regression
 - Linear Support Vector Machines (SVM)
 - k-Nearest Neighbors
- Overfitting and generalization
- Training, testing, validation
- Balanced datasets
- Measuring performance of a classifier

Confusion Matrix

• If n classes, n x n matrix comparing actual versus predicted classes.



Example: Handwritten Digit Recognition

Class	1	2	3	4	5	6	7	8	9	0	Error Rate	
1	191	0	0	5	1	0	1	0	2	0	4.5	
2	0	188	2	0	0	1	3	0	0	6	6.0	
3	0	3	191	0	1	0	2	0	3	0	4.5	
4	8	0	0	187	4	0	1	0	0	0	6.5	
5	0	0	0	0	193	0	0	0	7	0	3.5	
6	0	0	0	0	1	196	0	2	0	1	2.0	
7	2	2	0	2	0	1	190	0	1	2	5.0	
8	0	1	0	0	1	2	2	196	0	0	2.0	
9	5	0	2	0	8	0	3	0	179	3	10.5	
0	1	4	0	0	0	1	1	0	1	192	4.5	L
Overall error rate 4.85%												

Slide from Nelson Morgan at ICSI / Berkeley

Example: Handwritten Digit Recognition

Visualization from MathWorks

Confusion Matrix

- For binary classifier, can call one class positive, the other negative.
- Should we call cats positive or negative?

Class

Actual



Class Predicted



```
Photo from [1]
```
- For binary classifier, can call one class positive, the other negative.
- Cats are cuter, so cats = positive.

		Class Predicted by Model	
		Positive	Negative
Class	Positive	9	1
Actual	Negative	4	6

• For binary classifier, can call one class positive, the other negative.

		Class Predicted by Model	
		Positive	Negative
ISS		True	
Cla	Positive	positive	1
ctual			True
Ac	Negative	4	negative

• For binary classifier, can call one class positive, the other negative.

		Class Predicted by Model	
		Positive	Negative
ISS		True	
Cla	Positive	positive	?
tual			True
Ac	Negative	?	negative

• For binary classifier, can call one class positive, the other negative.

		Class Predicted by Model	
		Positive	Negative
ISS		True	False
Cla	Positive	positive	negative
tual		False	True
Ac	Negative	positive	negative

Classifier Performance

- Accuracy: (TP + TN) / (Population Size)
- **Precision**: TP / (Predicted Positives) = (TP + FP)
- **Recall**: TP / (Actual Positives) = TP / (TP + FN) (also known as **sensitivity, true positive rate**)
- Specificity: TN / (Actual Negatives) = TN / (TN + FP) (also known as true negative rate) Class F
 by F

Class Predicted by Model

	Positive	Negative
	True	False
Positive	positive	negative
	False	True
Vegative	positive	negative

Actual Class

- Accuracy: (TP + TN) / (Population Size)
- **Precision**: TP / (Predicted Positives) = (TP + FP)
- **Recall**: TP / (Actual Positives) = TP / (TP + FN) (also known as **sensitivity, true positive rate**)
- Specificity: TN / (Actual Negatives) = TN / (TN + FP) (also known as true negative rate)
 Class Predicted by Model

by ModelPositiveNegativePositiveTP: 99FN: 0NegativeFP: 1TN: 0

ass

ctual

- Accuracy: (TP + TN) / (Population Size) = ?
- **Precision**: TP / (Predicted Positives) = (TP + FP)
- **Recall**: TP / (Actual Positives) = TP / (TP + FN) (also known as **sensitivity, true positive rate**)
- Specificity: TN / (Actual Negatives) = TN / (TN + FP) (also known as true negative rate) Class Pr by M

Class Predicted
by ModelPositiveNegativePositiveTP: 99FN: 0

TN: 0

FP: 1

ass

ctual

Negative

- Accuracy: (TP + TN) / (Population Size) = 99%
- **Precision**: TP / (Predicted Positives) = (TP + FP) = ?
- **Recall**: TP / (Actual Positives) = TP / (TP + FN) (also known as **sensitivity**, **true positive rate**)
- Specificity: TN / (Actual Negatives) = TN / (TN + FP) (also known as true negative rate)
 Class Predicted by Model



ass

ctual

- Accuracy: (TP + TN) / (Population Size) = 99%
- **Precision**: TP / (Predicted Positives) = (TP + FP) = 99%
- Recall: TP / (Actual Positives) = TP / (TP + FN) (also known as sensitivity, true positive rate)
- Specificity: TN / (Actual Negatives) = TN / (TN + FP) (also known as true negative rate)
 Class Predicted by Model

- Accuracy: (TP + TN) / (Population Size) = 99%
- **Precision**: TP / (Predicted Positives) = (TP + FP) = 99%
- Recall: TP / (Actual Positives) = TP / (TP + FN) = ?
 (also known as sensitivity, true positive rate)
- Specificity: TN / (Actual Negatives) = TN / (TN + FP) (also known as true negative rate)
 Class Predicted by Model

- Accuracy: (TP + TN) / (Population Size) = 99%
- **Precision**: TP / (Predicted Positives) = (TP + FP) = 99%
- **Recall**: TP / (Actual Positives) = TP / (TP + FN) = 100% (also known as **sensitivity**, **true positive rate**)
- Specificity: TN / (Actual Negatives) = TN / (TN + FP) (also known as true negative rate)
 Class Predicted by Model



- Accuracy: (TP + TN) / (Population Size) = 99%
- **Precision**: TP / (Predicted Positives) = (TP + FP) = 99%
- **Recall**: TP / (Actual Positives) = TP / (TP + FN) = 100% (also known as **sensitivity**, **true positive rate**)
- Specificity: TN / (Actual Negatives) = 0% (also known as true negative rate)



ass

ctual

Summary

- Supervised, unsupervised, and reinforcement learning
- Simple learning models
 - Clustering
 - Linear regression
 - Linear Support Vector Machines (SVM)
 - k-Nearest Neighbors
- Overfitting and generalization
- Training, testing, validation
- Balanced datasets
- Measuring performance of a classifier