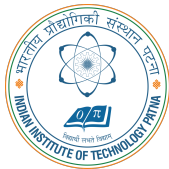


CS365: Deep Learning

Neural Networks-II



Arijit Mondal

Dept. of Computer Science & Engineering

Indian Institute of Technology Patna

`arijit@iitp.ac.in`

Machine Learning

- A form of applied statistics with
 - Increased emphasis on the use of computers to statistically estimate complicated function
 - Decreased emphasis on proving confidence intervals around these functions
- Two primary approaches
 - Frequentist estimators
 - Bayesian inference
- A ML/DL algorithm is an algorithm that is able to learn from data
- Mitchell (1997)
 - A computer program is said to learn from experience E with respect to some class of task T and performance measure P , if its performance at task in T as measured by P , improves with experience E .

Task

- A ML/DL task is usually described in terms of how the system should process an example
 - Example is a collection of features that have been quantitatively measured from some objects or events that we want the learning system process
 - Represented as $\mathbf{x} \in \mathbb{R}^n$ where x_i is a feature
 - Feature of an image — pixel values

Typical tasks

- Classification

- Need to predict which of the k categories some input belongs to
- Need to have a function $f: \mathbb{R}^n \rightarrow \{1, 2, \dots, k\}$
- $y = f(x)$ input x is assigned a category identified by y
- Examples
 - Object identification
 - Face recognition

- Regression

- Need to predict numeric value for some given input
- Need to have a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$
- Examples
 - Energy consumption
 - Amount of insurance claim

Typical tasks (contd.)

- Classification with missing inputs
 - Need to have a set of functions
 - Each function corresponds to classifying x with different subset of inputs missing
 - Examples
 - Medical diagnosis (expensive or invasive)

Typical tasks (contd.)

- Classification with missing inputs
 - Need to have a set of functions
 - Each function corresponds to classifying x with different subset of inputs missing
 - Examples
 - Medical diagnosis (expensive or invasive)
- Transcription
 - Need to convert relatively unstructured data into discrete, textual form
 - Optical character recognition
 - Speech recognition

Typical tasks (contd.)

- Classification with missing inputs
 - Need to have a set of functions
 - Each function corresponds to classifying x with different subset of inputs missing
 - Examples
 - Medical diagnosis (expensive or invasive)
- Transcription
 - Need to convert relatively unstructured data into discrete, textual form
 - Optical character recognition
 - Speech recognition
- Machine translation
 - Conversion of sequence of symbols in one language to some other language
 - Natural language processing (English to Spanish conversion)

Typical tasks (contd.)

- Structured output
 - Output is a vector with important relationship between the different elements
 - Mapping natural language sentence into a tree that describes grammatical structure
 - Pixel based image segmentation (eg. identify roads)

Typical tasks (contd.)

- Structured output
 - Output is a vector with important relationship between the different elements
 - Mapping natural language sentence into a tree that describes grammatical structure
 - Pixel based image segmentation (eg. identify roads)
- Anomaly detection
 - Observes a set of events or objects and flags if some of them are unusual
 - Fraud detection in credit card

Typical tasks (contd.)

- Structured output
 - Output is a vector with important relationship between the different elements
 - Mapping natural language sentence into a tree that describes grammatical structure
 - Pixel based image segmentation (eg. identify roads)
- Anomaly detection
 - Observes a set of events or objects and flags if some of them are unusual
 - Fraud detection in credit card
- Synthesis and sampling
 - Generate new example similar to past examples
 - Useful for media application
 - Text to speech

Performance measure

- Accuracy is one of the key measures
 - The proportion of examples for which the model produces correct outputs
 - Similar to error rate
 - Error rate often referred as expected 0-1 loss
- Mostly interested how DL algorithm performs on unseen data
- Choice of performance measure may not be straight forward
 - Transcription
 - Accuracy of the system at transcribing entire sequence
 - Any partial credit for some elements of the sequence are correct

Experience

- Kind of experience allowed during learning process
 - Supervised
 - Unsupervised

Supervised learning

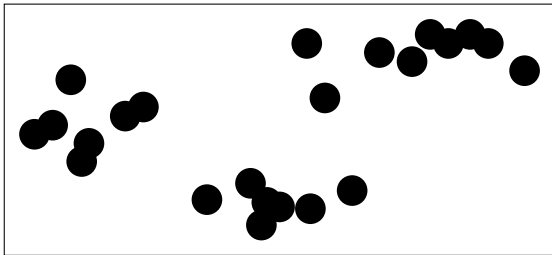
- Allowed to use labeled dataset
- Example — Iris
 - Collection of measurements of different parts of Iris plant
 - Each plant means each example
 - Features
 - Sepal length/width, petal length/width
 - Also record which species the plant belong to

Supervised learning (contd.)

- A set of labeled examples $\langle x_1, x_2, \dots, x_n, y \rangle$
 - x_i are input variables
 - y output variable
- Need to find a function $f: X_1 \times X_2 \times \dots \times X_n \rightarrow Y$
- Goal is to minimize error/loss function
 - Like to minimize over all dataset
 - We have limited dataset

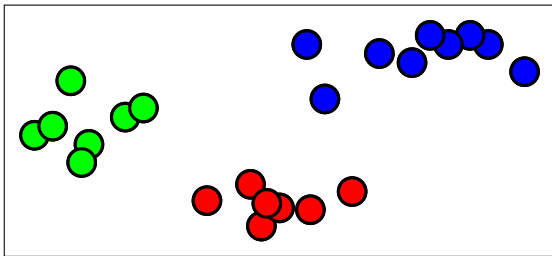
Unsupervised learning

- Learns useful properties of the structure of data set
- Unlabeled data
 - Tries to learn entire probability distribution that generated the dataset
 - Examples
 - Clustering, dimensionality reduction



Unsupervised learning

- Learns useful properties of the structure of data set
- Unlabeled data
 - Tries to learn entire probability distribution that generated the dataset
 - Examples
 - Clustering, dimensionality reduction



Supervised vs Unsupervised learning

- Unsupervised attempts to learn implicitly or explicitly probability distribution of $p(\mathbf{x})$
- Supervised tries to predict y from \mathbf{x} ie. $p(y|\mathbf{x})$

Supervised vs Unsupervised learning

- Unsupervised attempts to learn implicitly or explicitly probability distribution of $p(\mathbf{x})$
- Supervised tries to predict y from \mathbf{x} ie. $p(y|\mathbf{x})$
- Unsupervised learning can be decomposed as supervised learning

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, x_2, \dots, x_{i-1})$$

Supervised vs Unsupervised learning

- Unsupervised attempts to learn implicitly or explicitly probability distribution of $p(\mathbf{x})$
- Supervised tries to predict y from \mathbf{x} ie. $p(y|\mathbf{x})$
- Unsupervised learning can be decomposed as supervised learning

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, x_2, \dots, x_{i-1})$$

- Solving supervised learning using traditional unsupervised learning

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}, y)}{\sum_{y'} p(\mathbf{x}, y')}$$

Multi layer neural network

- Pre-activation in layer

$$k > 0 \quad (h^{(0)}(x) = x)$$

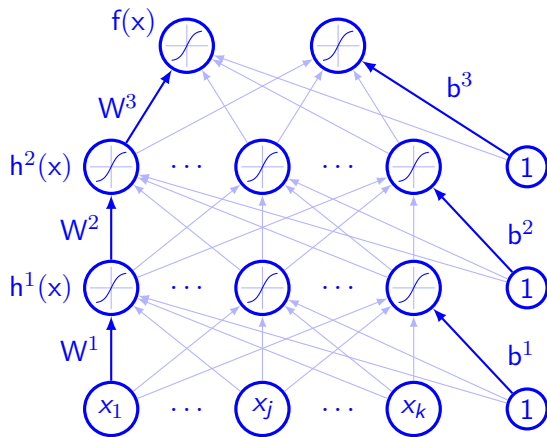
$$a^{(k)}(x) = b^{(k)} + W^{(k)}h^{(k-1)}_x$$

- Hidden layer activation

$$h^{(k)}(x) = g(a^{(k)}(x))$$

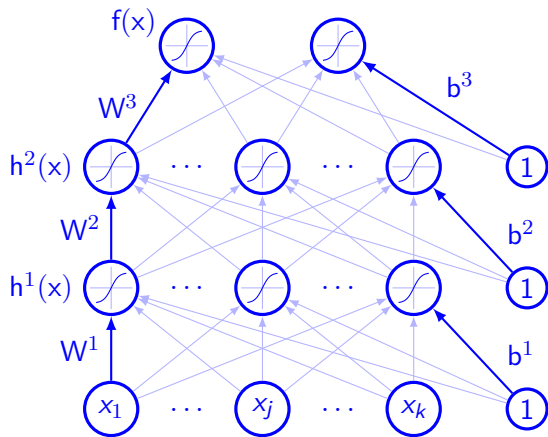
- Output layer activation

$$h^{(L+1)}(x) = o(a^{(L+1)}(x)) = f(x)$$



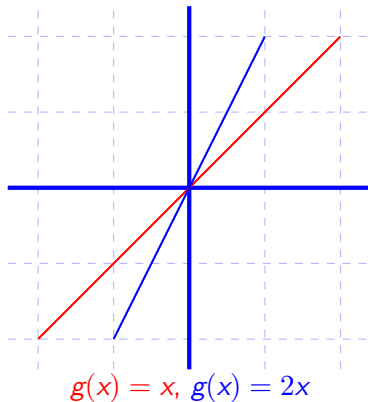
Multi layer neural network

- Design issues
 - Number of layers
 - Number of neurons in each layer
 - Activation function
 - Output function
 - Loss function
 - Optimizer



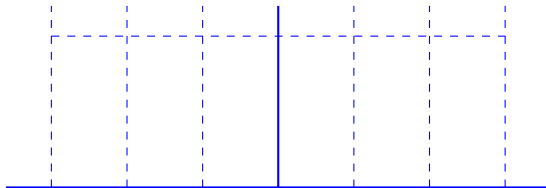
Activation function

- Linear activation function
 - Not very interesting
 - No change in values
 - Huge range



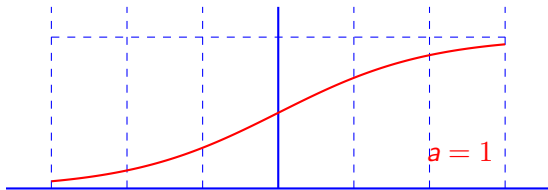
Activation function

- Sigmoid function
 - Values lie between 0 and 1
 - Strictly increasing function
 - Bounded



Activation function

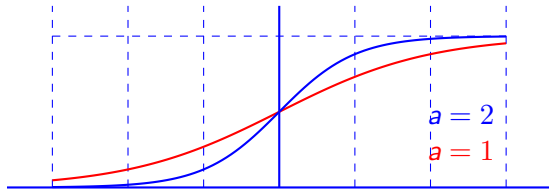
- Sigmoid function
 - Values lie between 0 and 1
 - Strictly increasing function
 - Bounded



$$\sigma(ax) = \frac{1}{1 + \exp(-a \times x)}$$

Activation function

- Sigmoid function
 - Values lie between 0 and 1
 - Strictly increasing function
 - Bounded

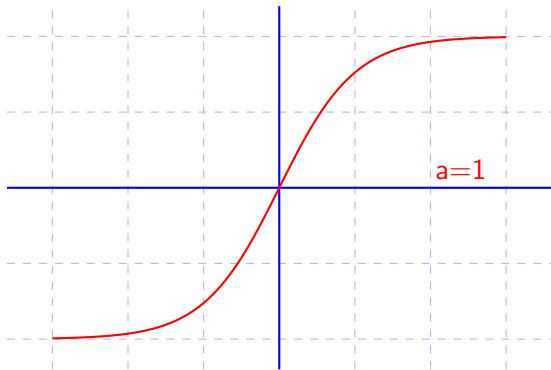


$$\sigma(ax) = \frac{1}{1 + \exp(-a \times x)}$$

Activation function

- Hyperbolic Tangent (Tanh) function

- Can be positive or negative
- Values lie between -1 and 1
- Strictly increasing function
- Bounded

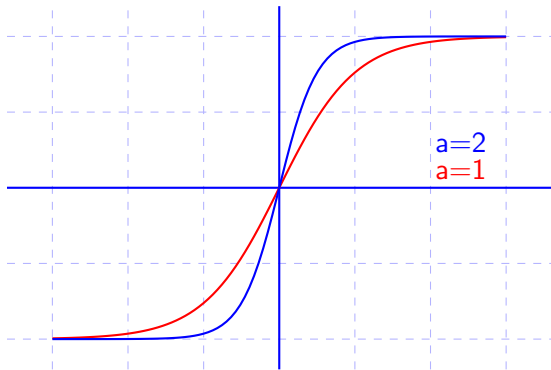


$$g(ax) = \tanh(ax) = \frac{\exp(ax) - \exp(-ax)}{\exp(ax) + \exp(-ax)}$$

Activation function

- Hyperbolic Tangent (Tanh) function

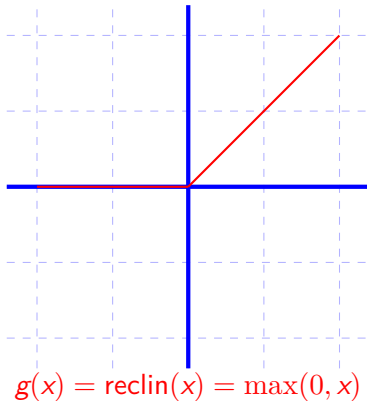
- Can be positive or negative
- Values lie between -1 and 1
- Strictly increasing function
- Bounded



$$g(ax) = \tanh(ax) = \frac{\exp(ax) - \exp(-ax)}{\exp(ax) + \exp(-ax)}$$

Activation function

- Rectified linear activation function (ReLU)
 - Bounded below by 0
 - Strictly increasing function
 - Not upper bounded



Generalization of ReLU

- ReLU is defined as $g(z) = \max\{0, z\}$
- Using non-zero slope, $h_i = g(z, \alpha)_i = \max(0, z_i) + \alpha_i \min(0, z_i)$
 - Absolute value rectification will make $\alpha_i = -1$ and $g(z) = |z|$
- Leaky ReLU assumes very small values for α_i
- Parametric ReLU tries to learn α_i parameters
- Maxout unit $g(z)_i = \max_{j \in \mathbb{G}^{(i)}} z_j$
 - Suitable for learning piecewise linear function

Logistic sigmoid & hyperbolic tangent

- Logistic sigmoid $g(z) = \sigma(z)$
- Hyperbolic tangent $g(z) = \tanh(z)$
 - $\tanh(z) = 2\sigma(2z) - 1$
- Widespread saturation of sigmoidal unit is an issue for gradient based learning
 - Usually discouraged to use as hidden units
- Usually, hyperbolic tangent function performs better where sigmoidal function must be used
 - Behaves linearly at 0
 - Sigmoidal activation function are more common in settings other than feedforward network

Other activation functions

- Differentiable functions are usually preferred
- Activation function $h = \cos(Wx + b)$ performs well for MNIST data set
- Sometimes no activation function helps in reducing the number of parameters
- Radial Basis Function - $\phi(x, c) = \phi(\|x - c\|)$
 - Gaussian - $\exp(-(\varepsilon r)^2)$
- Softplus - $g(x) = \zeta(x) = \log(1 + \exp(x))$
- Hard tanh - $g(x) = \max(-1, \min(1, x))$
- Hidden unit design is an active area of research

Hidden units

- Active area of research and does not have good guiding theoretical principle
- Usually rectified linear unit (ReLU) is chosen in most of the cases
- Design process consists of trial and error, then the suitable one is chosen
- Some of the activation functions are not differentiable (eg. ReLU)
 - Still gradient descent performs well
 - Neural network does not converge to local minima but reduces the value of cost function to a very small value

Output units

- Choice of output function usually depends on the type of problem being solved
- Usually **linear** function is chosen for regression and **sigmoid** for classification problems
- Any kind of output unit can be used as hidden unit

Linear units

- Suited for Gaussian output distribution
- Given features \mathbf{h} , linear output unit produces $\hat{y} = \mathbf{W}^T \mathbf{h} + \mathbf{b}$
- This can be treated as conditional probability $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}, \mathbf{I})$
- Maximizing log-likelihood is equivalent to minimizing mean square error

Sigmoid unit

- Mostly suited for binary classification problem that is Bernoulli output distribution
- The neural networks need to predict $p(y = 1|x)$
 - If linear unit has been chosen, $p(y = 1|x) = \max\{0, \min\{1, W^T h + b\}\}$
 - Gradient?
- Model should have strong gradient whenever the answer is wrong
- Let us assume unnormalized log probability is linear with $z = W^T h + b$
- Therefore, $\log \tilde{P}(y) = yz \Rightarrow \tilde{P}(y) = \exp(yz) \Rightarrow P(y) = \frac{\exp(yz)}{\sum_{y' \in \{0,1\}} \exp(y'z)}$
 - It can be written as $P(y) = \sigma((2y - 1)z)$
- The loss function for maximum likelihood is
$$J(\theta) = -\log P(y|x) = -\log \sigma((2y - 1)z) = \zeta((1 - 2y)z)$$

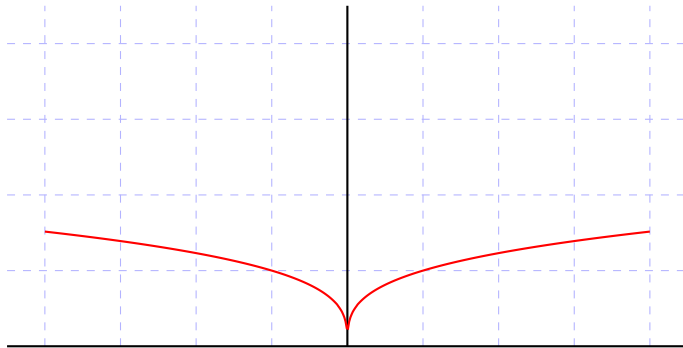
Softmax unit

- Similar to sigmoid. Mostly suited for multinoulli distribution
- We need to predict a vector \hat{y} such that $\hat{y}_i = P(Y = i|x)$
- A linear layer predicts unnormalized probabilities $z = W^T h + b$ that is $z_i = \log \tilde{P}(y = i|x)$
- Formally, $\text{softmax}(z)_i = \frac{\exp z_i}{\sum_j \exp(z_j)}$
- Log in log-likelihood can undo exp $\log \text{softmax}(z)_i = z_i - \log \sum_j \exp(z_j)$
 - Does it saturate?
 - What about incorrect prediction?
- Invariant to addition of some scalar to all input variables ie.
 $\text{softmax}(z) = \text{softmax}(z + c)$

Loss function

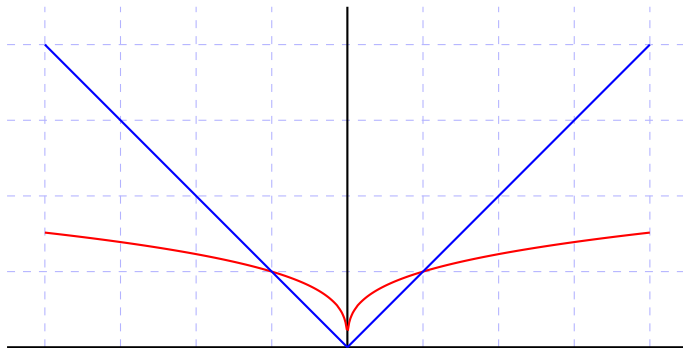
- Need to compare $\hat{y} = f(x)$ with the true label y for an input x
- For a single input example loss will be measured as $\mathcal{L}(y, f(x))$
- Average loss over a set of examples will be $\frac{1}{m} \sum_{i=1}^m \mathcal{L}(y_i, \hat{y}_i)$
- Target is to minimize the loss function
- Given the weights of the network W , the forward propagation yields $\hat{y}_i = f(x, W)$
- Our goal is as follows: $\underset{W}{\text{minimize}} \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y_i, f(x_i, W))$
- Generic loss function can have the following form $|y - a|^p$
- Euclidean norm $p = 2$

Loss curve



$$|y - a|^{0.3}, y=0$$

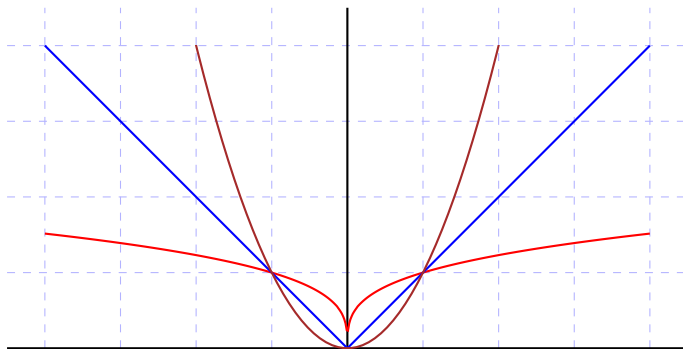
Loss curve



$$|y - a|^{1.0}, y=0$$

$$|y - a|^{0.3}, y=0$$

Loss curve

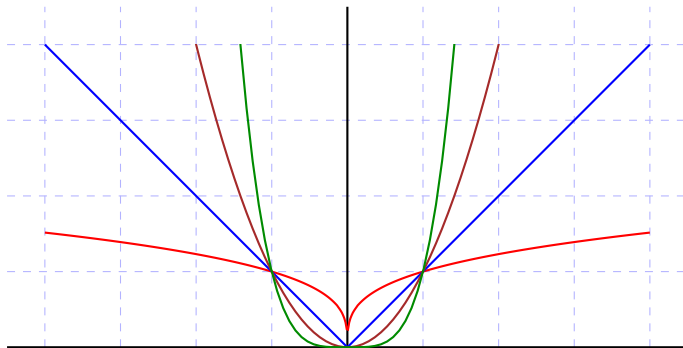


$$|y - a|^{2.0}, y=0$$

$$|y - a|^{1.0}, y=0$$

$$|y - a|^{0.3}, y=0$$

Loss curve



$$|y - a|^{4.0}, y=0$$

$$|y - a|^{2.0}, y=0$$

$$|y - a|^{1.0}, y=0$$

$$|y - a|^{0.3}, y=0$$

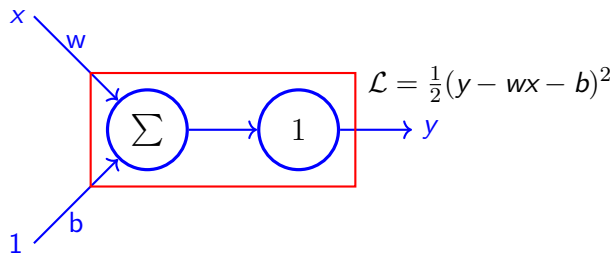
Linear regression

- Prediction of the value of a continuous variable
 - Example — price of a house, solar power generation in photo-voltaic cell, etc.

Linear regression

- Prediction of the value of a continuous variable
 - Example — price of a house, solar power generation in photo-voltaic cell, etc.
- Takes a vector $\mathbf{x} \in \mathbb{R}^n$ and predict scalar $y \in \mathbb{R}$
 - Predicted value will be represented as $\hat{y} = \mathbf{w}^T \mathbf{x}$ where \mathbf{w} is a vector of parameters
 - x_i receives positive weight — Increasing the value of the feature will increase the value of y
 - x_i receives negative weight — Increasing the value of the feature will decrease the value of y
 - Weight value is very high/large — Large effect on prediction

Linear regression using neural network



Performance

- Assume, we have m examples not used for training
 - This is known as test set

Performance

- Assume, we have m examples not used for training
 - This is known as test set
- Design matrix of inputs is $X^{(\text{test})}$ and target output is a vector $y^{(\text{test})}$
 - Performance is measured by Mean Square Error (MSE)

$$\text{MSE}_{(\text{test})} = \frac{1}{m} \sum_i \left(\hat{y}^{(\text{test})} - y^{(\text{test})} \right)_i^2 = \frac{1}{m} \|\hat{y}^{(\text{test})} - y^{(\text{test})}\|_2^2$$

- Error increases when the Euclidean distance between target and prediction increases

Performance

- Assume, we have m examples not used for training
 - This is known as test set
- Design matrix of inputs is $X^{(\text{test})}$ and target output is a vector $y^{(\text{test})}$
 - Performance is measured by Mean Square Error (MSE)

$$\text{MSE}_{(\text{test})} = \frac{1}{m} \sum_i \left(\hat{y}^{(\text{test})} - y^{(\text{test})} \right)_i^2 = \frac{1}{m} \|\hat{y}^{(\text{test})} - y^{(\text{test})}\|_2^2$$

- Error increases when the Euclidean distance between target and prediction increases
- The learning algorithm is allowed to gain experience from training set $(X^{(\text{train})}, y^{(\text{train})})$
- One of the common ideas is to minimize $\text{MSE}_{(\text{train})}$ for training set

Minimization of MSE

- We have the following now

$$\nabla_w \text{MSE}_{(\text{train})} = 0$$

Minimization of MSE

- We have the following now

$$\nabla_w \text{MSE}_{(\text{train})} = 0$$

$$\Rightarrow \nabla_w \frac{1}{m} \|\hat{y}^{(\text{train})} - y^{(\text{train})}\|_2^2 = 0$$

Minimization of MSE

- We have the following now

$$\nabla_w \text{MSE}_{(\text{train})} = 0$$

$$\Rightarrow \nabla_w \frac{1}{m} \|\hat{y}^{(\text{train})} - y^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \frac{1}{m} \nabla_w \|X^{(\text{train})} w - y^{(\text{train})}\|_2^2 = 0$$

Minimization of MSE

- We have the following now

$$\nabla_w \text{MSE}_{(\text{train})} = 0$$

$$\Rightarrow \nabla_w \frac{1}{m} \|\hat{y}^{(\text{train})} - y^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \frac{1}{m} \nabla_w \|X^{(\text{train})}_w - y^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \nabla_w (X^{(\text{train})}_w - y^{(\text{train})})^T (X^{(\text{train})}_w - y^{(\text{train})}) = 0$$

Minimization of MSE

- We have the following now

$$\nabla_w \text{MSE}_{(\text{train})} = 0$$

$$\Rightarrow \nabla_w \frac{1}{m} \|\hat{y}^{(\text{train})} - y^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \frac{1}{m} \nabla_w \|X^{(\text{train})} w - y^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \nabla_w (X^{(\text{train})} w - y^{(\text{train})})^T (X^{(\text{train})} w - y^{(\text{train})}) = 0$$

$$\Rightarrow \nabla_w (w^T X^{(\text{train})} X^{(\text{train})} w - 2w^T X^{(\text{train})} T y^{(\text{train})} + y^{(\text{train})} T y^{(\text{train})}) = 0$$

Minimization of MSE

- We have the following now

$$\nabla_w \text{MSE}_{(\text{train})} = 0$$

$$\Rightarrow \nabla_w \frac{1}{m} \|\hat{y}^{(\text{train})} - y^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \frac{1}{m} \nabla_w \|X^{(\text{train})} w - y^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \nabla_w (X^{(\text{train})} w - y^{(\text{train})})^T (X^{(\text{train})} w - y^{(\text{train})}) = 0$$

$$\Rightarrow \nabla_w (w^T X^{(\text{train})} X^{(\text{train})} w - 2w^T X^{(\text{train})} T y^{(\text{train})} + y^{(\text{train})} T y^{(\text{train})}) = 0$$

$$\Rightarrow 2X^{(\text{train})} X^{(\text{train})} w - 2X^{(\text{train})} T y^{(\text{train})} = 0$$

Minimization of MSE

- We have the following now

$$\nabla_w \text{MSE}_{(\text{train})} = 0$$

$$\Rightarrow \nabla_w \frac{1}{m} \|\hat{y}^{(\text{train})} - y^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \frac{1}{m} \nabla_w \|\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \nabla_w (\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})})^T (\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}) = 0$$

$$\Rightarrow \nabla_w (\mathbf{w}^T \mathbf{X}^{(\text{train})}{}^T \mathbf{X}^{(\text{train})} \mathbf{w} - 2 \mathbf{w}^T \mathbf{X}^{(\text{train})}{}^T \mathbf{y}^{(\text{train})} + \mathbf{y}^{(\text{train})}{}^T \mathbf{y}^{(\text{train})}) = 0$$

$$\Rightarrow 2 \mathbf{X}^{(\text{train})}{}^T \mathbf{X}^{(\text{train})} \mathbf{w} - 2 \mathbf{X}^{(\text{train})}{}^T \mathbf{y}^{(\text{train})} = 0$$

$$\Rightarrow \mathbf{w} = (\mathbf{X}^{(\text{train})}{}^T \mathbf{X}^{(\text{train})})^{-1} \mathbf{X}^{(\text{train})}{}^T \mathbf{y}^{(\text{train})}$$

Minimization of MSE

- We have the following now

$$\nabla_w \text{MSE}_{(\text{train})} = 0$$

$$\Rightarrow \nabla_w \frac{1}{m} \|\hat{y}^{(\text{train})} - y^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \frac{1}{m} \nabla_w \|\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \nabla_w (\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})})^T (\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}) = 0$$

$$\Rightarrow \nabla_w (\mathbf{w}^T \mathbf{X}^{(\text{train})}{}^T \mathbf{X}^{(\text{train})} \mathbf{w} - 2 \mathbf{w}^T \mathbf{X}^{(\text{train})}{}^T \mathbf{y}^{(\text{train})} + \mathbf{y}^{(\text{train})}{}^T \mathbf{y}^{(\text{train})}) = 0$$

$$\Rightarrow 2 \mathbf{X}^{(\text{train})}{}^T \mathbf{X}^{(\text{train})} \mathbf{w} - 2 \mathbf{X}^{(\text{train})}{}^T \mathbf{y}^{(\text{train})} = 0$$

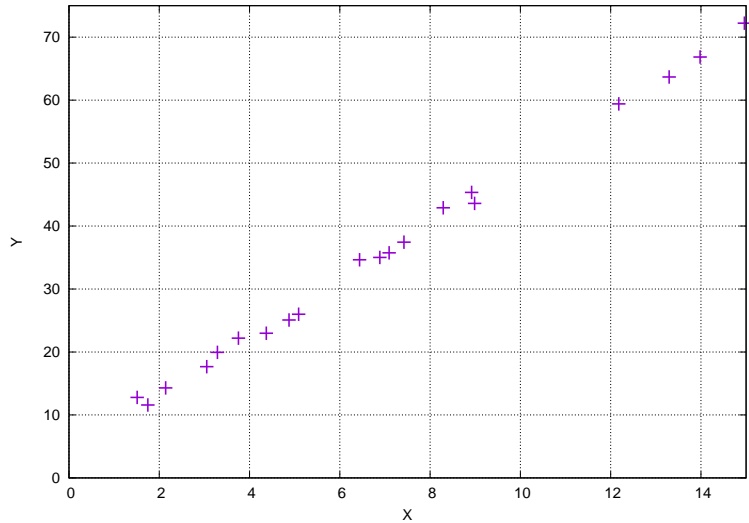
$$\Rightarrow \mathbf{w} = (\mathbf{X}^{(\text{train})}{}^T \mathbf{X}^{(\text{train})})^{-1} \mathbf{X}^{(\text{train})}{}^T \mathbf{y}^{(\text{train})}$$

- Linear regression with bias term $\hat{y} = [\mathbf{w}^T \quad w_0][\mathbf{x} \quad 1]^T$

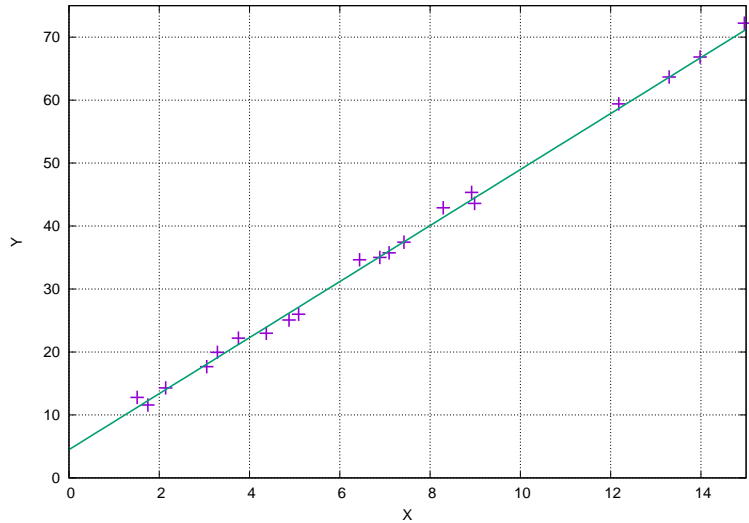
Moore-Penrose Pseudoinverse

- Let $A \in \mathbb{R}^{n \times m}$
- Every A has pseudoinverse $A^+ \in \mathbb{R}^{m \times n}$ and it is unique
 - $AA^+A = A$
 - $A^+AA^+ = A^+$
 - $(AA^+)^T = AA^+$
 - $(A^+A)^T = A^+A$
- $A^+ = \lim_{\alpha \rightarrow 0} (A^T A + \alpha I)^{-1} A^T$
- Example
 - If $A = \begin{bmatrix} 1 & 2 \end{bmatrix}^T$ then $A^+ = \begin{bmatrix} \frac{1}{5} & \frac{2}{5} \end{bmatrix}$
 - If $A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 1 & 5 \end{bmatrix}$ then $A^+ = \begin{bmatrix} 0.121212 & 0.515152 & -0.151515 \\ 0.030303 & -0.121212 & 0.212121 \end{bmatrix}$

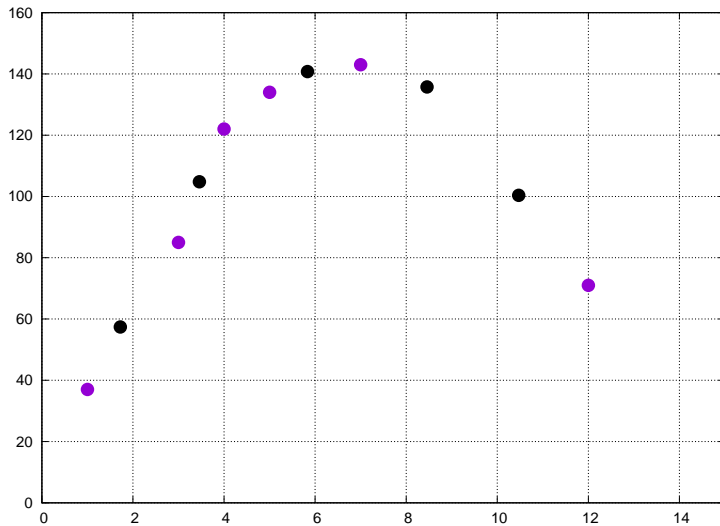
Regression example



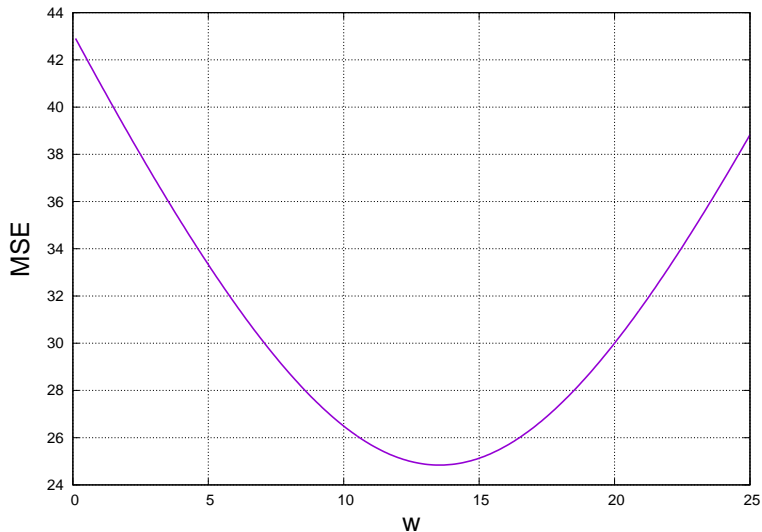
Regression example



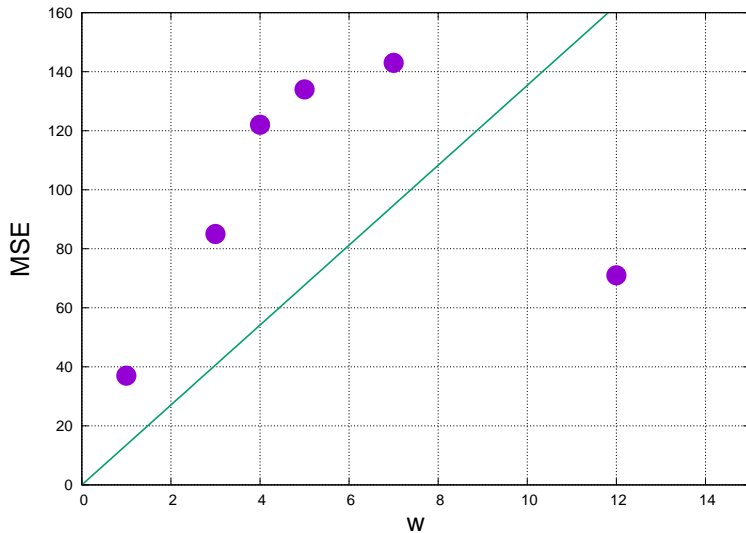
Example



Example: Variation of MSE wrt w



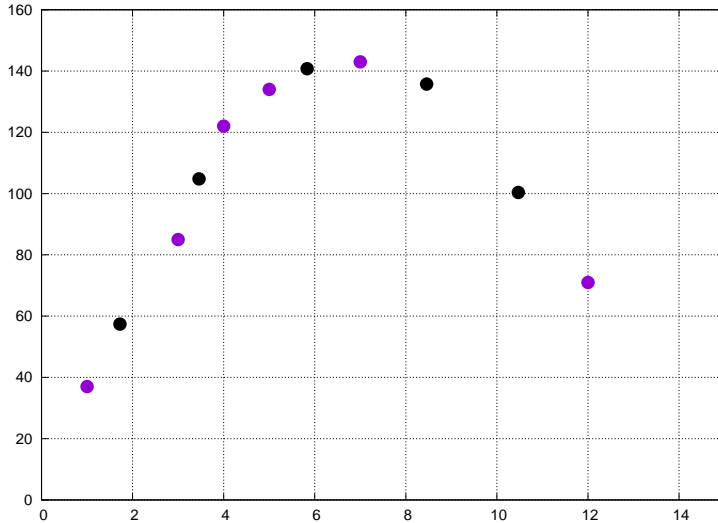
Example: Best fit



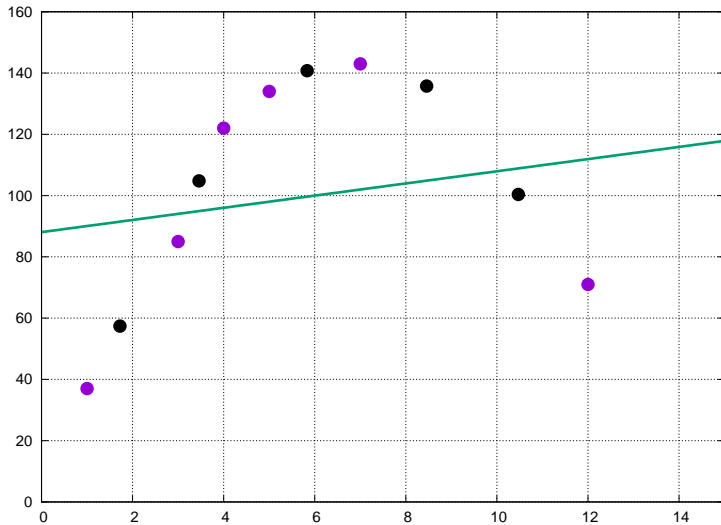
Error

- Training error - Error obtained on a training set
- Generalization error - Error on unseen data
- Data assumed to be independent and identically distributed (iid)
 - Each data set are independent of each other
 - Train and test data are identically distributed
- Expected training and test error will be the same
- It is more likely that the test error is greater than or equal to the expected value of training error
- Target is to make the training error is small. Also, to make the gap between training and test error smaller

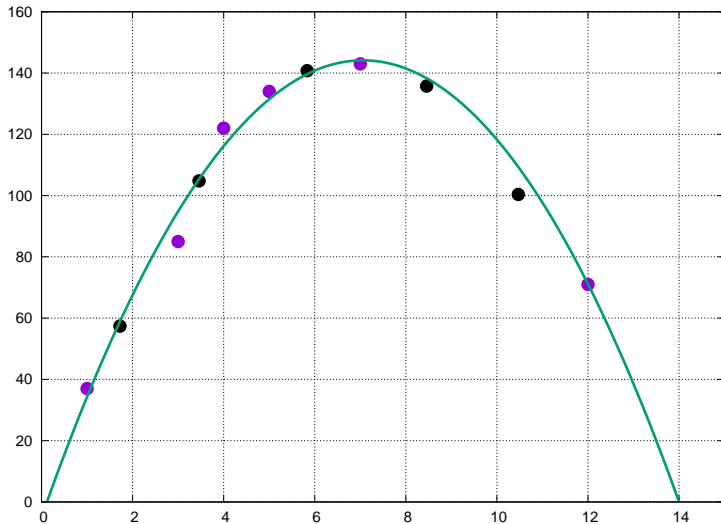
Regression example



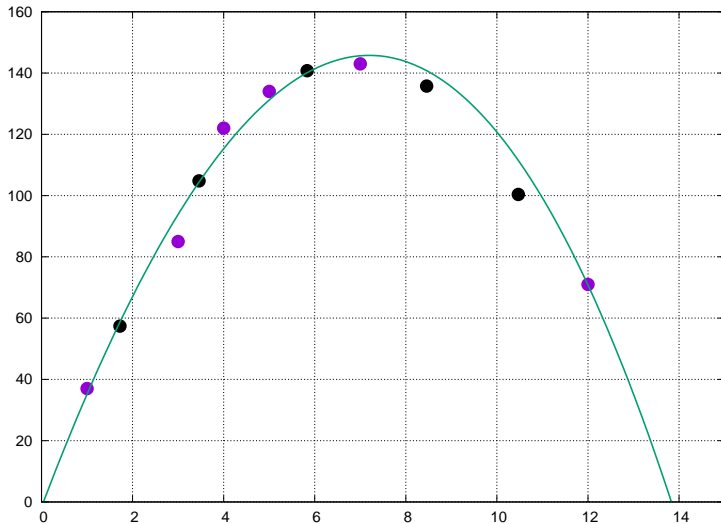
Regression example: degree 1



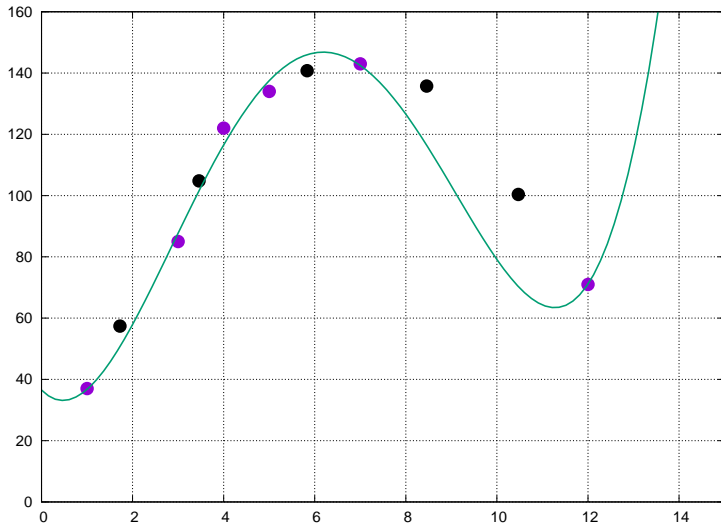
Regression example: degree 2



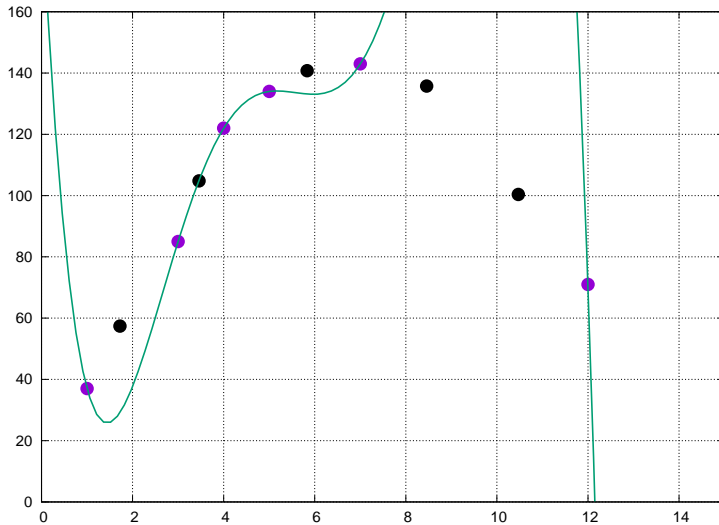
Regression example: degree 3



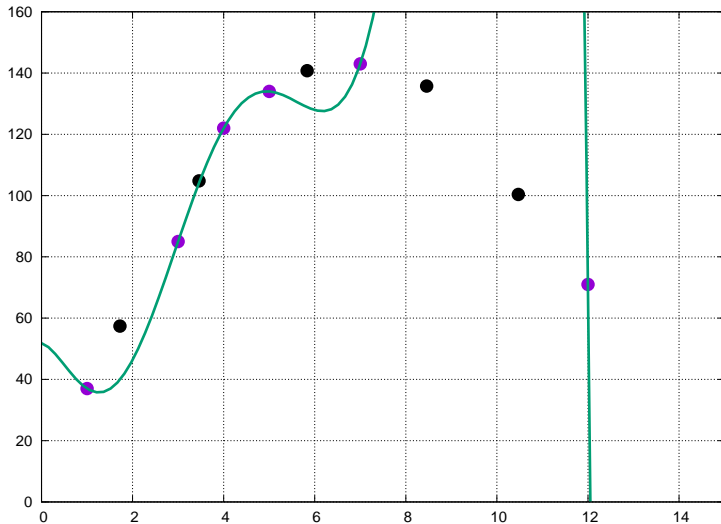
Regression example: degree 4



Regression example: degree 5



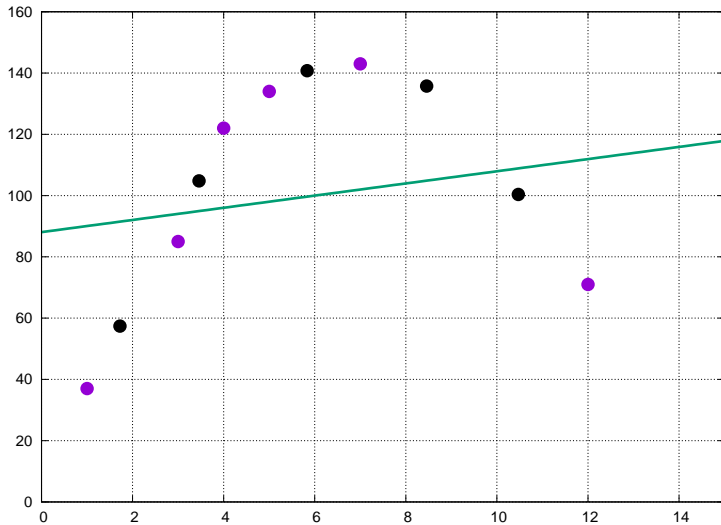
Regression example: degree 6



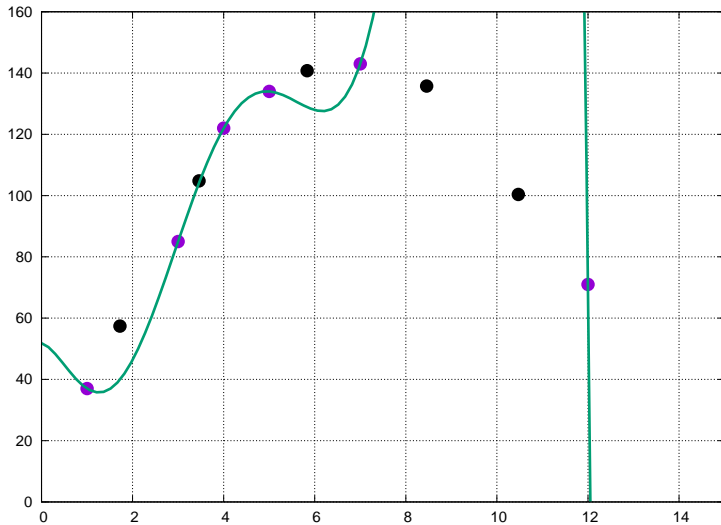
Underfitting & Overfitting

- Underfitting
 - When the model is not able to obtain sufficiently low error value on the training set
- Overfitting
 - When the gap between training set and test set error is too large

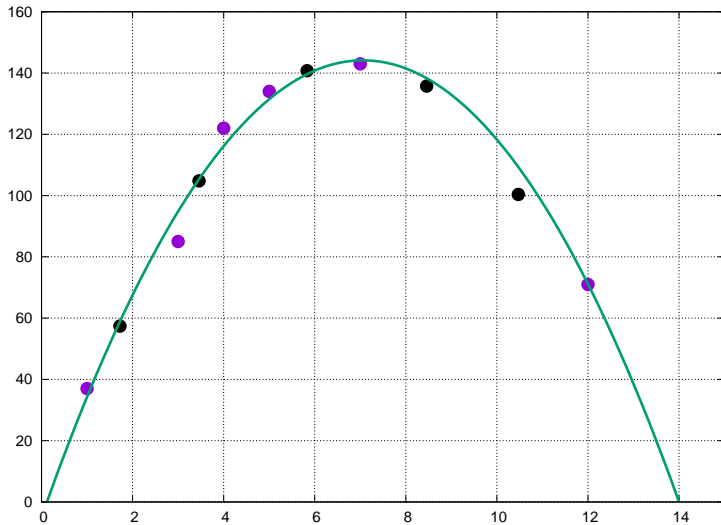
Underfitting example



Overfitting example



Better fit



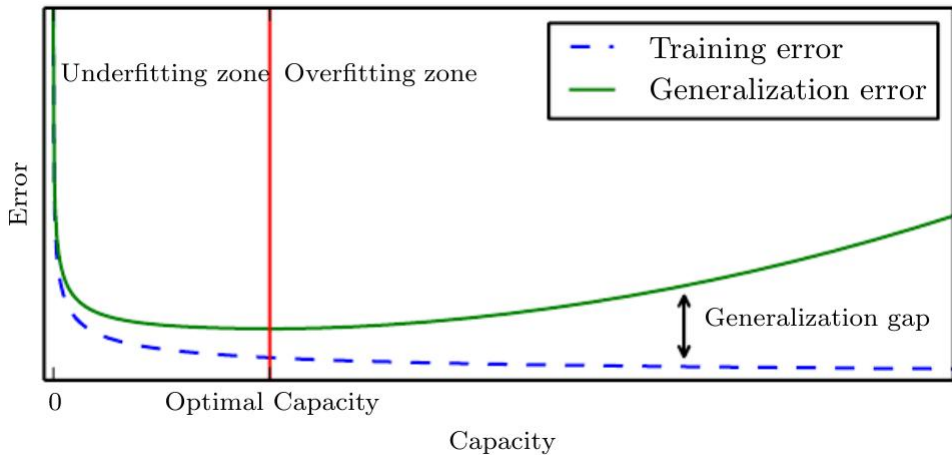
Capacity

- Ability to fit wide variety of functions
 - Low capacity will struggle to fit the training set
 - High capacity will can overfit by memorizing the training set
- Capacity can be controlled by choosing hypothesis space
 - A polynomial of degree 1 gives linear regression $\hat{y} = b + wx$
 - By adding x^2 term, it can learn quadratic curve $\hat{y} = b + w_1x + w_2x^2$
 - Output is still a linear function of parameters
- Capacity is determined by the choice of model (Representational capacity)
- Finding best function is very difficult optimization problem
 - Learning algorithm does not find the best function but reduces the training error
 - Imperfection in optimization algorithm can further reduce the capacity of model (effective capacity)

Capacity (contd.)

- Occam's razor
 - Among equally well hypotheses, choose the simplest one
- Vapnik-Chervonenski dimension - Capacity for binary classifier
 - Largest possible value of m for which a training set of m different \mathbf{x} points that the classifier can label arbitrarily
- Training and test error is bounded from above by a quantity that grows as model capacity grows but shrinks as the number of training example increases
 - Bounds are usually provided for ML algorithm and rarely provided for DL
 - Capacity of deep learning model is difficult as the effective capacity is limited by optimization algorithm
 - Little knowledge on non-convex optimization

Error vs Capacity



Non-parametric model

- Parametric model learns a function described by a parameter vector
 - Size of vector is finite and fixed
- Nearest neighbor regression
 - Finds out the nearest entry in training set and returns the associated value as the predicted one
 - Mathematically, for a given point x , $\hat{y} = y_i$ where $i = \arg \min \|X_{i,:} - x\|_2^2$
- Wrapping parametric algorithm inside another algorithm

Bayes error

- Ideal model is an oracle that knows the true probability distribution for data generation
- Such model can make error because of noise
 - Supervised learning
 - Mapping of x to y may be stochastic
 - y may be deterministic but x does not have all variables
- Error by an oracle in predicting from the true distribution is known as Bayes error

Note

- Training and generalization error varies as the size of training set varies
- Expected generalization error can never increase as the number of training example increases
- Any fixed parametric model with less than the optimal capacity will asymptote to an error value that exceeds the Bayes error
- It is possible to have optimal capacity but have large gap between training and generalization error
 - Need more training examples

No free lunch

- Averaged over all possible data generating distribution, every classification algorithm has same error rate when classifying unseen points
- No machine learning algorithm is universally any better than any other