Support Vector Machines and Kernel Methods

Weinan Zhang Shanghai Jiao Tong University http://wnzhang.net

http://wnzhang.net/teaching/cs420/index.html

References and Acknowledgement





- A large part of slides in this lecture are originally from Prof. Andrew Ng's lecture at Stanford University
 - http://cs229.stanford.edu/notes/cs229-notes3.pdf
 - http://www.andrewng.org/

• For linear separable cases, we have multiple decision boundaries



• For linear separable cases, we have multiple decision boundaries



• Ruling out some separators by considering data noise

• For linear separable cases, we have multiple decision boundaries



• The intuitive optimal decision boundary: the largest margin

Review: Logistic Regression

• Logistic regression is a binary classification model

$$p_{\theta}(y=1|x) = \sigma(\theta^{\top}x) = \frac{1}{1+e^{-\theta^{\top}x}}$$
$$p_{\theta}(y=0|x) = \frac{e^{-\theta^{\top}x}}{1+e^{-\theta^{\top}x}}$$



Cross entropy loss function

$$\mathcal{L}(y, x, p_{\theta}) = -y \log \sigma(\theta^{\top} x) - (1 - y) \log(1 - \sigma(\theta^{\top} x))$$

Gradient

$$\frac{\partial \mathcal{L}(y, x, p_{\theta})}{\partial \theta} = -y \frac{1}{\sigma(\theta^{\top} x)} \sigma(z)(1 - \sigma(z))x - (1 - y) \frac{-1}{1 - \sigma(\theta^{\top} x)} \sigma(z)(1 - \sigma(z))x$$
$$= (\sigma(\theta^{\top} x) - y)x$$
$$\theta \leftarrow \theta + \eta(y - \sigma(\theta^{\top} x))x$$
$$\boxed{\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z))}$$

Label Decision

• Logistic regression provides the probability

$$p_{\theta}(y=1|x) = \sigma(\theta^{\top}x) = \frac{1}{1+e^{-\theta^{\top}x}}$$
$$p_{\theta}(y=0|x) = \frac{e^{-\theta^{\top}x}}{1+e^{-\theta^{\top}x}}$$

- The final label of an instance is decided by setting a threshold \boldsymbol{h}

$$\hat{y} = \begin{cases} 1, & p_{\theta}(y=1|x) > h \\ 0, & \text{otherwise} \end{cases}$$

Logistic Regression Scores



The higher score, the larger distance to the decision boundary, the higher confidence

Example from Andrew Ng

• The intuitive optimal decision boundary: the highest confidence



Notations for SVMs

- Feature vector \boldsymbol{x}
- Class label $y \in \{-1, 1\}$
- Parameters
 - Intercept b
 - Feature weight vector w
- Label prediction

$$h_{w,b}(x) = g(w^{\top}x + b)$$
$$g(z) = \begin{cases} +1 & z \ge 0\\ -1 & \text{otherwise} \end{cases}$$

Logistic Regression Scores



The higher score, the larger distance to the separating hyperplane, the higher confidence Example from Andrew Ng

Margins



• Functional margin

$$\hat{\gamma}^{(i)} = y^{(i)}(w^{\top}x^{(i)} + b)$$

 Note that the separating hyperplane won't change with the magnitude of (w, b)

 $g(w^{\top}x+b) = g(2w^{\top}x+2b)$

Geometric margin

$$\gamma^{(i)} = y^{(i)}(w^{\top}x^{(i)} + b)$$

where $||w||^2 = 1$

Margins



• Decision boundary

$$w^{\top} \left(x^{(i)} - \gamma^{(i)} y^{(i)} \frac{w}{\|w\|} \right) + b = 0$$

$$\gamma^{(i)} = y^{(i)} \frac{w^{\top} x^{(i)} + b}{\|w\|}$$
$$= y^{(i)} \left(\left(\frac{w}{\|w\|} \right)^{\top} x^{(i)} + \frac{b}{\|w\|} \right)$$

Given a training set

$$S = \{(x_i, y_i)\}_{i=1...m}$$

the smallest geometric margin

$$\gamma = \min_{i=1\dots m} \gamma^{(i)}$$

Objective of an SVM

• Find a separable hyperplane that maximizes the minimum geometric margin

$$egin{aligned} &\max&\gamma\ \gamma,w,b&\gamma\ ext{s.t.}&y^{(i)}(w^ op x^{(i)}+b)\geq\gamma,\,\,i=1,\ldots,m\ &\|w\|=1 \quad ext{(non-convex constraint)} \end{aligned}$$

• Equivalent to normalized functional margin

$$\begin{array}{ll} \max_{\hat{\gamma},w,b} & \frac{\hat{\gamma}}{\|w\|} & \text{(non-convex objective)} \\ \text{s.t.} & y^{(i)}(w^{\top}x^{(i)}+b) \geq \hat{\gamma}, \, \, i=1,\ldots,m \end{array}$$

Objective of an SVM

- Functional margin scales w.r.t. (*w,b*) without changing the decision boundary.
 - Let's fix the functional margin at 1.

$$\hat{\gamma} = 1$$

• Objective is written as

$$\max_{\substack{w,b \\ w,b}} \frac{1}{\|w\|}$$

s.t. $y^{(i)}(w^{\top}x^{(i)}+b) \ge 1, \ i = 1, \dots, m$

• Equivalent with

$$\min_{\substack{w,b \ x,b \ x}} \frac{1}{2} \|w\|^2$$

s.t. $y^{(i)}(w^{\top} x^{(i)} + b) \ge 1, \ i = 1, \dots, m$

This optimization problem can be efficiently solved by quadratic programming

A Digression of Lagrange Duality in Convex Optimization

Boyd, Stephen, and Lieven Vandenberghe. Convex optimization. Cambridge university press, 2004.

Lagrangian for Convex Optimization

A convex optimization problem

$$\min_{w} f(w)$$

s.t. $h_i(w) = 0, \quad i = 1, \dots, l$

• The Lagrangian of this problem is defined as

$$\mathcal{L}(w,\beta) = f(w) + \sum_{i=1}^{l} \beta_i h_i(w)$$
 Lagrangian multipliers

• Solving

$$\frac{\partial \mathcal{L}(w,\beta)}{\partial w} = 0 \qquad \frac{\partial \mathcal{L}(w,\beta)}{\partial \beta} = 0$$

yields the solution of the original optimization problem.

Lagrangian for Convex Optimization



With Inequality Constraints

• A convex optimization problem

$$\min_{w} f(w)$$
s.t. $g_i(w) \le 0, \quad i = 1, \dots, k$
 $h_i(w) = 0, \quad i = 1, \dots, l$

• The Lagrangian of this problem is defined as

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^{k} \alpha_i g_i(w) + \sum_{i=1}^{l} \beta_i h_i(w)$$

Lagrangian multipliers

Primal Problem

- A convex optimization
 - $egin{array}{ll} \min_w & f(w) \ {
 m s.t.} & g_i(w) \leq 0, & i=1,\ldots,k \ & h_i(w)=0, & i=1,\ldots,l \end{array}$

$$\mathcal{L}(w,lpha,eta)=f(w)+\sum_{i=1}^klpha_ig_i(w)+\sum_{i=1}^leta_ih_i(w)$$

The Lagrangian

• The primal problem

$$\theta_{\mathcal{P}}(w) = \max_{\alpha,\beta:\alpha_i \ge 0} \mathcal{L}(w,\alpha,\beta)$$

• If a given \boldsymbol{w} violates any constraints, i.e.,

 $g_i(w) > 0$ or $h_i(w) \neq 0$

• Then $\theta_{\mathcal{P}}(w) = +\infty$

Primal Problem

- A convex optimization
 - $egin{aligned} \min_w & f(w) \ ext{s.t.} & g_i(w) \leq 0, & i=1,\ldots,k \ & h_i(w)=0, & i=1,\ldots,l \end{aligned}$

$$\mathcal{L}(w,lpha,eta)=f(w)+\sum_{i=1}^klpha_i g_i(w)+\sum_{i=1}^leta_i h_i(w)$$

The Lagrangian

• The primal problem

$$\theta_{\mathcal{P}}(w) = \max_{\alpha,\beta:\alpha_i \ge 0} \mathcal{L}(w,\alpha,\beta)$$

- Conversely, if all constraints are satisfied for $\,w\,$
- Then $\theta_{\mathcal{P}}(w) = f(w)$

 $\theta_{\mathcal{P}}(w) = \begin{cases} f(w) & \text{if } w \text{ satisfies primal constraints} \\ +\infty & \text{otherwise} \end{cases}$

Primal Problem

 $\theta_{\mathcal{P}}(w) = \begin{cases} f(w) & \text{if } w \text{ satisfies primal constraints} \\ +\infty & \text{otherwise} \end{cases}$

• The minimization problem

$$\min_{w} \theta_{\mathcal{P}}(w) = \min_{w} \max_{\alpha,\beta:\alpha_i \ge 0} \mathcal{L}(w,\alpha,\beta)$$

is the same as the original problem

$$\min_{w} f(w)$$
s.t. $g_i(w) \le 0, \quad i = 1, \dots, k$
 $h_i(w) = 0, \quad i = 1, \dots, l$

• Define the value of the primal problem $p^* = \min_w heta_\mathcal{P}(w)$

Dual Problem

• A slightly different problem

$$\theta_{\mathcal{D}}(\alpha,\beta) = \min_{w} \mathcal{L}(w,\alpha,\beta)$$

• Define the dual optimization problem

$$\max_{\alpha,\beta:\alpha_i\geq 0}\theta_{\mathcal{D}}(\alpha,\beta) = \max_{\alpha,\beta:\alpha_i\geq 0}\min_{w}\mathcal{L}(w,\alpha,\beta)$$

• Min & Max exchanged compared to the primal problem

$$\min_{w} \theta_{\mathcal{P}}(w) = \min_{w} \max_{\alpha,\beta:\alpha_i \ge 0} \mathcal{L}(w,\alpha,\beta)$$

• Define the value of the dual problem

$$d^* = \max_{\alpha,\beta:\alpha_i \ge 0} \min_{w} \mathcal{L}(w,\alpha,\beta)$$

Primal Problem vs. Dual Problem

- $d^* = \max_{\alpha,\beta:\alpha_i \ge 0} \min_{w} \mathcal{L}(w,\alpha,\beta) \le \min_{w} \max_{\alpha,\beta:\alpha_i \ge 0} \mathcal{L}(w,\alpha,\beta) = p^*$
 - Proof

$$\min_{w} \mathcal{L}(w, \alpha, \beta) \leq \mathcal{L}(w, \alpha, \beta), \forall w, \alpha \geq 0, \beta$$

$$\Rightarrow \max_{\alpha, \beta: \alpha \geq 0} \min_{w} \mathcal{L}(w, \alpha, \beta) \leq \max_{\alpha, \beta: \alpha \geq 0} \mathcal{L}(w, \alpha, \beta), \forall w$$

$$\Rightarrow \max_{\alpha, \beta: \alpha \geq 0} \min_{w} \mathcal{L}(w, \alpha, \beta) \leq \min_{w} \max_{\alpha, \beta: \alpha \geq 0} \mathcal{L}(w, \alpha, \beta)$$

• But under certain condition $d^* = p^*$

Karush-Kuhn-Tucker (KKT) Conditions

- If f and g'_i 's are convex and h'_i 's are affine, and suppose g'_i 's are all strictly feasible
- then there must exist w^* , α^* , β^*
 - w^{*} is the solution of the primal problem
 - α^* , β^* are the solutions of the dual problem
 - and the values of the two problems are equal $p^* = d^* = \mathcal{L}(w^*, \alpha^*, \beta^*)$
- And w^* , α^* , β^* satisfy the KKT conditions

$$\frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \ i = 1, \dots, n$$
$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \ i = 1, \dots, l$$

KKT dual

complementarity $\longrightarrow \alpha_i^* g_i(w^*) = 0, \ i = 1, \dots, k$ condition $g_i(w^*) \le 0, \ i = 1, \dots, k$

 $\alpha^* \ge 0, \ i = 1, \dots, k$

- Moreover, if some w^* , α^* , β^* satisfy the KKT conditions, then it is also a solution to the primal and dual problems.
- More details please refer to Boyd "Convex optimization" 2004.

Now Back to SVM Problem

Objective of an SVM

• SVM objective: finding the optimal margin classifier

$$\min_{\substack{w,b \ x,b \ x,b$$

• Re-wright the constraints as

$$g_i(w) = -y^{(i)}(w^{\top}x^{(i)} + b) + 1 \le 0$$

so as to match the standard optimization form

$$\min_{w} f(w)$$
s.t. $g_i(w) \le 0, \quad i = 1, \dots, k$
 $h_i(w) = 0, \quad i = 1, \dots, l$

Equality Cases



Objective of an SVM

• SVM objective: finding the optimal margin classifier

$$\min_{\substack{w,b \ w,b}} \frac{1}{2} \|w\|^2$$

s.t. $-y^{(i)}(w^{\top}x^{(i)}+b) + 1 \le 0, \ i = 1, \dots, m$

• Lagrangian

$$\mathcal{L}(w,b,\alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)}(w^\top x^{(i)} + b) - 1]$$

• No β or equality constraints in SVM problem

Solving

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)}(w^\top x^{(i)} + b) - 1]$$

• Derivatives

$$\frac{\partial}{\partial w}\mathcal{L}(w,b,\alpha) = w - \sum_{i=1}^{m} \alpha_i y^{(i)} x^{(i)} = 0 \quad \Rightarrow \quad w = \sum_{i=1}^{m} \alpha_i y^{(i)} x^{(i)}$$
$$\frac{\partial}{\partial b}\mathcal{L}(w,b,\alpha) = \sum_{i=1}^{m} \alpha_i y^{(i)} = 0$$

• Then Lagrangian is re-written as

$$\mathcal{L}(w,b,\alpha) = \frac{1}{2} \left\| \sum_{i=1}^{m} \alpha_i y^{(i)} x^{(i)} \right\|^2 - \sum_{i=1}^{m} \alpha_i [y^{(i)} (w^\top x^{(i)} + b) - 1]$$
$$= \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)} {}^\top x^{(j)} \left| - b \sum_{i=1}^{m} \alpha_i y^{(i)} \right| = 0$$

Solving α^*

• Dual problem

$$\max_{\alpha \ge 0} \theta_{\mathcal{D}}(\alpha) = \max_{\alpha \ge 0} \min_{w,b} \mathcal{L}(w, b, \alpha)$$
$$\max_{\alpha} \quad W(\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)^\top} x^{(j)}$$
s.t. $\alpha_i \ge 0, \ i = 1, \dots, m$
$$\sum_{i=1}^{m} \alpha_i y^{(i)} = 0$$

- To solve α^* with some methods e.g. SMO
 - We will get back to this solution later

Solving w^* and b^*

• With α^* solved, w^* is obtained by

$$w = \sum_{i=1}^{m} \alpha_i y^{(i)} x^{(i)}$$

- Only supporting vectors with α > 0
- With w^{*} solved, b^{*} is obtained by

$$b^* = -\frac{\max_{i:y^{(i)}=-1} w^{*\top} x^{(i)} + \min_{i:y^{(i)}=1} w^{*\top} x^{(i)}}{2}$$

Predicting Values

• With the solutions of *w*^{*}and *b*^{*}, the predicting value (i.e. functional margin) of each instance is

$$w^{*\top}x + b^* = \left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}\right)^\top x + b^*$$
$$= \sum_{i=1}^m \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b^*$$

• We only need to calculate the inner product of *x* with the supporting vectors



- The derivation of the SVM as presented so far assumes that the data is linearly separable.
- More practical cases are linearly non-separable.





Dealing with Non-Separable Cases

• Add slack variables

$$\min_{w,b} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \quad \leftarrow \quad \text{L1 regularization}$$

s.t.
$$y^{(i)}(w^\top x^{(i)} + b) \ge 1 - \xi_i, \ i = 1, \dots, m$$

$$\xi_i \ge 0, \ i = 1, \dots, m$$

• Lagrangian

$$\mathcal{L}(w,b,\xi,\alpha,r) = \frac{1}{2}w^{\top}w + C\sum_{i=1}^{m}\xi_i - \sum_{i=1}^{m}\alpha_i[y^{(i)}(x^{\top}w+b) - 1 + \xi_i] - \sum_{i=1}^{m}r_i\xi_i$$

Dual problem

SVM Hinge Loss vs. LR Loss

SVM Hinge loss

• LR log loss

 $\frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \max(0, 1 - y_i(w^\top x_i + b)) - y_i \log \sigma(w^\top x_i + b) - (1 - y_i) \log(1 - \sigma(w^\top x_i + b))$



Now Back to Solve α^*

• Dual problem

$$\max_{\alpha} \quad W(\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)^{\top}} x^{(j)}$$

s.t. $0 \le \alpha_i \le C, \ i = 1, \dots, m$
$$\sum_{i=1}^{m} \alpha_i y^{(i)} = 0$$

• With α^* solved, w and b are solved easily

Coordinate Ascent (Descent)

• For the optimization problem

 $\max_{\alpha} W(\alpha_1, \alpha_2, \ldots, \alpha_m)$

Coordinate ascent algorithm

```
Loop until convergence: {

For i = 1, ..., m {

\alpha_i := \arg \max_{\hat{\alpha}_i} W(\alpha_1, ..., \alpha_{i-1}, \hat{\alpha}_i, \alpha_{i+1}, ..., \alpha_m)

}
```

Coordinate Ascent (Descent)

A two-dimensional coordinate ascent example

- SMO: sequential minimal optimization
- SVM optimization problem

$$\begin{aligned} \max_{\alpha} \quad W(\alpha) &= \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)^{\top}} x^{(j)} \\ \text{s.t.} \quad 0 \leq \alpha_i \leq C, \ i = 1, \dots, m \\ \sum_{i=1}^{m} \alpha_i y^{(i)} &= 0 \end{aligned}$$

• Cannot directly apply coordinate ascent algorithm because

$$\sum_{i=1}^m lpha_i y^{(i)} = 0 \; \Rightarrow \; lpha_i y^{(i)} = - \sum_{j
eq i} lpha_j y^{(j)}$$

}

• Update two variable each time

Loop until convergence {

- 1. Select some pair α_i and α_i to update next
- 2. Re-optimize $W(\alpha)$ w.r.t. α_i and α_i

- Convergence test: whether the change of $W(\alpha)$ is smaller than a predefined value (e.g. 0.01)
- Key advantage of SMO algorithm is the update of α_i and α_i (step 2) is efficient

$$\max_{\alpha} \quad W(\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)^{\top}} x^{(j)}$$
s.t. $0 \le \alpha_i \le C, \ i = 1, \dots, m$

$$\sum_{i=1}^{m} \alpha_i y^{(i)} = 0$$

• Without loss of generality, hold $\alpha_3 \dots \alpha_m$ and optimize $W(\alpha)$ w.r.t. α_1 and α_2 $\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = -\sum_{i=3}^m \alpha_i y^{(i)} = \zeta$ $\Rightarrow \alpha_2 = -\frac{y^{(1)}}{y^{(2)}} \alpha_1 + \frac{\zeta}{y^{(2)}}$ $\alpha_1 = (\zeta - \alpha_2 y^{(2)}) y^{(1)}$

• With $\alpha_1 = (\zeta - \alpha_2 y^{(2)}) y^{(1)}$, the objective is written as

$$W(\alpha_1, \alpha_2, \dots, \alpha_m) = W((\zeta - \alpha_2 y^{(2)}) y^{(1)}, \alpha_2, \dots, \alpha_m)$$

• Thus the original optimization problem

$$\begin{aligned} \max_{\alpha} \quad W(\alpha) &= \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)^{\top}} x^{(j)} \\ \text{s.t.} \quad 0 \leq \alpha_i \leq C, \ i = 1, \dots, m \\ \sum_{i=1}^{m} \alpha_i y^{(i)} &= 0 \end{aligned}$$

is transformed into a quadratic optimization problem w.r.t. α_2

$$\max_{\alpha_2} \quad W(\alpha_2) = a\alpha_2^2 + b\alpha_2 + c$$

s.t. $0 \le \alpha_2 \le C$

• Optimizing a quadratic function is much efficient

$$\max_{\alpha_2} W(\alpha_2) = a\alpha_2^2 + b\alpha_2 + c$$

s.t. $0 \le \alpha_2 \le C$

Kernel Methods

• More practical cases are linearly non-separable.

Linearly separable case

May be solved by slack variables

Cannot be solved by slack variables

- More practical cases are linearly non-separable.
- Solution: mapping feature vectors to a higher-dimensional space $\phi(x)$
- An example

• More generally, mapping feature vectors to a different space

Feature Mapping Functions

• SVM only cares about the inner products

$$W(\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)^{\top}} x^{(j)}$$

• With the feature mapping function $\phi(x)$

$$W(\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} y^{(i)} y^{(j)} \alpha_i \alpha_j K(x^{(i)}, x^{(j)})$$

• Kernel
$$K(x^{(i)}, x^{(j)}) = \phi(x^{(i)})^{\top} \phi(x^{(j)})$$

Kernel

• With the example feature mapping function

$$\phi(x) = \begin{bmatrix} x \\ x^2 \\ x^3 \end{bmatrix}$$

• The corresponding kernel is

$$\begin{split} K(x^{(i)}, x^{(j)}) &= \phi(x^{(i)})^\top \phi(x^{(j)}) \\ &= x^{(i)} x^{(j)} + x^{(i)^2} x^{(j)^2} + x^{(i)^3} x^{(j)^3} \end{split}$$

[Kernel Trick]

- For lots of cases, we only need $K(x^{(i)}, x^{(j)})$, thus we can directly define $K(x^{(i)}, x^{(j)})$ without explicitly defining $\phi(x^{(i)})$
 - For example, suppose $x^{(i)}, x^{(j)} \in \mathbb{R}^n$

$$K(x^{(i)}, x^{(j)}) = (x^{(i)^{\top}} x^{(j)})^2$$

Kernel Example

• For example, suppose $x^{(i)}, x^{(j)} \in \mathbb{R}^n$

If n = 3, the mapping function is

$$\begin{split} K(x^{(i)}, x^{(j)}) &= (x^{(i)^{\top}} x^{(j)})^2 \\ &= \left(\sum_{k=1}^n x_k^{(i)} x_k^{(j)}\right) \left(\sum_{l=1}^n x_l^{(i)} x_l^{(j)}\right) \\ &= \sum_{k=1}^n \sum_{l=1}^n x_k^{(i)} x_k^{(j)} x_l^{(i)} x_l^{(j)} \qquad \Rightarrow \qquad \phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix} \end{split}$$

• Note that calculating $\phi(x)$ takes $O(n^2)$ time, while calculating $K(x^{(i)}, x^{(j)})$ only takes O(n) time

Kernel for Measuring Similarity

• Intuitively, for two instances x and z, if $\phi(x)$ and $\phi(z)$ are close together, then we expect

$$K(x,z) = \phi(x)^{\top} \phi(z)$$

to be large, and vice versa.

• Gaussian kernel (a very widely used kernel)

$$K(x,z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$$

- Also called radial basis function (RBF) kernel
- Then what is the feature mapping function for this kernel?

Kernel Matrix

- Consider a finite set of instances $\{x^{(1)}, \ldots, x^{(m)}\}$
- The corresponding Kernel Matrix K is defined as $\{K_{ij}\}_{i,j=1,...,m}$
- The kernel matrix K must be symmetric since

 $K_{ij} = K(x^{(i)}, x^{(j)}) = \phi(x^{(i)})^{\top} \phi(x^{(j)}) = \phi(x^{(j)})^{\top} \phi(x^{(i)}) = K(x^{(j)}, x^{(i)}) = K_{ji}$

• If we define $\phi_k(x)$ as the *k*-th coordinate of the vector $\phi(x)$, then for any vector $z \in \mathbb{R}^m$, we have

$$z^{\top}Kz = \sum_{i} \sum_{j} z_{i}K_{ij}z_{j}$$

= $\sum_{i} \sum_{j} z_{i}\phi(x^{(i)})^{\top}\phi(x^{(j)})z_{j} = \sum_{i} \sum_{j} z_{i} \sum_{k} \phi_{k}(x^{(i)})\phi_{k}(x^{(j)})z_{j}$
= $\sum_{k} \sum_{i} \sum_{j} z_{i}\phi_{k}(x^{(i)})\phi_{k}(x^{(j)})z_{j} = \sum_{k} \left(\sum_{i} z_{i}\phi_{k}(x^{(i)})\right)^{2} \ge 0$

• Therefore, K is semi-definite

Valid (Mercer) Kernel

James Mercer UK Mathematician 1883-1932

• Theorem (Mercer)

Let $K : \mathbb{R}^n \times \mathbb{R}^n \mapsto \mathbb{R}$ be given. Then for K to be a valid (Mercer) kernel, it is necessary and sufficient that for any $\{x^{(1)}, \ldots, x^{(m)}\}, m < \infty$, the corresponding kernel matrix is symmetric positive semi-definite.

- Example valid kernels
 - **RBF kernel** $K(x,z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$
 - Simple polynomial kernel $K(x,z) = (x^{\top}z)^d$
 - Cosine similarity kernel $K(x,z) = \frac{x^{\top}z}{\|x\| \cdot \|z\|}$

Sigmoid Kernel

$$K(x,z) = \tanh(\alpha x^{\top} z + c)$$

$$\tanh(b) = \frac{1 - e^{-2b}}{1 + e^{-2b}}$$

- Neural networks use sigmoid as activation function
- SVM with a sigmoid kernel is equivalent to a 2-layer perceptron

(We shall return to this after the study of neural networks)

Generalized Linear Models

Review: Linear Regression

$$\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}^{(1)} \\ \boldsymbol{x}^{(2)} \\ \vdots \\ \boldsymbol{x}^{(n)} \end{bmatrix} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} & \dots & x_d^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1^{(n)} & x_2^{(n)} & x_3^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{bmatrix} \quad \boldsymbol{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ \theta_d \end{bmatrix}$$

• Prediction $\hat{\boldsymbol{y}} = \boldsymbol{X}\boldsymbol{\theta} = \begin{bmatrix} \boldsymbol{x}^{(1)}\boldsymbol{\theta} \\ \boldsymbol{x}^{(2)}\boldsymbol{\theta} \\ \vdots \\ \boldsymbol{x}^{(n)}\boldsymbol{\theta} \end{bmatrix}$

• Objective
$$J(\boldsymbol{\theta}) = \frac{1}{2}(\boldsymbol{y} - \hat{\boldsymbol{y}})^{\top}(\boldsymbol{y} - \hat{\boldsymbol{y}}) = \frac{1}{2}(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta})^{\top}(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta})$$

Review: Matrix Form of Linear Reg.

• Objective

$$J(\boldsymbol{\theta}) = \frac{1}{2} (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta})^{\top} (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}) \quad \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

Gradient

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -\boldsymbol{X}^{\top}(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta})$$

• Solution

$$egin{aligned} rac{\partial J(oldsymbol{ heta})}{\partialoldsymbol{ heta}} &= oldsymbol{0} & o & oldsymbol{X}^ op(oldsymbol{y} - oldsymbol{X}oldsymbol{ heta}) = oldsymbol{0} \ & o & oldsymbol{X}^ opoldsymbol{y} = oldsymbol{X}^ opoldsymbol{X} + oldsymbol{x} oldsymbol{ heta} = oldsymbol{X}^ opoldsymbol{X} + oldsymbol{X} oldsymbol{ heta} = oldsymbol{0} \ & o & oldsymbol{ heta} = (oldsymbol{X}^ opoldsymbol{X})^{-1}oldsymbol{X}^ opoldsymbol{y} = oldsymbol{0} \ & o & oldsymbol{ heta} = (oldsymbol{X}^ opoldsymbol{X})^{-1}oldsymbol{X}^ opoldsymbol{y} = oldsymbol{0} \ & o & oldsymbol{ heta} = (oldsymbol{X}^ opoldsymbol{X})^{-1}oldsymbol{X}^ opoldsymbol{y}$$

Generalized Linear Models

• Dependence

$$y = f(\theta^{\top}\phi(x))$$

- Feature mapping function $\phi(x): \mathbb{R}^d \mapsto \mathbb{R}^h$
- Mapped feature matrix $\Phi_{n imes h}$

$$\Phi = \begin{bmatrix} \phi(x^{(1)}) \\ \phi(x^{(2)}) \\ \vdots \\ \phi(x^{(i)}) \\ \vdots \\ \phi(x^{(n)}) \end{bmatrix} = \begin{bmatrix} \phi_1(x^{(1)}) & \phi_2(x^{(1)}) & \cdots & \phi_h(x^{(1)}) \\ \phi_1(x^{(2)}) & \phi_2(x^{(2)}) & \cdots & \phi_h(x^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(x^{(i)}) & \phi_2(x^{(i)}) & \cdots & \phi_h(x^{(i)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(x^{(n)}) & \phi_2(x^{(n)}) & \cdots & \phi_h(x^{(n)}) \end{bmatrix}$$

Matrix Form of Kernel Linear Regression

• Objective

$$J(\boldsymbol{\theta}) = \frac{1}{2} (\boldsymbol{y} - \boldsymbol{\Phi} \boldsymbol{\theta})^{\top} (\boldsymbol{y} - \boldsymbol{\Phi} \boldsymbol{\theta}) \quad \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

Gradient

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -\boldsymbol{\Phi}^\top (\boldsymbol{y} - \boldsymbol{\Phi} \boldsymbol{\theta})$$

• Solution

$$egin{aligned} rac{\partial J(m{ heta})}{\partialm{ heta}} = m{0} & o & m{\Phi}^ op(m{y} - m{\Phi}m{ heta}) = m{0} \ & o & m{\Phi}^ opm{y} = m{\Phi}^ opm{\Phi}m{ heta} \ & o & m{ heta} = (m{\Phi}^ opm{\Phi})^{-1}m{\Phi}^ opm{y} \end{aligned}$$

Matrix Form of Kernel Linear Regression

• With the Algebra trick

 $(P^{-1} + B^{\top}R^{-1}B)^{-1}B^{\top}R^{-1} = PB^{\top}(BPB^{\top} + R)^{-1}$

• The optimal parameters with L2 regularization

$$\hat{oldsymbol{ heta}} = (oldsymbol{\Phi}^{ op} oldsymbol{\Phi} + \lambda oldsymbol{I}_h)^{-1} oldsymbol{\Phi}^{ op} oldsymbol{y} \ = oldsymbol{\Phi}^{ op} (oldsymbol{\Phi} oldsymbol{\Phi}^{ op} + \lambda oldsymbol{I}_n)^{-1} oldsymbol{y}$$

for prediction, we never actually need access ${f \Phi}$

$$\hat{oldsymbol{y}} = oldsymbol{\Phi} \hat{oldsymbol{ heta}} = oldsymbol{\Phi} oldsymbol{\Phi}^ op (oldsymbol{\Phi} oldsymbol{\Phi}^ op + \lambda oldsymbol{I}_n)^{-1} oldsymbol{y}$$

 $= oldsymbol{K} (oldsymbol{K} + \lambda oldsymbol{I}_n)^{-1} oldsymbol{y}$

where the kernel matrix $\boldsymbol{K} = \{K(\boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)})\}$

[http://www.ics.uci.edu/~welling/classnotes/papers_class/Kernel-Ridge.pdf]