



Introduction to Machine Learning

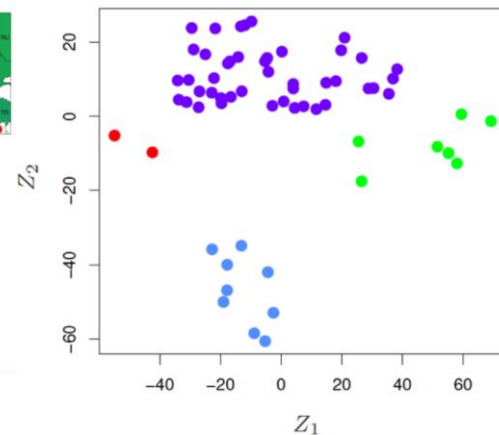
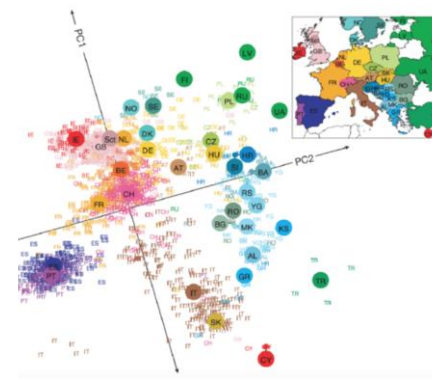
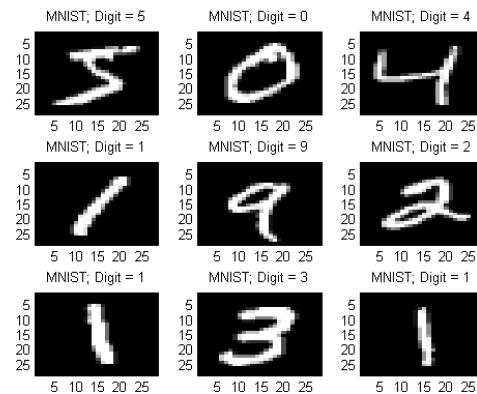
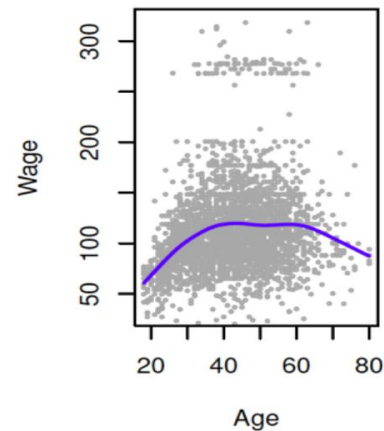


Outline

- Machine Learning Overview
 - Tasks of machine learning
 - Linear regression (regression)
 - Logistic regression (classification)
 - Support vector machine and kernels (classification)
 - Perceptron and neural network (classification or regression)
 - Decision tree (classification)
 - Unsupervised learning

Tasks of machine learning

- Parameterized regression
- Classification
- Dimensionality reduction
- Exploit data structure (unsupervised clustering)
- Many other methods: RNN, reinforcement learning and etc.



Machine learning key ideas

- Maximum likelihood $\mathbb{P}(\text{data}|\text{parameters})$
- Maximum posterior $\mathbb{P}(\text{parameters}|\text{data})$
- Bayesian rule

$$\underset{\text{posterior}}{\mathbb{P}(\text{parameters}|\text{data})} \propto \underset{\text{likelihood}}{\mathbb{P}(\text{data}|\text{parameters})} \times \underset{\text{prior}}{\mathbb{P}(\text{parameters})}$$

- Problem: Variance bias tradeoff (prevent over fitting)
 - Training data (find model)
 - Validation data (identify model)
 - Test data (make prediction and report result)

$$\underset{\text{loss}}{\mathbb{E}[(Y - \hat{Y}(X))^2]} = \underset{\text{inherent loss}}{\mathbb{E}[(Y - \mathbb{E}[Y|X])^2]} + \underset{(\text{bias})^2}{\mathbb{E}[(\mathbb{E}[Y|X] - \mathbb{E}[\hat{Y}(X)])^2]} + \underset{\text{variance}}{\mathbb{E}[(\hat{Y}(X) - \mathbb{E}[\hat{Y}(X)])^2]}$$

- Gradient descent $\mathbf{w}^{t+1} = \mathbf{w}^t - \alpha_t \nabla g(\mathbf{w}^t)$
 $\alpha_t \geq 0, \lim_{t \rightarrow \infty} \alpha_t = 0, \sum_t \alpha_t = \infty$

Linear Regression

- Approximate $f(\mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$ as

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + \text{constant} \\ &= w_{\text{TV}} x^{\text{TV}} + w_{\text{Radio}} x^{\text{Radio}} + w_{\text{NewsPaper}} x^{\text{NewsPaper}} + w_0 \end{aligned}$$

- Or, more generally

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_{i=0}^p w_i x_i, \text{ with } x_0 = 1$$

- Regression: find \mathbf{w} that minimizes

$$\sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \lambda \mathbf{w}^T \mathbf{w}$$

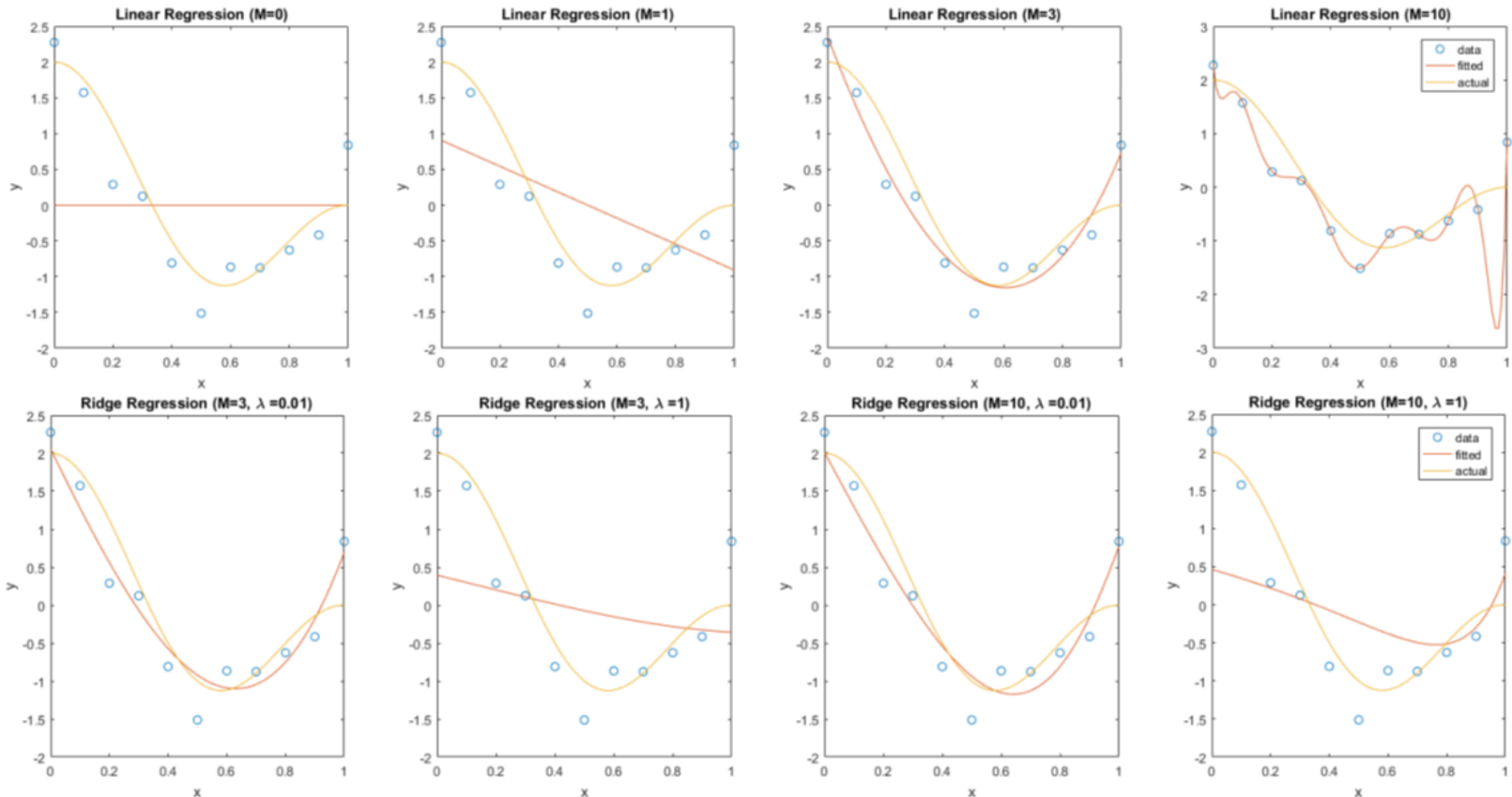
Let $Y^N = [y_n]^T$ and $X^N = [\mathbf{x}_n]^T$ regularization

Results

$$\begin{aligned} \nabla L(\mathbf{w}) &= -2X^T Y + 2(X^T X)\mathbf{w} \\ \mathbf{w}^* &= (X^T X)^{-1} X^T Y \end{aligned}$$

$$\begin{aligned} \nabla g(\mathbf{w}) &= -2X^T Y + 2X^T X\mathbf{w} + 2\lambda \mathbf{w} \\ \mathbf{w} &= (X^T X + \lambda \mathbf{I})^{-1} X^T Y \end{aligned}$$

Effect of Complex Model and Regularization



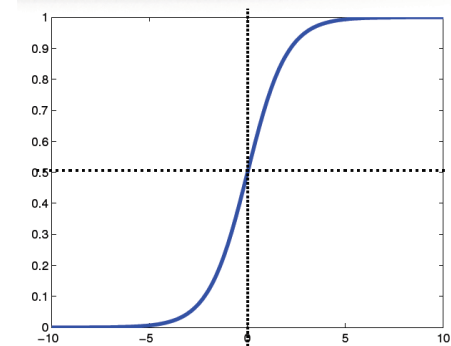
Logistic Regression

- Model probability with linear combinations

$$\log \frac{p(x)}{1 - p(x)} = \beta_0 + x \cdot \beta$$

- Obtain sigmoid function for prediction

$$\sigma(z) := \frac{1}{1 + e^{-z}} = \frac{e^z}{1 + e^z}$$



Sigmoid / logistic function

- Define cross entropy loss function

$$-\sum_{i=1}^N [y_i \log \sigma(w^T x_i) + (1 - y_i) \log (1 - \sigma(w^T x_i))]$$

- Transform to using label 1 and -1 to get easier loss function

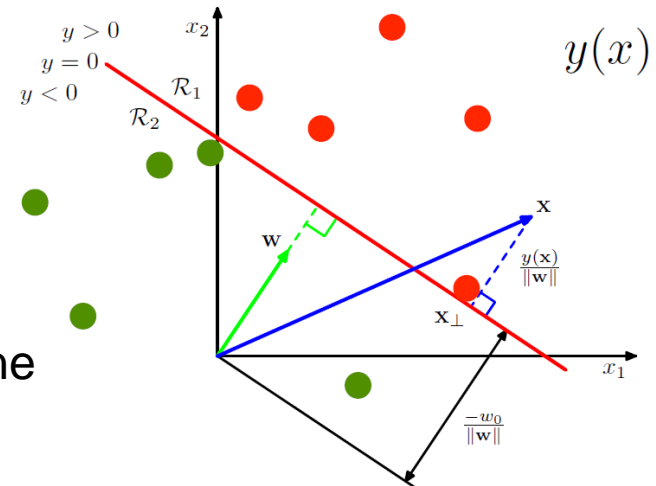
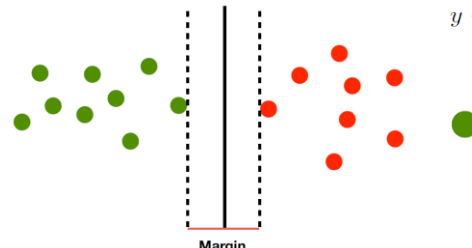
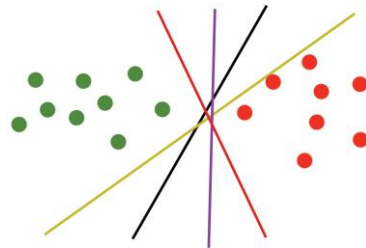
$$p(y = 1|x) = \frac{1}{1 + e^{-w^T x}}, \quad p(y = -1|x) = \frac{1}{1 + e^{+w^T x}} \quad \mathcal{L}(w) = \sum_{i=1}^N \log (1 + \exp(-y_i w^T x_i))$$

- Apply gradient descent and note that

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

Support Vector Machine

- Which line should we choose for separation?
- Idea: maximize the margin

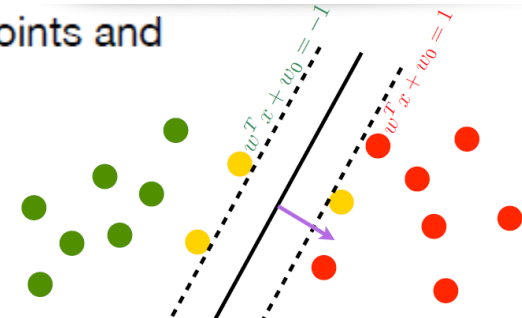


- Signed distance of a point to a hyperplane

$$\gamma = \frac{w^T x + w_0}{\|w\|}$$

There exists (w, w_0) such that $w^T x + w_0 > 0$ for positive points and $w^T x + w_0 < 0$ for negative points (assume data are 'bounded')

$$\min_{1 \leq i \leq N} |w^T x_i + w_0| = 1$$



Support Vector Machine

- How to maximize the margin?

Hard problem

$$\begin{aligned} \max_{w, w_0} \quad & \frac{1}{\|w\|} \\ \min_{1 \leq i \leq N} \quad & y_i(w^T x_i + w_0) = 1 \end{aligned}$$

Easier formulation

$$\begin{aligned} \min_{w, w_0} \quad & \frac{1}{2} \|w\|^2 \\ y_i(w^T x_i + w_0) & \geq 1, \quad 1 \leq i \leq N \end{aligned}$$

- Constrained optimization problem with Lagrange multiplier

$$L(w, w_0, \alpha) := \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i [y_i(w^T x_i + w_0) - 1].$$

(stationarity)

$$\frac{\partial L}{\partial w} = 0, \quad \frac{\partial L}{\partial w_0} = 0$$

(complementarity)

$$\alpha_i [y_i(w^T x_i + w_0) - 1] = 0, \quad \forall i.$$

(primal feasibility)

$$y_i(w^T x_i + w_0) \geq 1, \quad \forall i$$

(dual feasibility)

$$\alpha_i \geq 0, \quad \forall i$$

- Alternative dual form SVM

SVM dual

$$\begin{aligned} \max_{\alpha \geq 0} \quad & g(\alpha) := \min_{w, w_0} L(w, w_0, \alpha) \\ & = -\frac{1}{2} \left\| \sum_i \alpha_i y_i x_i \right\|^2 + \sum_i \alpha_i \\ & \sum_i y_i \alpha_i = 0. \end{aligned}$$

SVM Inseparable Case

- Include slack variable for inseparable case

Soft margin SVM

$$\begin{aligned} \min_{w, w_0, \xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_i \xi_i \\ y_i(w^T x_i + w_0) & \geq 1 - \xi_i, \quad 1 \leq i \leq N \\ \xi_i & \geq 0, \quad 1 \leq i \leq N \end{aligned}$$

Dual form

$$\begin{aligned} \max_{\alpha} \quad & -\frac{1}{2} \left\| \sum_i \alpha_i y_i x_i \right\|^2 + \sum_i \alpha_i \\ \sum_i y_i \alpha_i & = 0 \\ 0 & \leq \alpha \leq C. \end{aligned}$$

- How to identify support vectors

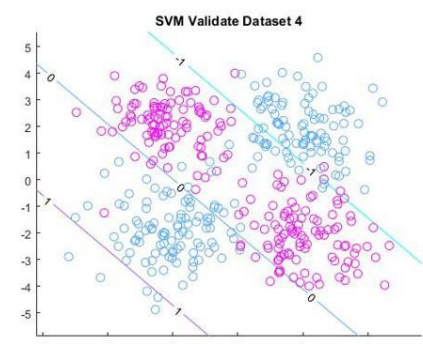
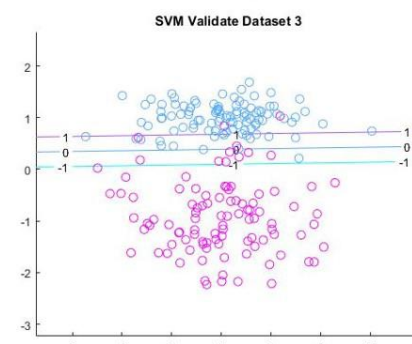
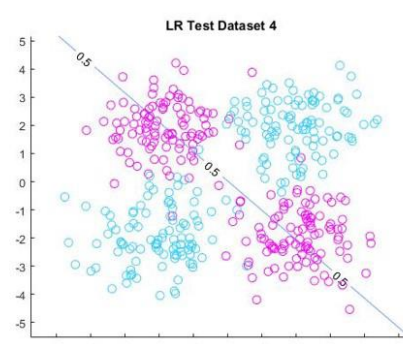
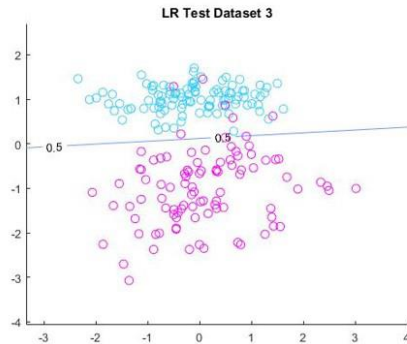
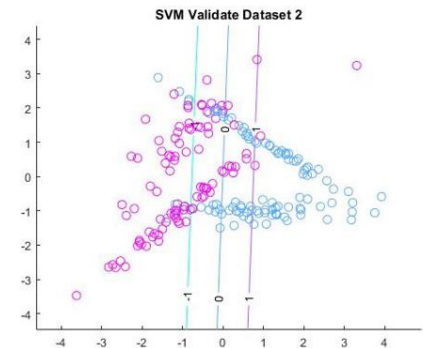
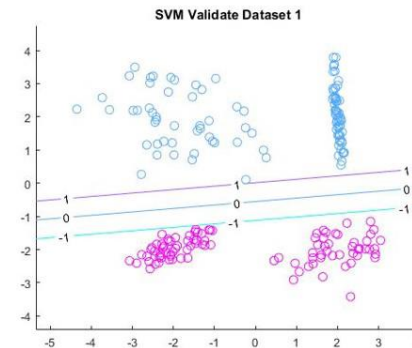
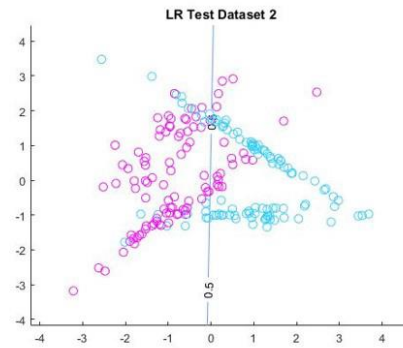
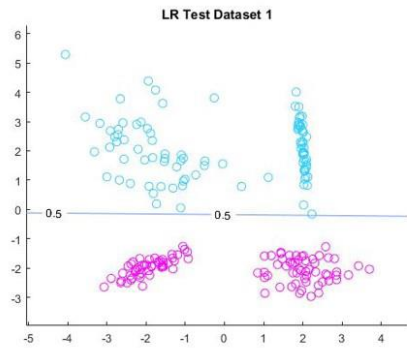
$$\begin{aligned} \alpha_i = 0 & \implies y_i(w^T x_i + w_0) \geq 1 && \text{(correctly classified)} \\ \alpha_i = C & \implies y_i(w^T x_i + w_0) \leq 1 && \text{(margin violation)} \\ 0 < \alpha_i < C & \implies y_i(w^T x_i + w_0) = 1 && \text{(support vector)} \end{aligned}$$

- Solution takes the form

$$w = \sum_i \alpha_i y_i \phi(x_i)$$



Classification of LR and SVM



Logistic regression classification

Linear SVM classification

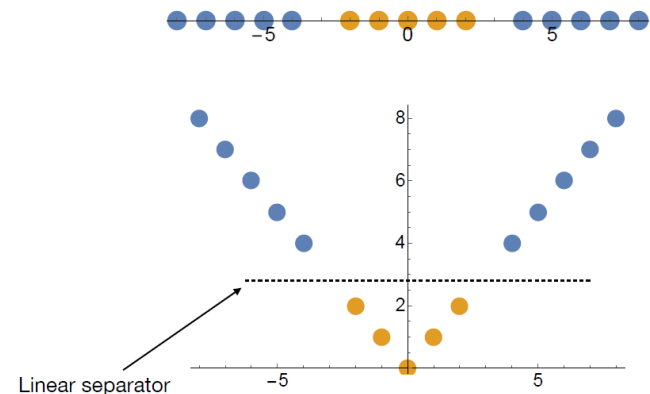
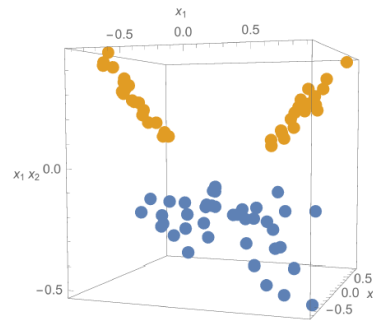
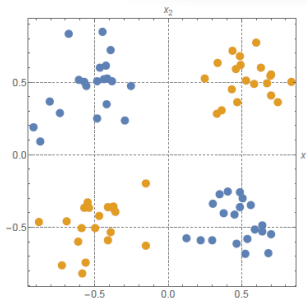
Kernel Trick

- Intuition: project data to high dimension and include non-linearity
- Key idea: Decision depends on just inner-products

Decision function

$$h(x) = \text{sgn}(\langle w, \phi(x) \rangle + w_0) \quad \langle w, \phi(x) \rangle = \sum_i \alpha_i y_i \langle \phi(x_i), \phi(x) \rangle$$

$$w = \sum_i \alpha_i y_i \phi(x_i) \quad k(x, x') = \langle \phi(x), \phi(x') \rangle$$

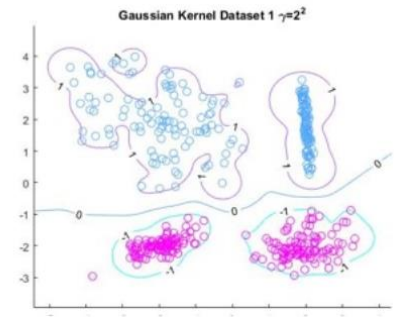
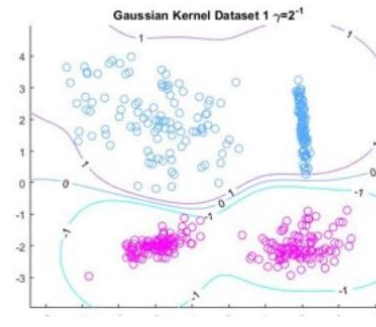
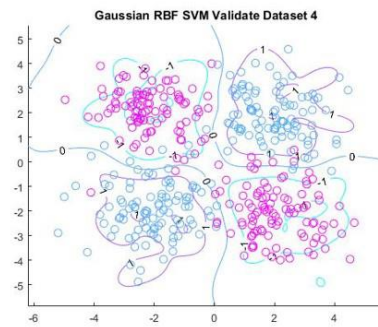
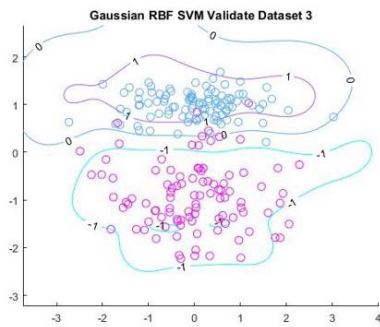
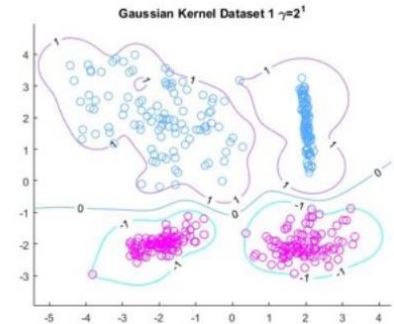
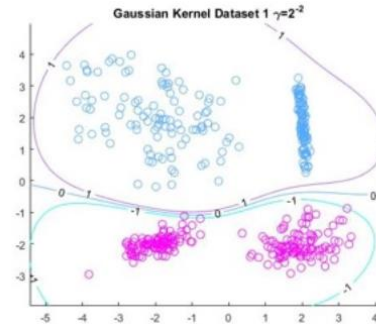
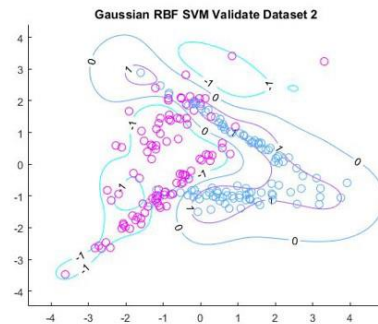
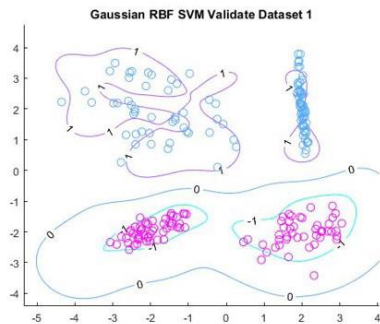


$$\phi(x) : (x_1, x_2) \mapsto (x_1, x_1x_2, x_2)$$

$$\phi(x) : x \mapsto (x, |x|)$$

Gaussian RBF Kernel

- Kernel formulation: $\exp(-\lambda \|x - x'\|^2)$

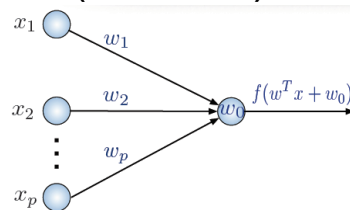


Performance on 4 datasets

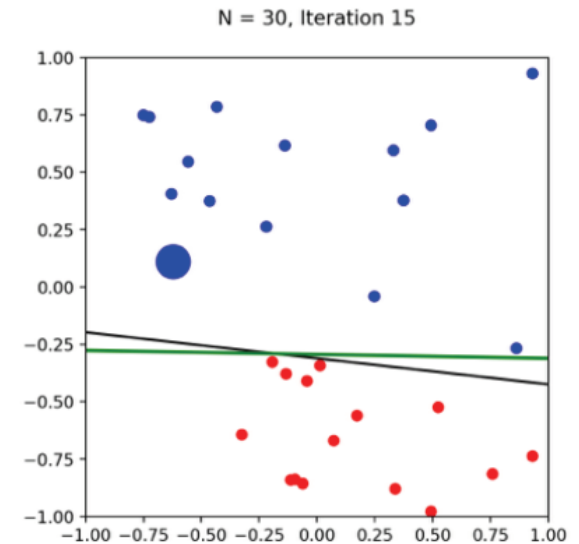
Effect of different λ

Perceptron Algorithm

- Idea: keep adjusting when miss classified
- 1 neuron neural (network)



- Converges when data is linearly separable



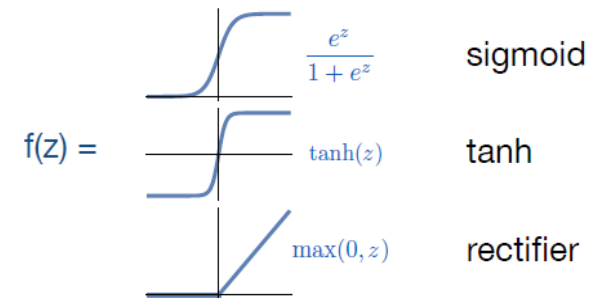
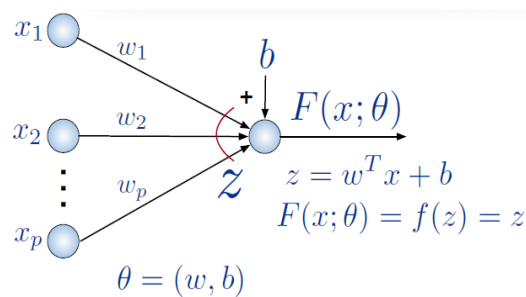
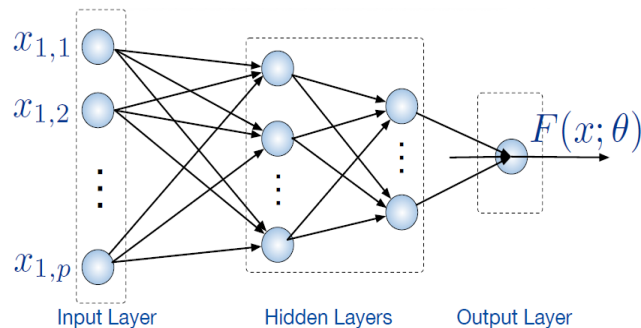
Algorithm:

1. Initialize parameters; set iteration counter $t = 1$.
2. Cycle through training data $(x_1, y_1), \dots, (x_N, y_N)$ and update

$$\left. \begin{array}{l} \text{if } y_i \neq h(x_i; w^t, w_0^t), \text{ then} \\ w^{t+1} = w^t + y_i x_i \\ w_0^{t+1} = w_0^t + y_i \end{array} \right\} y_i(x_i^T w^t + w_0^t) \leq 0$$

Neural Network

- Architecture: input layer, hidden layers, output layer, loss function
- Individual neuron: weighed sum, bias, activation function (non-linear)
- Very powerful tool that can accomplish complex tasks
- Can do both classification and regression
- Multiple variations available: Convolution NN, Recurrent NN, etc.



Back Propagation

- Back propagation to compute gradient with chain rule
- Can be hard with many weights to compute
- Back propagation saves derivative results in the middle
- We trade memory storage for speed

$$\frac{\partial \ell}{\partial W^l} = \frac{\partial z^l}{\partial W^l} \left[\frac{\partial \ell}{\partial z^l} \right] \quad \delta^l = \frac{\partial \ell}{\partial z^l} = \frac{\partial z^{l+1}}{\partial z^l} \cdot \frac{\partial \ell}{\partial z^{l+1}} = \frac{\partial z^{l+1}}{\partial z^l} \cdot \delta^{l+1} \quad \frac{\partial z^{l+1}}{\partial z^l} = \text{Diag}[f'(z^l)]W^{l+1}$$

$$\delta^l = \text{Diag}[f'(z^l)]W^{l+1} \text{Diag}[f'(z^{l+1})]W^{l+2} \dots W^L \delta^L$$

1. **Input:** (x, y) set activations for a^1 input layer

2. **Feedforward:**

for each layer $l = 2, \dots, L$ **do**

$$z^l = (W^l)^T a^{l-1} + b^l$$

$$a^l = f(z^l)$$

3. **Output layer:** $\delta^L = \mathcal{D}[f'(z^L)] \nabla_a \text{loss}$

4. **Backpropagation:**

for each $l = L - 1, \dots, 2$ **do**

$$\delta^l = \mathcal{D}[f'(z^l)]W^{l+1} \delta^{l+1}$$

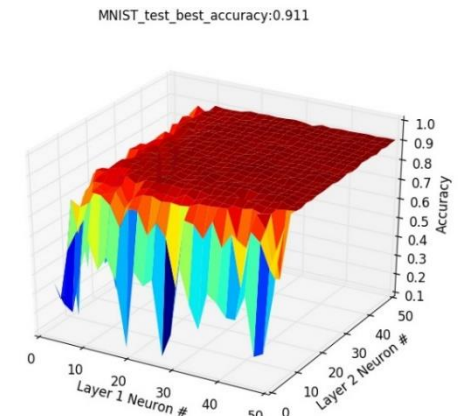
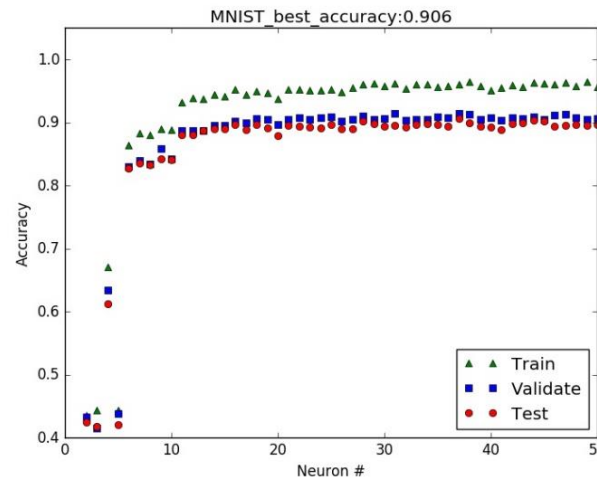
5. **Final gradients:**

$$\frac{\partial \text{loss}}{\partial W^l} = a^{l-1} \delta^l, \quad \frac{\partial \text{loss}}{\partial b^l} = \delta^l$$

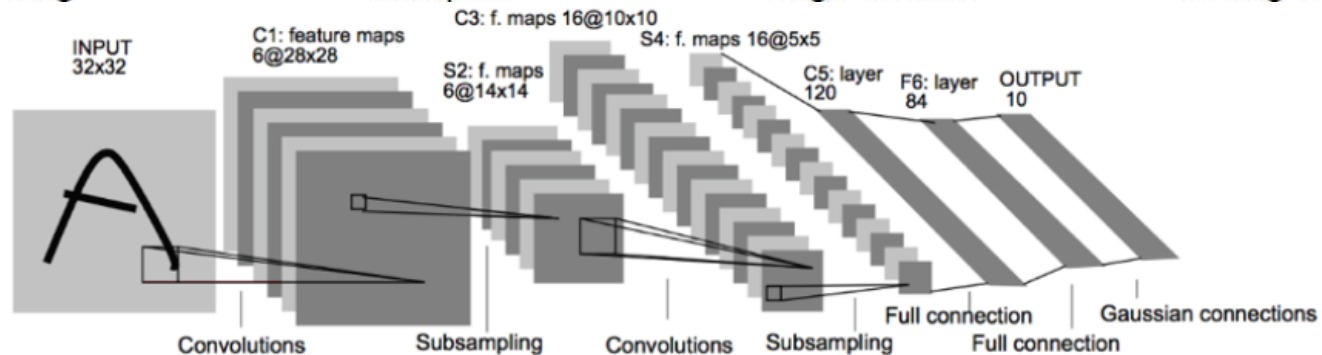
NN Performance on MNIST

- Classification for MNIST at 90.6% accuracy with 1 layer
- 2 layer accuracy: 91.1%
- Exercise: use python NN implementation for MNIST dataset

MNIST Samples

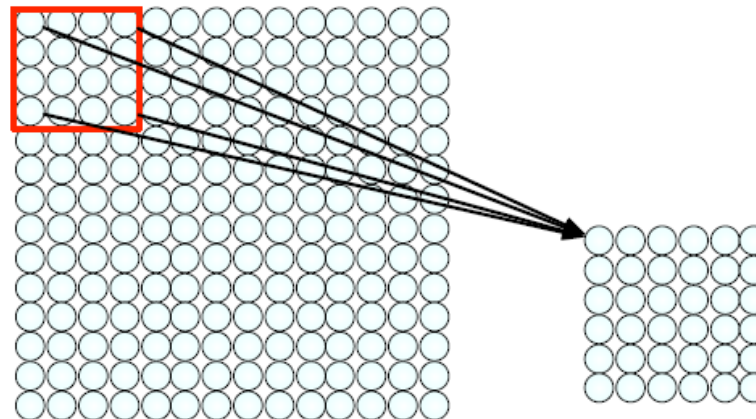


Convolutional Neural Network



Convolutional Neural Network

- Majorly used in image classification applications
- Utilize special information from image to reduce parameters
- Convolution: apply filter moving through each pixel
- Filter can be constructed manually to achieve various effects
- Learn filter parameters using CNN instead of manually constructing

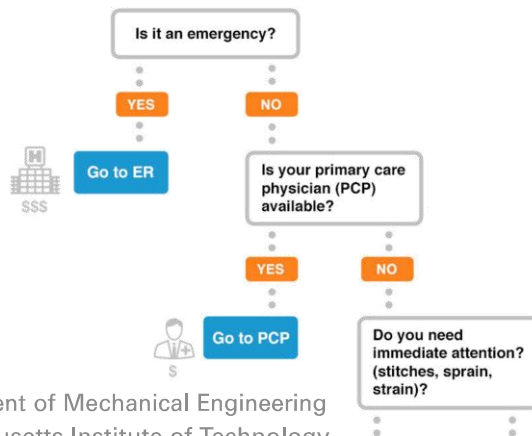


Decision Tree

- Simple heuristic in real life for making decisions
- Need formalization on how we split data (greedy algorithm)
- Define (conditional) entropy and maximize information gain

$$H(Y | X) = - \sum_{j=1}^v P(X = x_j) \sum_{i=1}^k P(Y = y_i | X = x_j) \log_2 P(Y = y_i | X = x_j)$$

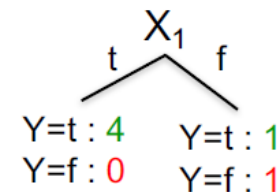
$$H(Y) = - \sum_{i=1}^k P(Y = y_i) \log_2 P(Y = y_i) \quad IG(X) = H(Y) - H(Y | X)$$



Example:

$$P(X_1=t) = 4/6$$

$$P(X_1=f) = 2/6$$



$$H(Y|X_1) = - 4/6 (1 \log_2 1 + 0 \log_2 0) - 2/6 (1/2 \log_2 1/2 + 1/2 \log_2 1/2) = 2/6$$

X ₁	X ₂	Y
T	T	T
T	F	T
T	T	T
T	F	T
F	T	T
F	F	F

Unsupervised learning

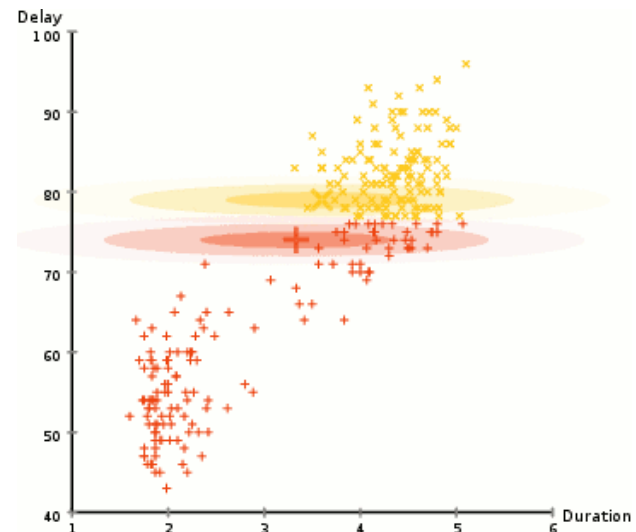
- When no data labels are available
- K-means clustering
- EM Algorithm: compute expectation based on parameter to assign labels, maximize probability by adjusting parameter values

– **Initialize:** Pick K random points as cluster centers

– **Alternate:**

1. Assign data points to closest cluster center
2. Change the cluster center to the average of its assigned points

– **Stop** when no points' assignments change

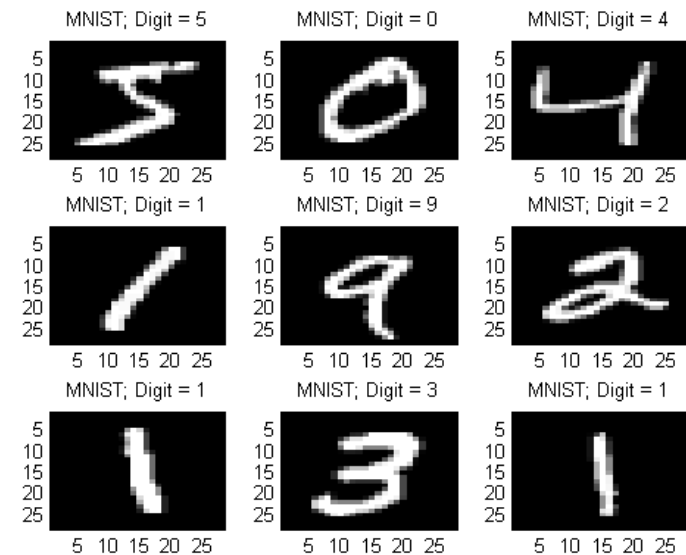
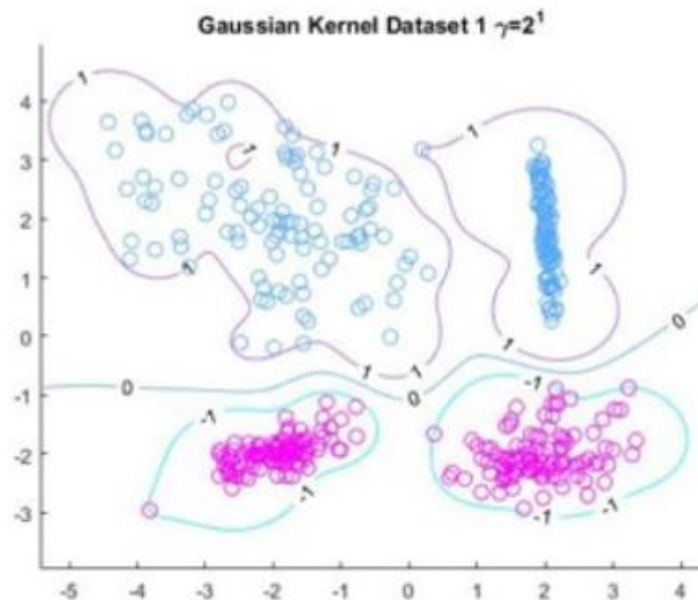


Machine Learning Summary

- Very active research field for both algorithm and hardware
- Many other techniques exist
 - Recursive neural network for language processing
 - Latent Dirichlet allocation for topic group modeling
 - Reinforcement learning for game play
 - Matrix data principle component analysis for dimensionality reduction
- Many good packages: Scikit-learn, TensorFlow, Pytorch, etc.
- Many dataset available online (e.g. [UC Ivrine Repository](#))
- Fun to plan around and apply to various tasks
- Many resources online to learn in more details
 - Udacity, Coursera, Lynda, Packt Publishing
 - [MIT OCW](#), [Stanford](#),

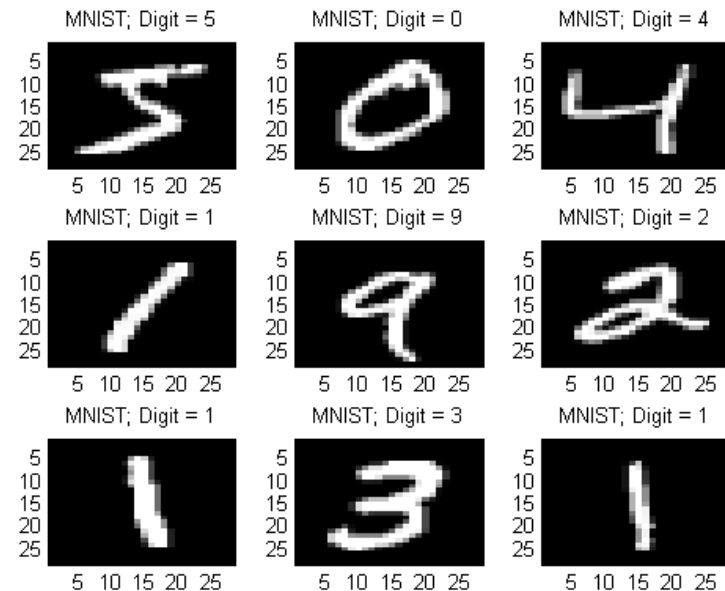
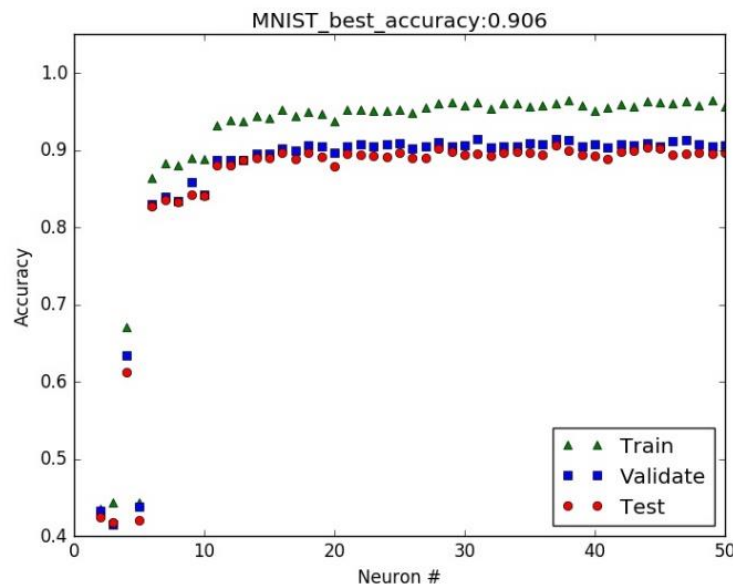
Demo

- Matlab SVM for 2 Dataset and MNIST classification
- Python Neural Network for MNIST classification



Exercise

- Implement 1 layer neural network for MNIST multiclass classification
- Report your result for 1 hidden layer NN with 0 to 50 neurons





Thank You!