

Chapter 5. Machine Learning Basics

Shansuo Liang

August 25th, 2018

Outline

- 1 Definition and Examples
- 2 Capacity, Overfitting and Underfitting
- 3 Hyperparameters and Validation Sets
- 4 Supervised & Unsupervised Learning
- 5 Stochastic Gradient Descent
- 6 From Machine Learning to Deep Learning

Definition by T. Mitchell

- T. Mitchell: A computer program is said to learn from experience E with respect to some task T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .

Definition by T. Mitchell

- T. Mitchell: A computer program is said to learn from experience E with respect to some task T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .
- The Task, T : Classification, Regression, Transcription, Structured output, Anomaly detection, Denoising, Density estimation, etc.
- The Performance Measure, P : Error rate (the expected 0-1 loss), MSE, etc.
- The Experience, E : Supervised learning, Unsupervised learning, Reinforcement learning (fixed dataset or not), etc.

An Example: Linear Regression

- *T*: To build a linear system that can take a vector $\mathbf{x} \in \mathbb{R}^n$ as input and predict the value of a scalar $y \in \mathbb{R}$ as its output as

$$\hat{y} = \mathbf{w}^T \mathbf{x}, \quad (1)$$

where $\mathbf{w} \in \mathbb{R}^n$ is a vector of parameters.

- *P*: To compute MSE of the model on the test set:

$$\text{MSE}_{\text{test}} = \frac{1}{m} \|\hat{\mathbf{y}}^{\text{test}} - \mathbf{y}^{\text{test}}\|_2^2. \quad (2)$$

- *E*: To design an algorithm that will improve the weights \mathbf{w} in a way that reduce MSE_{test} when the algorithm can observe a training set $\{\mathbf{X}^{\text{test}}, \mathbf{y}^{\text{test}}\}$. (Solved by the normal equations)

Outline

- 1 Definition and Examples
- 2 Capacity, Overfitting and Underfitting
- 3 Hyperparameters and Validation Sets
- 4 Supervised & Unsupervised Learning
- 5 Stochastic Gradient Descent
- 6 From Machine Learning to Deep Learning

Capacity, Overfitting and Underfitting

- Training error & Test error (expected):
 - ① IID assumptions: Training error = Test error,
 - ② The parameters are not fixed: Test error \geq Training error.

Capacity, Overfitting and Underfitting

- Training error & Test error (expected):
 - ① IID assumptions: Training error = Test error,
 - ② The parameters are not fixed: Test error \geq Training error.
- Two factors correspond to underfitting and overfitting:
 - ① Make the training error small,
 - ② Make the gap between training and test error small.

Capacity, Overfitting and Underfitting

- Training error & Test error (expected):
 - ① IID assumptions: Training error = Test error,
 - ② The parameters are not fixed: Test error \geq Training error.
- Two factors correspond to underfitting and overfitting:
 - ① Make the training error small,
 - ② Make the gap between training and test error small.
- Capacity is defined by the ability to fit a wide variety of functions. We can control whether a model is more likely to be overfit or underfit by altering its capacity.

Capacity, Overfitting and Underfitting

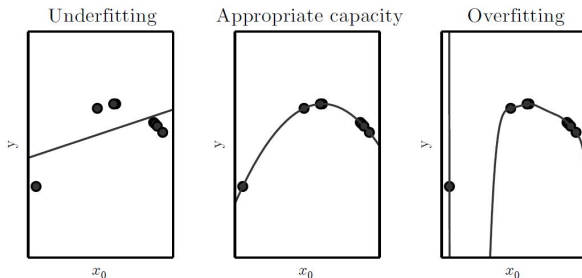


Figure: (Left) A linear function fit to the data suffers from underfitting. (Center) A quadratic function fit to the data generalizes well to unseen points. (Right) A polynomial of degree 9 fit to the data suffers from overfitting.

Capacity, Overfitting and Underfitting

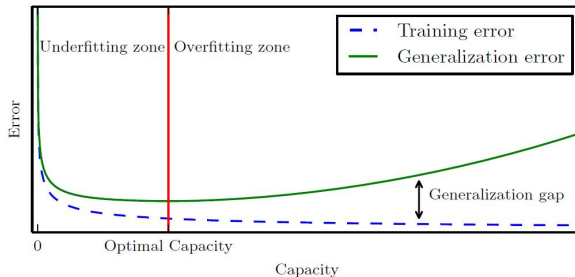


Figure: Typical relationship between capacity and error.

Outline

- 1 Definition and Examples
- 2 Capacity, Overfitting and Underfitting
- 3 Hyperparameters and Validation Sets
- 4 Supervised & Unsupervised Learning
- 5 Stochastic Gradient Descent
- 6 From Machine Learning to Deep Learning

Hyperparameters

- Hyperparameters: Settings that we can use to control the algorithm's behavior. The values of hyperparameters are not adapted by the learning algorithm itself.

Hyperparameters

- Hyperparameters: Settings that we can use to control the algorithm's behavior. The values of hyperparameters are not adapted by the learning algorithm itself.

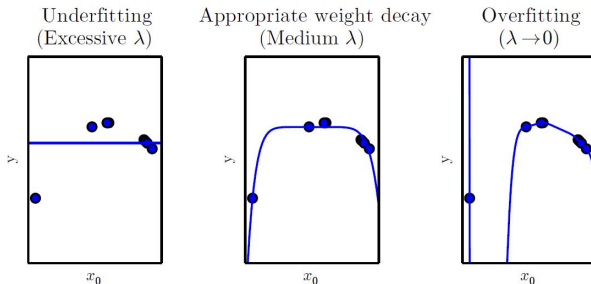


Figure: Example: Fit a high-degree polynomial regression model. The true function is quadratic, but here we use only models with degree 9. We vary the amount of weight decay to prevent these high-degree models from overfitting.

Validation Sets

- A setting is chosen to be a hyperparameter since it is not appropriate to learn that on the training set. If learned on the training set, such hyperparameters would always choose the maximum possible model capacity, resulting in overfitting.

Validation Sets

- A setting is chosen to be a hyperparameter since it is not appropriate to learn that on the training set. If learned on the training set, such hyperparameters would always choose the maximum possible model capacity, resulting in overfitting.
- Validation set: Split the training data into two disjoint subsets, and one of these subsets is used to learn the parameters.

Validation Sets

- A setting is chosen to be a hyperparameter since it is not appropriate to learn that on the training set. If learned on the training set, such hyperparameters would always choose the maximum possible model capacity, resulting in overfitting.
- Validation set: Split the training data into two disjoint subsets, and one of these subsets is used to learn the parameters.
- Cross-validation: When the dataset is small, repeating the training and testing computation on different randomly chosen subsets or splits of the original sets.

Outline

- 1 Definition and Examples
- 2 Capacity, Overfitting and Underfitting
- 3 Hyperparameters and Validation Sets
- 4 Supervised & Unsupervised Learning
- 5 Stochastic Gradient Descent
- 6 From Machine Learning to Deep Learning

Supervised Learning

- Definition: Roughly speaking, learning algorithms that learn to associate some inputs \mathbf{x} and some outputs \mathbf{y} .
In many cases, the outputs \mathbf{y} may be difficult to collect and must be provided by a human “supervisor”.

Supervised Learning

- Definition: Roughly speaking, learning algorithms that learn to associate some inputs \mathbf{x} and some outputs \mathbf{y} .
In many cases, the outputs \mathbf{y} may be difficult to collect and must be provided by a human “supervisor”.
- Probabilistic interpretation: Most supervised learning algorithms are based on estimating a probability distribution $p(\mathbf{y}|\mathbf{x})$. Using ML estimation to find the best parameter vector θ for a parameter family of distribution $p(\mathbf{y}|\mathbf{x}; \theta)$.

Supervised Learning

- Definition: Roughly speaking, learning algorithms that learn to associate some inputs \mathbf{x} and some outputs \mathbf{y} .
In many cases, the outputs \mathbf{y} may be difficult to collect and must be provided by a human “supervisor”.
- Probabilistic interpretation: Most supervised learning algorithms are based on estimating a probability distribution $p(\mathbf{y}|\mathbf{x})$. Using ML estimation to find the best parameter vector θ for a parameter family of distribution $p(\mathbf{y}|\mathbf{x}; \theta)$.
- Examples: Logistic Regression, Support Vector Machines, $k - NN$, Decision Tree.

Logistic Regression

- Logistic regression or classification: For those with the output $y \in \{0, 1\}$ or other discrete values.

Logistic Regression

- Logistic regression or classification: For those with the output $y \in \{0, 1\}$ or other discrete values.

- Probabilistic interpretation:

Classes $y \in \{0, 1\}$;

Function $z = \mathbf{w}^T \mathbf{x}$;

Hypothesis logistic sigmoid function $\sigma(z) = \frac{1}{1+e^{-z}}$.

- Decision rules:

$$\hat{y} = \begin{cases} 1, & z > 0 \\ 0, & z \leq 0 \end{cases}.$$

Logistic Regression

- Define $\pi_i = \sigma(\mathbf{w}^T \mathbf{x}_i)$, then $p(y_i | \mathbf{x}_i, \mathbf{w}) = \pi_i^{y_i} \cdot (1 - \pi_i)^{1-y_i}$.
- Loglikelihood of all data: $\ell(\mathbf{w}) = \sum_i \log p(y_i | \mathbf{x}_i, \mathbf{w})$.
- MLE: $\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} \ell(\mathbf{w})$.
 - ① No close-form solution;
 - ② Gradient ascent / SGD (first-order);
 - ③ Newton-Raphson method (second-order).

Support Vector Machines

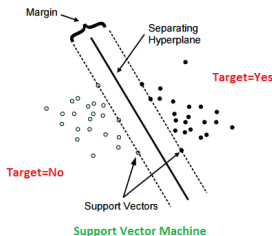
- Non-probabilistic interpretation:

Classes $y \in \{-1, 1\};$

Function $z = \mathbf{w}^T \mathbf{x} + b;$

Decision rules $\hat{y} = \text{sign}(z).$

[http://www.saedsayad.com/
support_vector_machine.htm](http://www.saedsayad.com/support_vector_machine.htm)

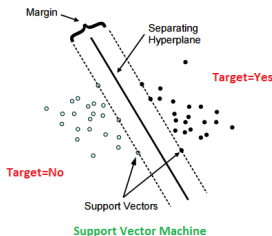


Support Vector Machines

- Non-probabilistic interpretation:

Classes $y \in \{-1, 1\}$;
Function $z = \mathbf{w}^T \mathbf{x} + b$;
Decision rules $\hat{y} = \text{sign}(z)$.

[http://www.saedsayad.com/
support_vector_machine.htm](http://www.saedsayad.com/support_vector_machine.htm)



- Assumption: Data set $\{(\mathbf{x}_i, y_i)\}$ are linearly separable.

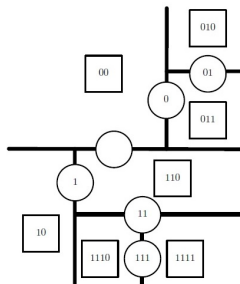
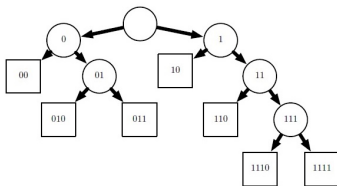
Margin: Distance from boundary to the closet point in training sets as

$$\gamma = \min_i \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|} = \min_i \frac{y_i(\mathbf{w}^T \mathbf{x}_i + b)}{\|\mathbf{w}\|}. \quad (3)$$

Idea: To maximize the margin. Solved by Quadratic Programming.

k -NN and Decision Tree

- k -nearest neighbors (k -NN) algorithm: When we want to produce an output y for a new test input \mathbf{x} , we find the k -nearest neighbors to \mathbf{x} in the training data \mathbf{X} . We then return the average of the corresponding y values in the training set.
- Decision Tree:



Unsupervised Learning

- Definition: Informally, unsupervised learning refers to most attempts to extract information from a distribution that do not require human labor to annotate example.

Unsupervised Learning

- Definition: Informally, unsupervised learning refers to most attempts to extract information from a distribution that do not require human labor to annotate example.
- A classic unsupervised learning task is to find the “best” representation of the data.

Unsupervised Learning

- Definition: Informally, unsupervised learning refers to most attempts to extract information from a distribution that do not require human labor to annotate example.
- A classic unsupervised learning task is to find the “best” representation of the data.
- Three common representations
 - ① Low-dimensional representations: Compress as much information as possible;
 - ② Sparse representations: Entries are mostly zeros;
 - ③ Independent representations: Disentangle the source of variation underlying the data distribution such that the dimensions of the representation are statistically independent.

Principle component analysis

- PCA learns an orthogonal, linear transformation of the data that projects an input \mathbf{x} to a representation \mathbf{z} .

Principle component analysis

- PCA learns an orthogonal, linear transformation of the data that projects an input \mathbf{x} to a representation \mathbf{z} .

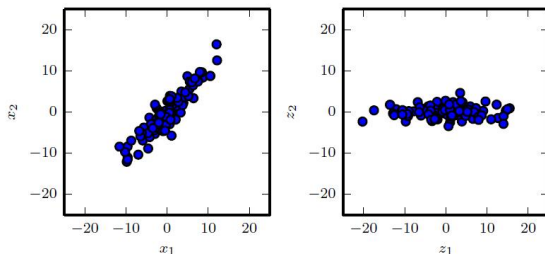


Figure: (Left) The original data consist of samples of \mathbf{x} . In this space, the variance might occur along directions that are not axis aligned. (Right) The transformed data $\mathbf{z} = \mathbf{x}^T \mathbf{W}$ now varies most along the axis z_1 . The direction of second-most variance is now along z_2 .

Principle Component Analysis

- Consider a zero-mean data matrix \mathbf{X} , the covariance matrix associated with \mathbf{X} is given by

$$\text{Var}[\mathbf{x}] = \frac{1}{m-1} \mathbf{X}^T \mathbf{X}. \quad (4)$$

PCA finds a representation $\mathbf{z} = \mathbf{W}^T \mathbf{x}$, where $\text{Var}[\mathbf{x}]$ is diagonal.

- The principle components of a design matrix \mathbf{X} are given by the eigenvectors of $\mathbf{X}^T \mathbf{X}$.

k -means Clustering

- The k -means clustering algorithm divides the training set into k different clusters of examples that are near each other.

k -means Clustering

- The k -means clustering algorithm divides the training set into k different clusters of examples that are near each other.
- Initializing k different centroids $\{\mu^{(1)}, \dots, \mu^{(k)}\}$, the k -means algorithm alternates between the following two different steps until convergence:
 - 1 Each training example is assigned to cluster i , where i is the index of the nearest centroid $\mu^{(i)}$.
 - 2 Each centroid $\mu^{(i)}$ is updated to the mean of all training examples $\mathbf{x}^{(i)}$ assigned to cluster i .

Outline

- 1 Definition and Examples
- 2 Capacity, Overfitting and Underfitting
- 3 Hyperparameters and Validation Sets
- 4 Supervised & Unsupervised Learning
- 5 Stochastic Gradient Descent
- 6 From Machine Learning to Deep Learning

Stochastic Gradient Descent

- The cost function used by a machine learning algorithm often decomposes as a sum over training examples of some per-example loss function, such as

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{data}} L(\mathbf{x}, y, \boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta}), \quad (5)$$

where L is the per-example loss $L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta}) = -\log p(y^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta})$.

Stochastic Gradient Descent

- The cost function used by a machine learning algorithm often decomposes as a sum over training examples of some per-example loss function, such as

$$J(\theta) = \mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{data}} L(\mathbf{x}, y, \theta) = \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}^{(i)}, y^{(i)}, \theta), \quad (5)$$

where L is the per-example loss $L(\mathbf{x}^{(i)}, y^{(i)}, \theta) = -\log p(y^{(i)} | \mathbf{x}^{(i)}, \theta)$.

- Gradient descent requires computing

$$\nabla_{\theta} J(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(\mathbf{x}^{(i)}, y^{(i)}, \theta). \quad (6)$$

As the size m grows to billions, a single gradient step becomes computational prohibitive.

Stochastic Gradient Descent

- The insight of SGD is that the gradient is an expectation that may be approximately estimated using a small set of samples.

Stochastic Gradient Descent

- The insight of SGD is that the gradient is an expectation that may be approximately estimated using a small set of samples.
- On each step of the algorithm, we can sample a minibatch of examples $\mathbb{B} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m')}\}$ drawn uniformly from the training set.
- The estimate of the gradient is formed as

$$\mathbf{g} = \frac{1}{m'} \sum_{i=1}^{m'} \nabla_{\theta} L(\mathbf{x}^{(i)}, y^{(i)}, \theta) \quad (7)$$

The SGD algorithm then follows $\theta \leftarrow \theta - \epsilon \mathbf{g}$ with ϵ being the learning rate.

Outline

- 1 Definition and Examples
- 2 Capacity, Overfitting and Underfitting
- 3 Hyperparameters and Validation Sets
- 4 Supervised & Unsupervised Learning
- 5 Stochastic Gradient Descent
- 6 From Machine Learning to Deep Learning

Building a Machine Learning Algorithm

- Recipe: Combine a specification of a dataset, a cost function, an optimization procedure and a model.

Building a Machine Learning Algorithm

- Recipe: Combine a specification of a dataset, a cost function, an optimization procedure and a model.
- For example, the linear regression combines a dataset consisting of \mathbf{X} and \mathbf{y} , the cost function

$$J(\mathbf{w}, b) = -\mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(y|\mathbf{x}), \quad (8)$$

the model specification $p_{\text{model}}(y|\mathbf{x}) = \mathcal{N}(y; \mathbf{x}^T \mathbf{w} + b, 1)$, and the optimization algorithm. is solving the normal equations for the gradient of the cost being zero.

Challenges Motivating Deep Learning

- The Curse of Dimensionality

Challenges Motivating Deep Learning

- The Curse of Dimensionality
- Local Constancy and Smoothness Regularization

Challenges Motivating Deep Learning

- The Curse of Dimensionality
- Local Constancy and Smoothness Regularization
- Manifold Learning

Thanks!

Q&A

Aug 25, 2018