

Linear Models

Li, Jia

DSAA 5002

The Hong Kong of Science and Technology (Guangzhou)

2025 Fall

Sep 15

Overview

- Instance-based classifiers, k-NN
- Perceptron Learning
- Linear Models

Instance-Based Classifiers

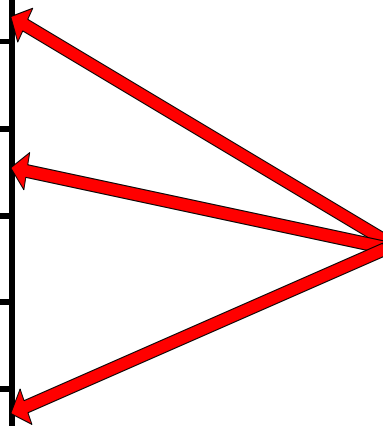
Set of Stored Cases

Atr1	AtrN	Class
			A
			B
			B
			C
			A
			C
			B

- Store the training records
- Use training records to predict the class label of unseen cases

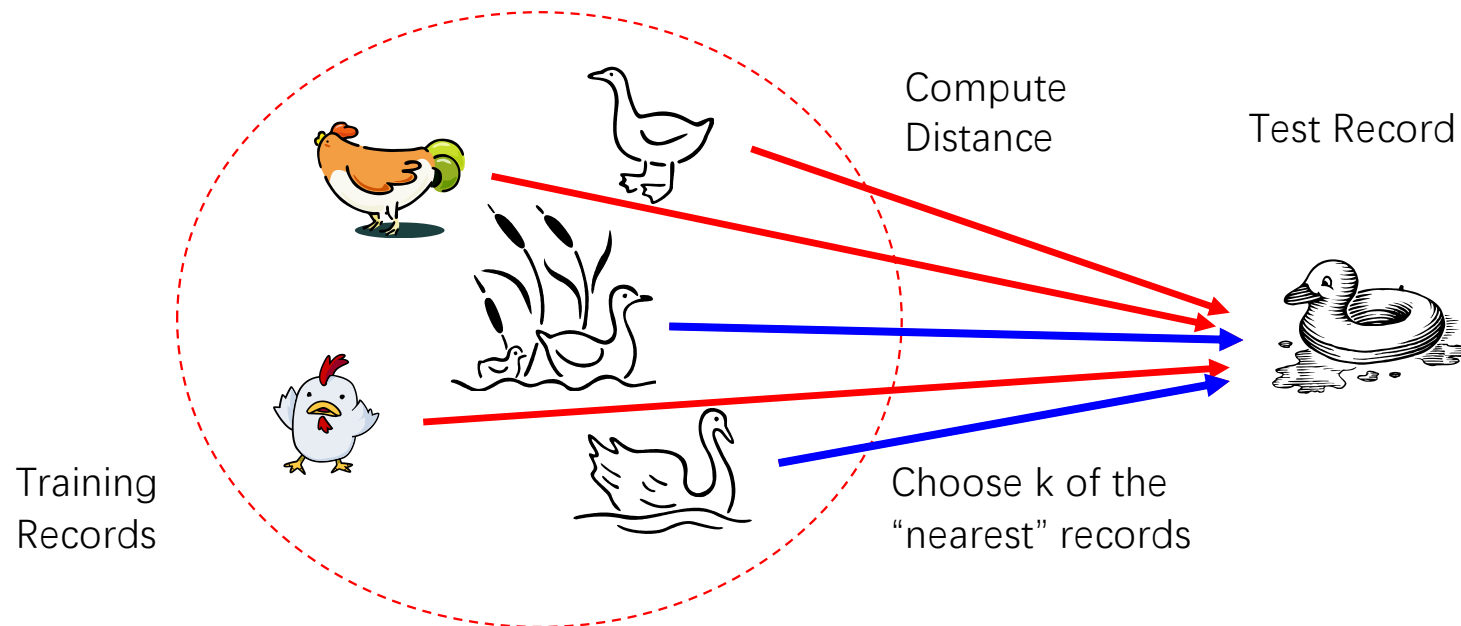
Unseen Case

Atr1	AtrN

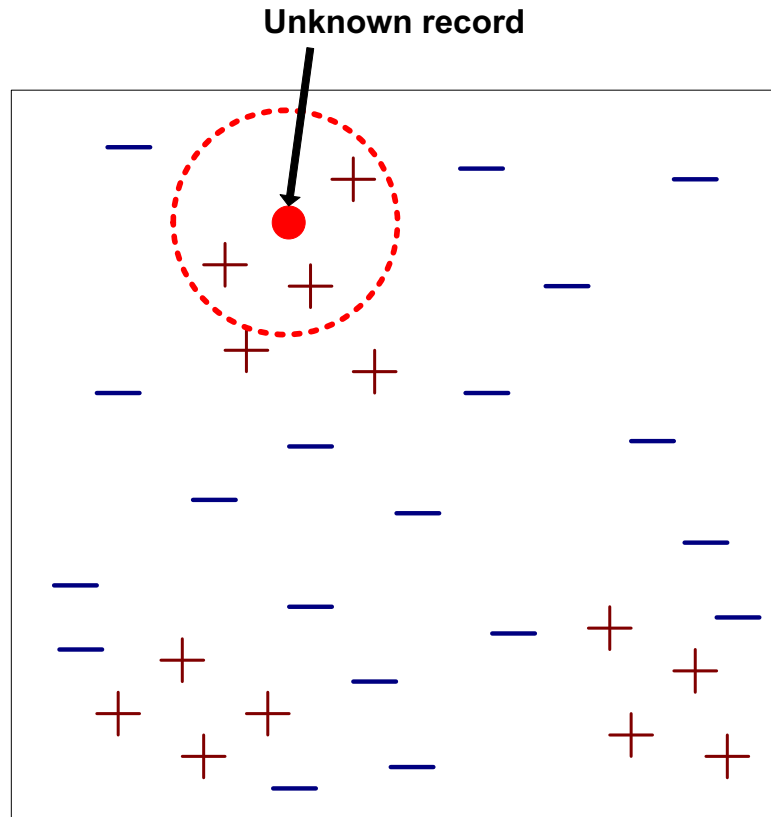


Nearest Neighbor Classifiers

- Basic idea:
 - If it walks like a duck, quacks like a duck, then it's probably a duck



Nearest-Neighbor Classifiers



- Three things
 - The set of stored records
 - Distance Metric to compute distance between records
 - The value of k , the number of nearest neighbors to retrieve
- To classify an unknown record:
 - Compute distance to other training records
 - Identify k nearest neighbors
 - Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

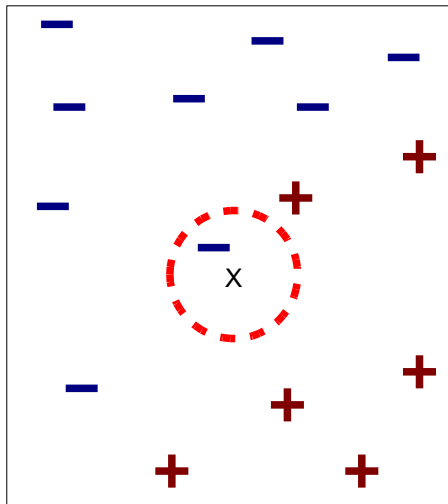
K Nearest Neighbor Classification

- Compute distance between two points:
 - Euclidean distance

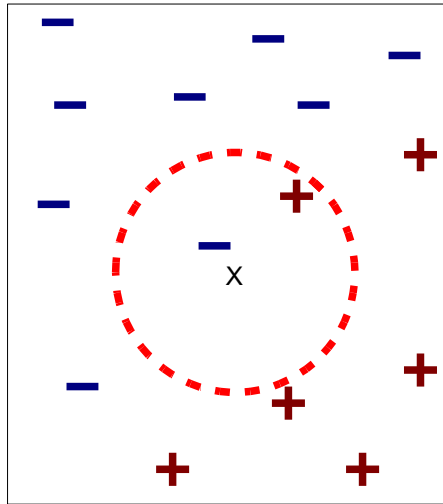
$$d(p, q) = \sqrt{\sum_i (p_i - q_i)^2}$$

- Determine the class from nearest neighbor list
 - take the majority vote of class labels among the k-nearest neighbors
 - Weigh the vote according to distance
 - weight factor, $w = 1/d^2$

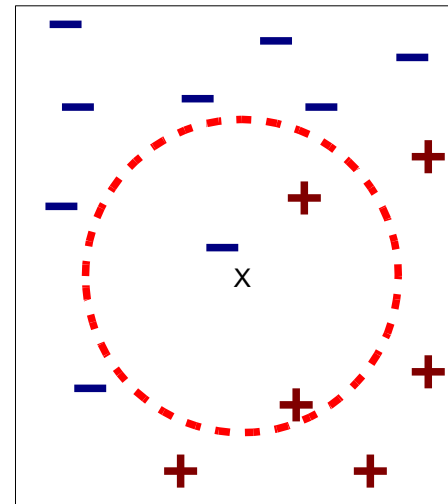
Definition of Nearest Neighbor



(a) 1-nearest neighbor



(b) 2-nearest neighbor

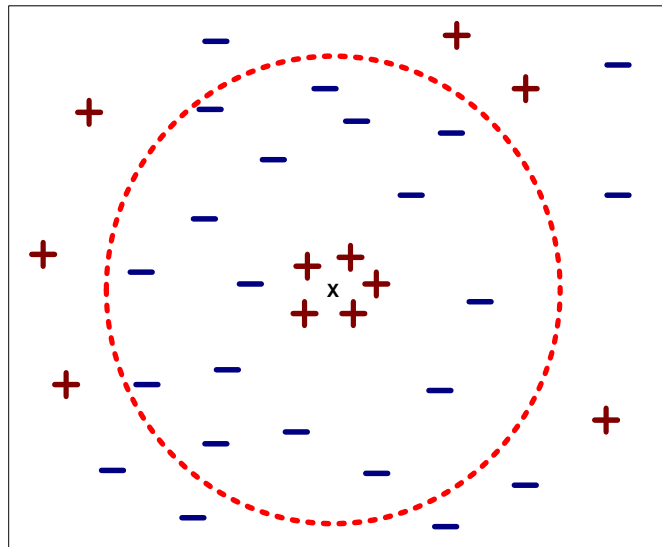


(c) 3-nearest neighbor

K-nearest neighbors of a record x are data points that have the k smallest distance to x

Nearest Neighbor Classification

- Choosing the value of k :
 - If k is too small, sensitive to noise points
 - If k is too large, neighborhood may include points from other classes



The KNN Algorithm

Algorithm 6.2 The k -nearest neighbor classifier.

- 1: Let k be the number of nearest neighbors and D be the set of training examples.
 - 2: **for** each test instance $z = (\mathbf{x}', y')$ **do**
 - 3: Compute $d(\mathbf{x}', \mathbf{x})$, the distance between z and every example, $(\mathbf{x}, y) \in D$.
 - 4: Select $D_z \subseteq D$, the set of k closest training examples to z .
 - 5: $y' = \underset{v}{\operatorname{argmax}} \sum_{(\mathbf{x}_i, y_i) \in D_z} I(v = y_i)$
 - 6: **end for**
-

- Here, v is a class label, y_i is the class label for one of the nearest neighbors, and $I(\cdot)$ is an indicator function that returns 1 if its argument is true and 0 otherwise.

Nearest Neighbor Classification

- Scaling issues
 - Attributes may have to be scaled to prevent distance measures from being dominated by one of the attributes
 - Example:
 - height of a person may vary from 1.5m to 1.8m
 - weight of a person may vary from 90lb to 300lb
 - income of a person may vary from \$10K to \$1M

The Impacts of Distance functions

- $D_{\text{sum}}(A,B) = D_{\text{gender}}(A,B) + D_{\text{age}}(A,B) + D_{\text{salary}}(A,B)$
- $D_{\text{norm}}(A,B) = D_{\text{sum}}(A,B)/\max(D_{\text{sum}})$
- $D_{\text{euclid}}(A,B) = \sqrt{D_{\text{gender}}(A,B)^2 + D_{\text{age}}(A,B)^2 + D_{\text{salary}}(A,B)^2}$

Table 9.9 New Customer

<i>Recnum</i>	<i>Gender</i>	<i>Age</i>	<i>Salary</i>
new	female	45	\$100,000

Table 9.10 Set of Nearest Neighbors for New Customer

	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>Neighbors</i>
d_{sum}	1.662	1.659	1.338	1.003	1.640	4,3,5,2,1
d_{norm}	0.554	0.553	0.446	0.334	0.547	4,3,5,2,1
d_{euclid}	0.781	1.052	1.251	0.494	1.000	4,1,5,2,3

Table 9.11 Customers with Attrition History

<i>Recnum</i>	<i>Gender</i>	<i>Age</i>	<i>Salary</i>	<i>Attriter</i>
1	female	27	\$19,000	no
2	male	51	\$64,000	yes
3	male	52	\$105,000	yes
4	female	33	\$55,000	yes
5	male	45	\$45,000	no
new	female	45	\$100,000	?

Table 9.12 Using MBR to Determine if the New Customer Will Attrite

	<i>Neighbors</i>	<i>Neighbor Attrition</i>	<i>k = 1</i>	<i>k = 2</i>	<i>k = 3</i>	<i>k = 4</i>	<i>k = 5</i>
d_{sum}	4,3,5,2,1	Y,Y,N,Y,N	yes	yes	yes	yes	yes
d_{Euclid}	4,1,5,2,3	Y,N,N,Y,Y	yes	?	no	?	yes

Table 9.13 Attrition Prediction with Confidence

	<i>k = 1</i>	<i>k = 2</i>	<i>k = 3</i>	<i>k = 4</i>	<i>k = 5</i>
d_{sum}	yes, 100%	yes, 100%	yes, 67%	yes, 75%	yes, 60%
d_{Euclid}	yes, 100%	yes, 50%	no, 67%	yes, 50%	yes, 60%


Remarks on Lazy vs. Eager Learning

- Instance-based learning: lazy evaluation
- Decision-tree and Bayesian classification: eager evaluation
- Key differences
 - Lazy method may consider query instance xq when deciding how to generalize beyond the training data D
 - Eager method cannot since they have already chosen global approximation when seeing the query
- Efficiency: Lazy - less time training but more time predicting
- Accuracy
 - Lazy method effectively uses a richer hypothesis space since it uses many local linear functions to form its implicit global approximation to the target function
 - Eager: must commit to a single hypothesis that covers the entire instance space

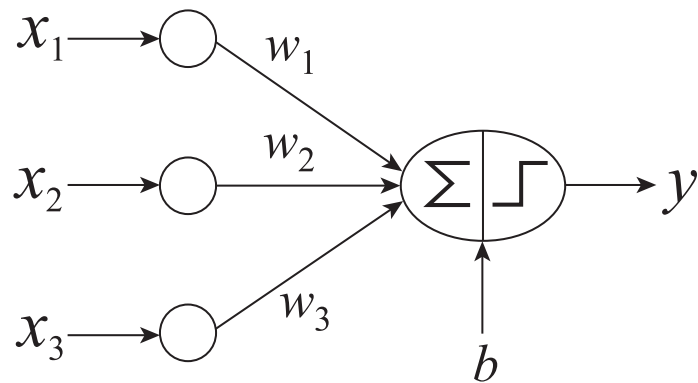
Artificial Neural Networks (ANN)

- **Basic Idea:** A complex non-linear function can be learned as a composition of simple processing units.
- ANN is a collection of simple processing units (nodes) that are connected by directed links (edges)
 - Every node receives signals from incoming edges, performs computations, and transmits signals to outgoing edges
 - Analogous to *human brain* where nodes are neurons and signals are electrical impulses
 - Weight of an edge determines the strength of connection between the nodes
- Simplest ANN: **Perceptron** (single neuron)

Why Do We Discuss Perceptron?

- We have discussed decision tree and k-NN.
 - For decision tree, we use some of its attributes (features) to build the model.
 - For k-NN, we use all its attributes/features/dimensions to classify.
 - Are there other ones?
 - Perceptron is to learn weights for features, and is to find a linear decision boundary.
- 

Basic Architecture of Perceptron



$$y = \begin{cases} 1, & \text{if } \mathbf{w}^T \mathbf{x} + b > 0. \\ -1, & \text{otherwise.} \end{cases}$$

$$\hat{y} = \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})$$

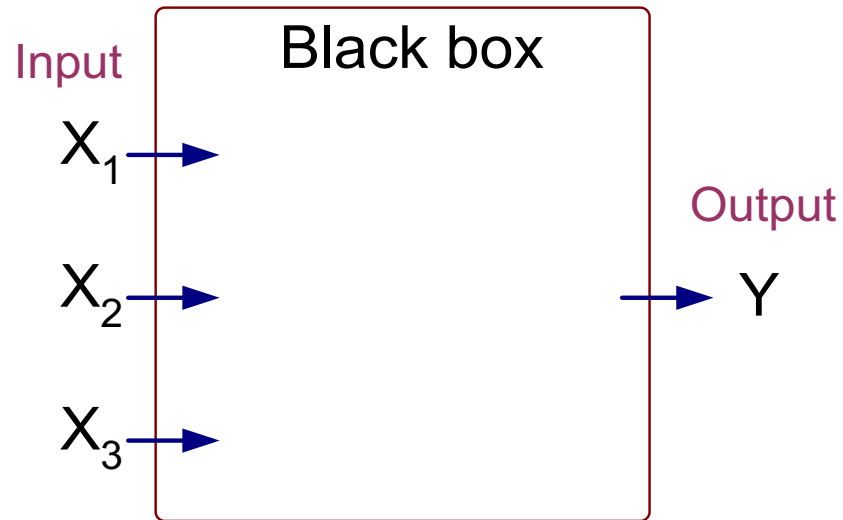
Activation Function

$$\tilde{\mathbf{w}} = (\mathbf{w}^T \ b)^T \quad \tilde{\mathbf{x}} = (\mathbf{x}^T \ 1)^T$$

- Learns linear decision boundaries
- Related to logistic regression (activation function is sign instead of sigmoid)

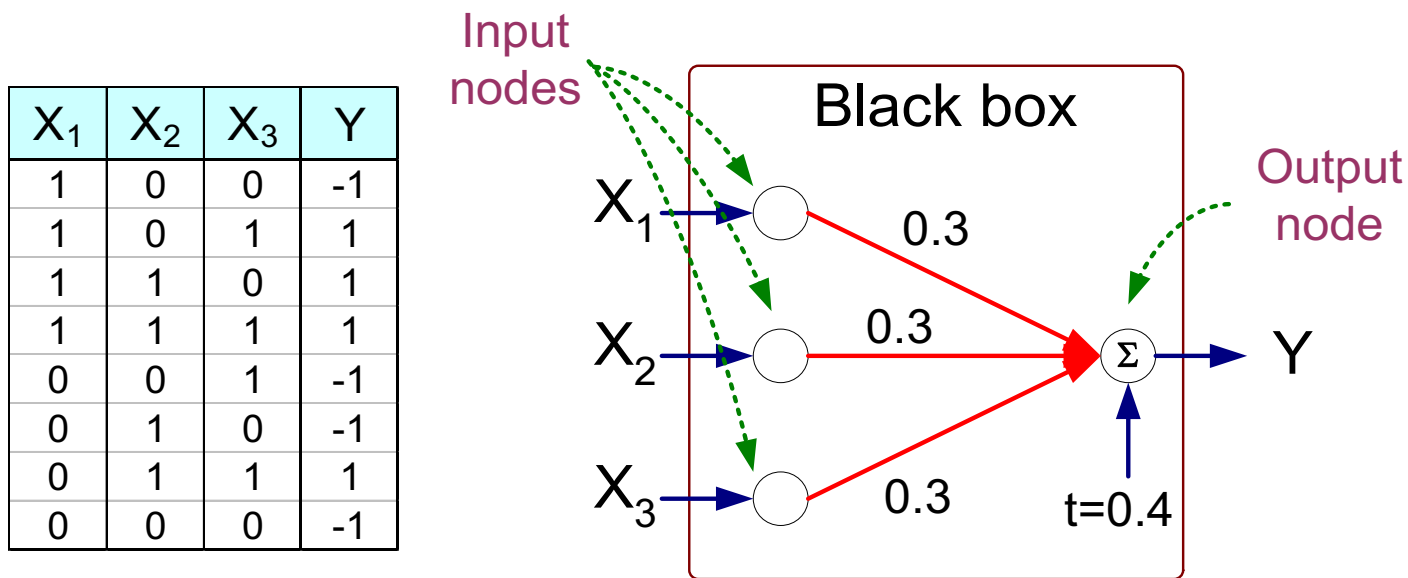
Perceptron Example

X_1	X_2	X_3	Y
1	0	0	-1
1	0	1	1
1	1	0	1
1	1	1	1
0	0	1	-1
0	1	0	-1
0	1	1	1
0	0	0	-1



Output Y is 1 if at least two of the three inputs are equal to 1.

Perceptron Example



$$Y = \text{sign}(0.3X_1 + 0.3X_2 + 0.3X_3 - 0.4)$$

$$\text{where } \text{sign}(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases}$$

Perceptron Learning Rule

- Initialize the weights (w_0, w_1, \dots, w_d)
- Repeat
 - For each training example (x_i, y_i)
 - Compute \hat{y}_i
 - Update the weights:

$$w_j^{(k+1)} = w_j^{(k)} + \lambda(y_i - \hat{y}_i^{(k)})x_{ij};$$

- Until stopping condition is met
- k: iteration number; λ : learning rate

Perceptron Learning Rule

- Weight update formula:

$$w_j^{(k+1)} = w_j^{(k)} + \lambda(y_i - \hat{y}_i^{(k)})x_{ij}$$

- Intuition:
 - Update weight based on error: $e = (y_i - \hat{y}_i)$
 - If $y = \hat{y}$, $e=0$: no update needed
 - If $y > \hat{y}$, $e=2$: weight must be increased (assuming X_{ij} is positive) so that \hat{y} will increase
 - If $y < \hat{y}$, $e=-2$: weight must be decreased (assuming X_{ij} is positive) so that \hat{y} will decrease

The Algorithm

Algorithm 6.3 Perceptron learning algorithm.

- 1: Let $D.train = \{(\tilde{\mathbf{x}}_i, y_i) \mid i = 1, 2, \dots, n\}$ be the set of training instances.
 - 2: Set $k \leftarrow 0$.
 - 3: Initialize the weight vector $\tilde{\mathbf{w}}^{(0)}$ with random values.
 - 4: **repeat**
 - 5: **for** each training instance $(\tilde{\mathbf{x}}_i, y_i) \in D.train$ **do**
 - 6: Compute the predicted output $\hat{y}_i^{(k)}$ using $\tilde{\mathbf{w}}^{(k)}$.
 - 7: **for** each weight component w_j **do**
 - 8: Update the weight, $w_j^{(k+1)} = w_j^{(k)} + \lambda(y_i - \hat{y}_i^{(k)})x_{ij}$.
 - 9: **end for**
 - 10: Update $k \leftarrow k + 1$.
 - 11: **end for**
 - 12: **until** $\sum_{i=1}^n |y_i - \hat{y}_i^{(k)}|/n$ is less than a threshold γ
-

Example of Perceptron Learning

$$\lambda = 0.1$$

X_1	X_2	X_3	Y
1	0	0	-1
1	0	1	1
1	1	0	1
1	1	1	1
0	0	1	-1
0	1	0	-1
0	1	1	1
0	0	0	-1

	w_0	w_1	w_2	w_3
0	0	0	0	0
1	-0.2	-0.2	0	0
2	0	0	0	0.2
3	0	0	0	0.2
4	0	0	0	0.2
5	-0.2	0	0	0
6	-0.2	0	0	0
7	0	0	0.2	0.2
8	-0.2	0	0.2	0.2

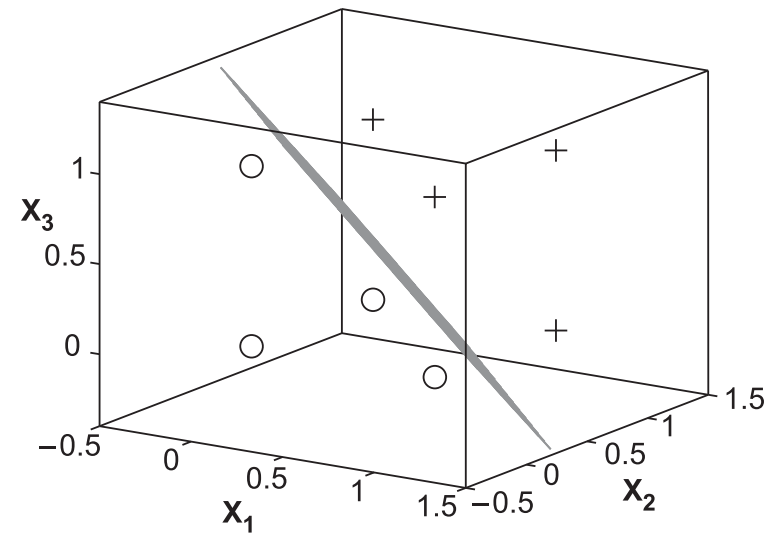
Weight updates over first epoch

Epoch	w_0	w_1	w_2	w_3
0	0	0	0	0
1	-0.2	0	0.2	0.2
2	-0.2	0	0.4	0.2
3	-0.4	0	0.4	0.2
4	-0.4	0.2	0.4	0.4
5	-0.6	0.2	0.4	0.2
6	-0.6	0.4	0.4	0.2


Weight updates over
all epochs

Perceptron Learning


- Here, y is a linear combination of input variables, decision boundary is linear.



Overview

- Linear models
 - Perceptron: **model** (linear classifier) and **learning algorithm** (updating rules) combined as one.
 - Is there a better way to learn linear models?
- To separate **models** and **learning algorithms**
 - Learning as optimization
 - Surrogate loss function
 - Regularization

model design

 - Gradient descent
 - Batch and online gradients
 - Support vector machines

optimization

Learning as Optimization

$$\min_{\mathbf{w}} \sum_n \mathbf{1}[y_n \mathbf{w}^T \mathbf{x}_n < 0]$$

↑
fewest mistakes

- The **perceptron algorithm** will find an optimal \mathbf{w} if the data is **separable**
 - efficiency depends on the **margin** and **norm** of the data
- However, if the data is not separable, optimizing this is **NP-hard**
 - i.e., there is no efficient way to minimize this unless **P=NP**

Learning as Optimization

$$\min_{\mathbf{w}} \sum_n \mathbf{1}[y_n \mathbf{w}^T \mathbf{x}_n < 0] + \lambda R(\mathbf{w})$$

↑
fewest mistakes

↓ hyperparameter

← simpler model

- In addition to minimizing **training error**, we want a **simpler model**
 - Remember our goal is to minimize **generalization error**
- We can add a **regularization** term $R(\mathbf{w})$ that prefers simpler models
 - For example we may prefer decision trees of shallow depth
- Here λ is a **hyperparameter** of optimization problem

Learning as Optimization

$$\min_{\mathbf{w}} \sum_n \mathbf{1}[y_n \mathbf{w}^T \mathbf{x}_n < 0] + \lambda R(\mathbf{w})$$

↑
fewest mistakes

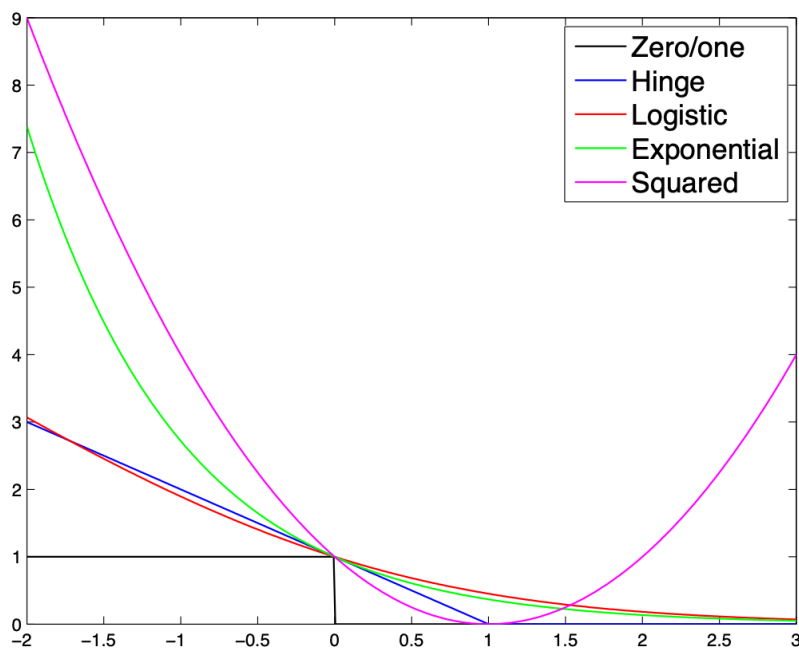
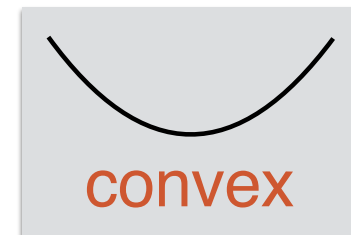
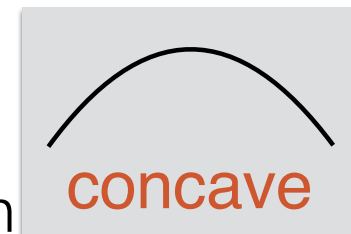
↓ hyperparameter

↙ simpler model

- The questions that remain are:
 - What are good ways to [adjust the optimization problem](#) so that there are efficient algorithms for solving it?
 - What are good [regularizations](#) $R(\mathbf{w})$ for hyperplanes?
 - Assuming that the optimization problem can be adjusted appropriately, what [algorithms](#) exist for solving the regularized optimization problem?

Convex Surrogate Loss Functions

- Zero/one loss is hard to optimize
 - Small changes in \mathbf{w} can cause large changes in the loss
- Surrogate loss: replace Zero/one loss by a smooth function
 - Easier to optimize if the surrogate loss is convex



$$y = +1 \quad \hat{y} \leftarrow \mathbf{w}^T \mathbf{x}$$

$$\text{Zero/one: } \ell^{(0/1)}(y, \hat{y}) = \mathbf{1}[y\hat{y} \leq 0]$$

$$\text{Hinge: } \ell^{(\text{hin})}(y, \hat{y}) = \max\{0, 1 - y\hat{y}\}$$

$$\text{Logistic: } \ell^{(\text{log})}(y, \hat{y}) = \frac{1}{\log 2} \log(1 + \exp[-y\hat{y}])$$

$$\text{Exponential: } \ell^{(\text{exp})}(y, \hat{y}) = \exp[-y\hat{y}]$$

$$\text{Squared: } \ell^{(\text{sqr})}(y, \hat{y}) = (y - \hat{y})^2$$

Weight Regularization

- What are good **regularization** functions $R(\mathbf{w})$ for hyperplanes?
- We would like the weights —
 - To be **small** —
 - Change in the features cause small change to the score
 - Robustness to noise
 - To be **sparse** —
 - Use as few features as possible
 - Similar to controlling the depth of a decision tree
- This is a form of **inductive bias**

Weight Regularization

- Just like the **surrogate loss function**, we would like $R(\mathbf{w})$ to be **convex**.
- **Small weights** regularization

$$R^{(\text{norm})}(\mathbf{w}) = \sqrt{\sum_d w_d^2}$$

$$R^{(\text{sqr})}(\mathbf{w}) = \sum_d w_d^2$$

- **Sparsity** regularization

$$R^{(\text{count})}(\mathbf{w}) = \sum_d \mathbf{1}[|w_d| > 0]$$

not convex

- Family of “**p-norm**” regularization

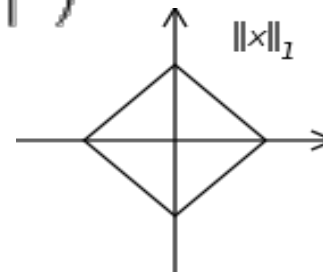
$$R^{(\text{p-norm})}(\mathbf{w}) = \left(\sum_d |w_d|^p \right)^{1/p}$$

Contours of p-norms

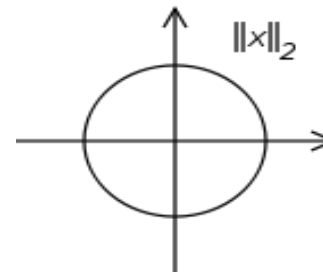
$$\|x\|_p = (|x_1|^p + |x_2|^p + \cdots + |x_n|^p)^{\frac{1}{p}}$$

convex for $p \geq 1$

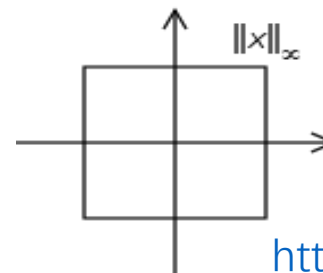
$$\|x\|_1 = \sum_{i=1}^n |x_i|$$



$$\|x\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2}$$



$$\|x\|_\infty = \max_{i=1, \dots, n} |x_i|$$

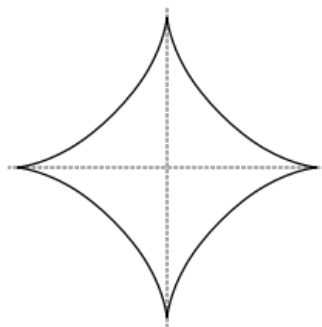


http://en.wikipedia.org/wiki/Lp_space

Contours of p-norms

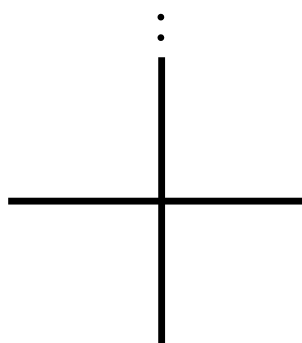
$$\|x\|_p = (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{\frac{1}{p}} \quad \text{not convex for } 0 \leq p < 1$$

$$p = \frac{2}{3}$$



Counting non-zeros:

$$p = 0$$



$$R^{(\text{count})}(\mathbf{w}) = \sum_d \mathbf{1}[|w_d| > 0]$$

http://en.wikipedia.org/wiki/Lp_space

General Optimization Framework

$$\min_{\mathbf{w}} \sum_n \ell(y_n, \mathbf{w}^T \mathbf{x}_n) + \lambda R(\mathbf{w})$$

Diagram illustrating the General Optimization Framework equation:

- The term $\ell(y_n, \mathbf{w}^T \mathbf{x}_n)$ is labeled as **surrogate loss** (indicated by an upward arrow).
- The term $\lambda R(\mathbf{w})$ is labeled as **regularization** (indicated by a diagonal arrow).
- The parameter λ is labeled as **hyperparameter** (indicated by a downward arrow).

- Select a suitable:
 - **convex surrogate loss**
 - **convex regularization**
- Select the **hyperparameter** λ
- Minimize the regularized objective with respect to \mathbf{w}
- This framework for optimization is called **Tikhonov regularization** or
- generally **Structural Risk Minimization (SRM)**

The Gradient-based Approaches

- It is like blindfolded mountain climbing.
- To find the max of a function $f(x)$, the optimizer keeps the current estimation about x , measures the gradient of the parameters of x , and take a step along the direction of the gradient as $x \leftarrow x + \eta g$, where η is the step size.

Gradient Descent

- A gradient is a multidimensional generalization of a derivative.
- Consider a function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ with a vector input $x = \langle x_1, x_2, \dots, x_D \rangle$ process a scalar value output y .
- To differentiate f with one of the input, x_i using $\frac{\partial f}{\partial x_i}$.
- The gradient of f is the vector with the derivative f with respect to every x_i .

The Gradient Descent Algorithm

Algorithm 21 GRADIENTDESCENT($\mathcal{F}, K, \eta_1, \dots$)

```
1:  $\mathbf{z}^{(0)} \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize variable we are optimizing
2: for  $k = 1 \dots K$  do
3:    $\mathbf{g}^{(k)} \leftarrow \nabla_{\mathbf{z}} \mathcal{F}|_{\mathbf{z}^{(k-1)}}$  // compute gradient at current location
4:    $\mathbf{z}^{(k)} \leftarrow \mathbf{z}^{(k-1)} - \eta^{(k)} \mathbf{g}^{(k)}$  // take a step down the gradient
5: end for
6: return  $\mathbf{z}^{(K)}$ 
```

An Example (1)

$$\min_{\mathbf{w}} \sum_n \mathbf{1}[y_n \mathbf{w}^T \mathbf{x}_n < 0] + \lambda R(\mathbf{w})$$

↑
↓

fewest mistakes
hyperparameter

←

simpler model

- Here, let the exponential loss $\ell^{(exp)}(y, \hat{y}) = \exp(-y\hat{y})$ be a loss function, and the 2-norm $R(w) = \frac{1}{2} (\sum_d |w_d|^2)^{1/2}$ be a regularizer.

$$\mathcal{L}(w, b) = \sum_n \exp [-y_n(w \cdot x_n + b)] + \frac{\lambda}{2} ||w||^2$$

An Example (2)

$$\mathcal{L}(w, b) = \sum_n \exp [-y_n(w \cdot x_n + b)] + \frac{\lambda}{2} ||w||^2$$

- The derivatives with respect to b is as follows, and update $b \leftarrow b - \eta \frac{\partial \mathcal{L}}{\partial b}$.

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial b} &= \frac{\partial}{\partial b} \sum_n \exp [-y_n(w \cdot x_n + b)] + \frac{\partial}{\partial b} \frac{\lambda}{2} ||w||^2 \\ &= \sum_n \frac{\partial}{\partial b} \exp [-y_n(w \cdot x_n + b)] + 0 \\ &= \sum_n \left(\frac{\partial}{\partial b} - y_n(w \cdot x_n + b) \right) \exp [-y_n(w \cdot x_n + b)] \\ &= - \sum_n y_n \exp [-y_n(w \cdot x_n + b)] \end{aligned}$$

An Example (3)

$$\mathcal{L}(w, b) = \sum_n \exp [-y_n(w \cdot x_n + b)] + \frac{\lambda}{2} \|w\|^2$$

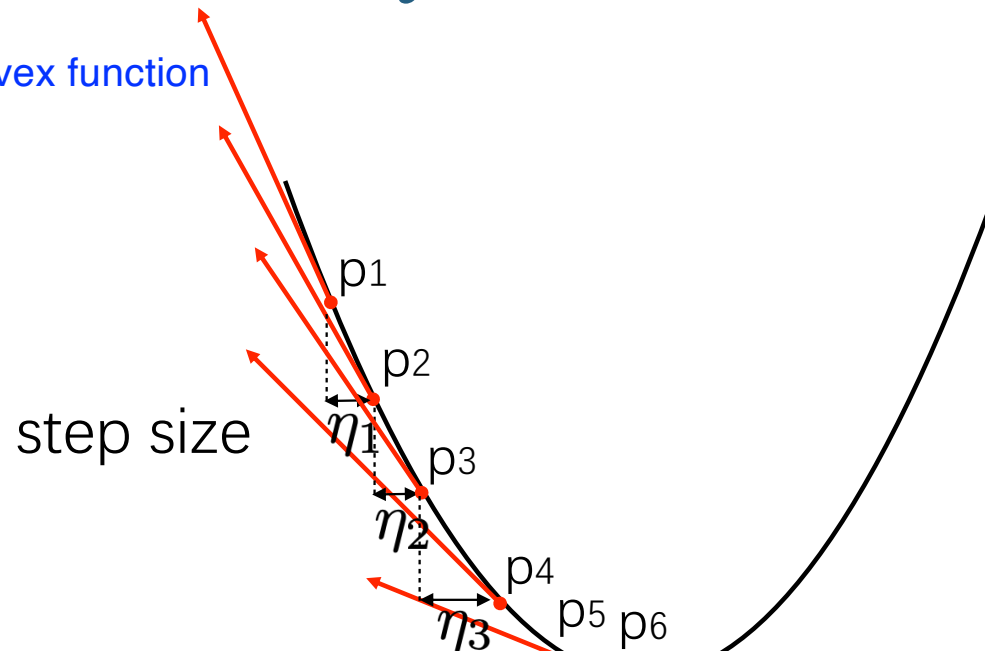
- The derivatives with respect to w is as follows, and update $w \leftarrow w - \eta \nabla_w \mathcal{L}$.

$$\begin{aligned} \nabla_w \mathcal{L} &= \nabla_w \sum_n \exp [-y_n(w \cdot x_n + b)] + \nabla_w \frac{\lambda}{2} \|w\|^2 \\ &= \sum_n (\nabla_w - y_n(w \cdot x_n + b)) \exp [-y_n(w \cdot x_n + b)] + \lambda w \\ &= - \sum_n y_n x_n \exp [-y_n(w \cdot x_n + b)] + \lambda w \end{aligned}$$

- For poorly classified points, the gradient points is in the direction $-y_n x_n$, so the update is $w \leftarrow w + c y_n x_n$ where c is a scalar value, the update for the part of the gradient related to the regularizer is $w \leftarrow (1 - \lambda)w$.

Optimization by Gradient Descent

Convex function



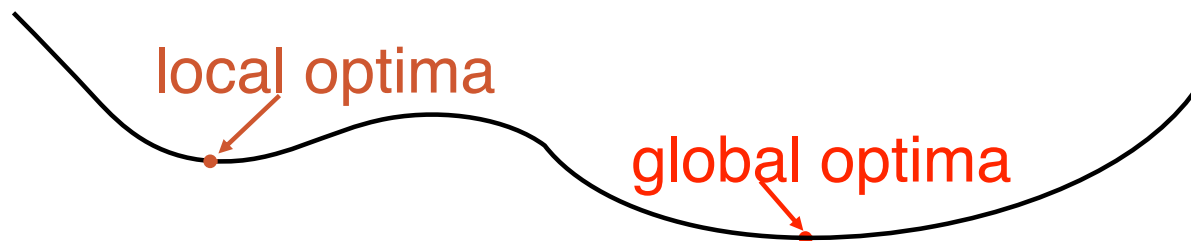
$$g^{(k)} \leftarrow \nabla_p F(p) \big|_{p_k}$$

compute gradient at the current location

$$p_{k+1} \leftarrow p_k - \eta_k g^{(k)}$$

take a step down the gradient

Non-convex function



local optima = global optima

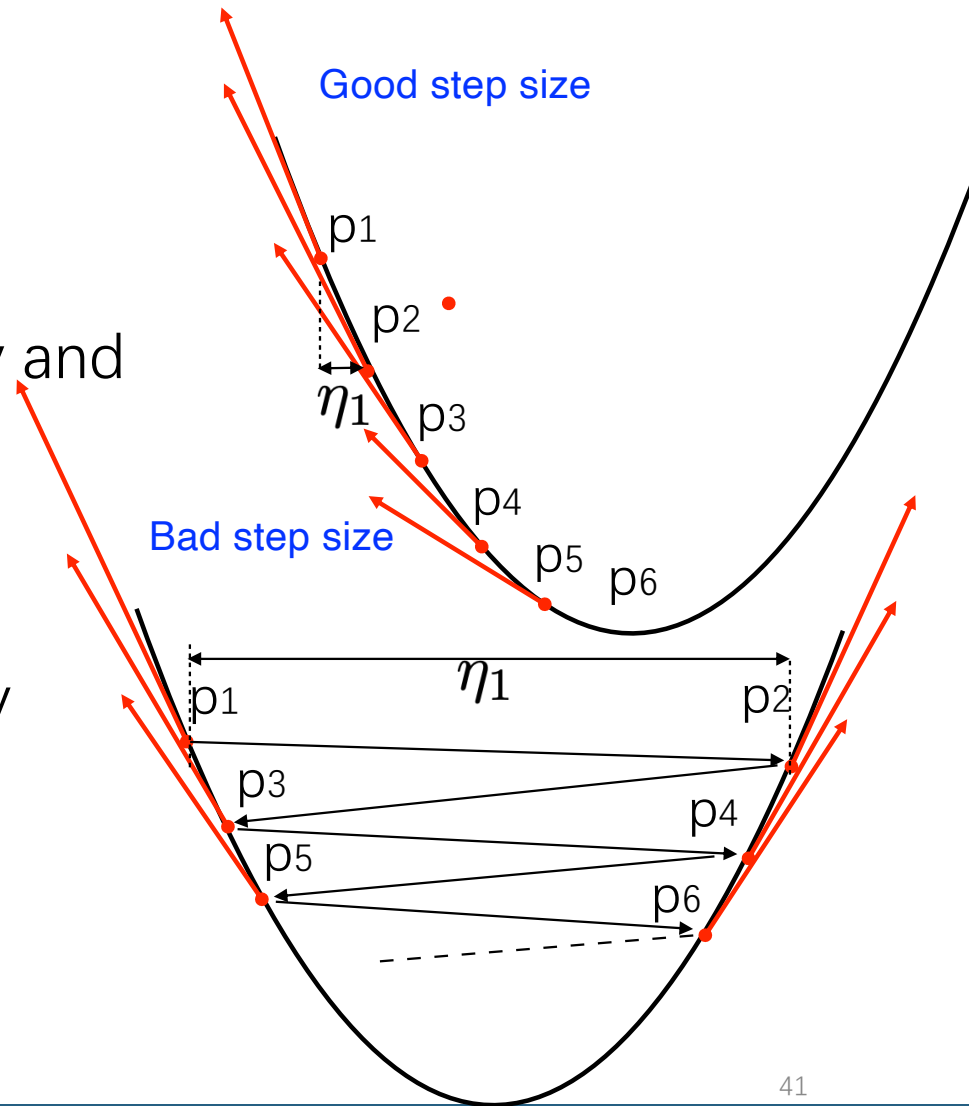
Choice of Step Size

- The step size is important —
 - **too small**: slow convergence
 - **too large**: no convergence
- A strategy is to use large step sizes initially and small step sizes later:

$$\eta_t \leftarrow \eta_0 / (t_0 + t)$$

- There are methods that converge faster by adapting step size to the **curvature** of the function
 - Field of **convex optimization**

<http://stanford.edu/~boyd/cvxbook/>



Exponential Loss

$$\mathcal{L}(\mathbf{w}) = \sum_n \exp(-y_n \mathbf{w}^T \mathbf{x}_n) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad \text{objective}$$

$$\frac{d\mathcal{L}}{d\mathbf{w}} = \sum_n -y_n \mathbf{x}_n \exp(-y_n \mathbf{w}^T \mathbf{x}_n) + \lambda \mathbf{w} \quad \text{gradient}$$

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \left(\sum_n -y_n \mathbf{x}_n \exp(-y_n \mathbf{w}^T \mathbf{x}_n) + \lambda \mathbf{w} \right) \quad \text{update}$$

loss term

$$\mathbf{w} \leftarrow \mathbf{w} + c y_n \mathbf{x}_n$$

↑
high for misclassified points

regularization term

$$\mathbf{w} \leftarrow (1 - \eta \lambda) \mathbf{w}$$

shrinks weights towards zero

Batch and Online Gradients

$$\mathcal{L}(\mathbf{w}) = \sum_n \mathcal{L}_n(\mathbf{w}) \quad \text{objective}$$

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{d\mathcal{L}}{d\mathbf{w}} \quad \text{gradient descent}$$

batch gradient

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \left(\sum_n \frac{d\mathcal{L}_n}{d\mathbf{w}} \right)$$

sum of n gradients

update weight after you see all points

online gradient

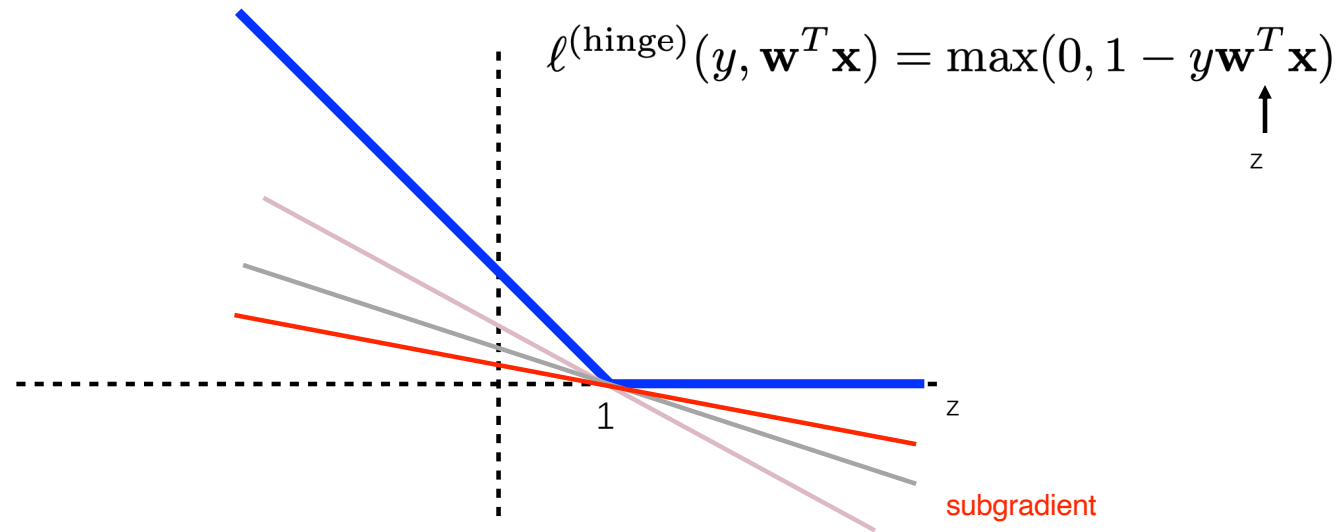
$$\mathbf{w} \leftarrow \mathbf{w} - \eta \left(\frac{d\mathcal{L}_n}{d\mathbf{w}} \right)$$

gradient at nth point

update weights after you see each point

Online gradients are the default method for multi-layer perceptrons

Subgradient



- ◆ The **hinge loss** is not differentiable at $z=1$
- ◆ **Subgradient** is any direction that is **below** the function
- ◆ For the **hinge loss** a possible **subgradient** is:

$$\frac{d\ell^{\text{hinge}}}{d\mathbf{w}} = \begin{cases} 0 & \text{if } y\mathbf{w}^T \mathbf{x} > 1 \\ -y\mathbf{x} & \text{otherwise} \end{cases}$$

Closed-form Optimization for Squared loss

- ◆ Gradient descent is a good, generic optimization algorithm.
- ◆ There are cases where we can do better than gradient descent.
- ◆ One such case is squared error loss function and 2-norm regularizer.
- ◆ We can get a closed-form solution!
- ◆ Let training data be a matrix \mathbf{X} of size $N \times D$, where $x_{n,d}$ is the value of the d th feature on the n th record, label as a vector \mathbf{Y} of dimension N , and weight as a column vector \mathbf{w} of size D .
- ◆ We have $\mathbf{a} = \mathbf{X}\mathbf{w}$, which has dimension N . That is $a_n = [\mathbf{X}\mathbf{w}]_n = \sum_d X_{n,d} w_d$.
- ◆ Here, we can take \mathbf{a} as \mathbf{Y} .
- ◆ We should minimize $\frac{1}{2} \sum_n (\hat{Y}_n - Y_n)^2$ following squared error. In vector form, it is $\frac{1}{2} \|\hat{\mathbf{Y}} - \mathbf{Y}\|^2$

Closed-form Optimization for Squared loss

$$\mathcal{L}(\mathbf{w}) = \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad \text{objective}$$

matrix notation

$$\underbrace{\begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,D} \\ x_{2,1} & x_{2,2} & \dots & x_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,1} & x_{N,2} & \dots & x_{N,D} \end{bmatrix}}_{\mathbf{X}} \underbrace{\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}}_{\mathbf{w}} = \underbrace{\begin{bmatrix} \sum_d x_{1,d} w_d \\ \sum_d x_{2,d} w_d \\ \vdots \\ \sum_d x_{N,d} w_d \end{bmatrix}}_{\hat{\mathbf{Y}}} \approx \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\mathbf{Y}}$$

equivalent loss

$$\min_w \mathcal{L}(w) = \frac{1}{2} \|\mathbf{X}w - \mathbf{Y}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Closed-form Optimization for Squared loss

$$\min_w \mathcal{L}(w) = \frac{1}{2} \|\mathbf{X}w - \mathbf{Y}\|^2 + \frac{\lambda}{2} \|w\|^2 \quad \text{objective}$$

$$\begin{aligned} \nabla_w \mathcal{L}(w) &= \mathbf{X}^\top (\mathbf{X}w - \mathbf{Y}) + \lambda w \\ &= \mathbf{X}^\top \mathbf{X}w - \mathbf{X}^\top \mathbf{Y} + \lambda w \\ &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) w - \mathbf{X}^\top \mathbf{Y} \end{aligned} \quad \text{gradient}$$

At optima the **gradient = 0**

$$\begin{aligned} &(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) w - \mathbf{X}^\top \mathbf{Y} = 0 \\ \iff &(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D) w = \mathbf{X}^\top \mathbf{Y} \\ \iff &w = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{Y} \end{aligned}$$

exact
closed-form
solution

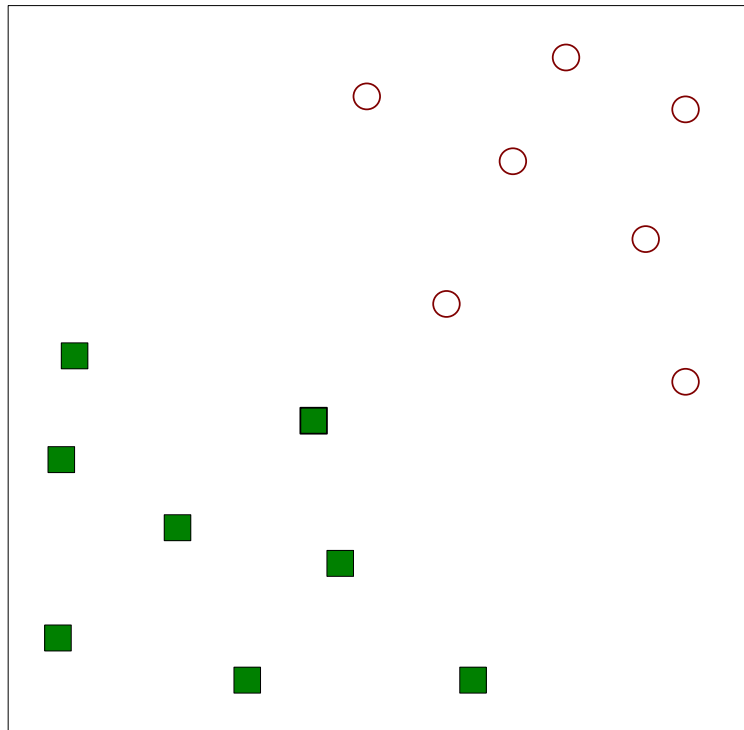
Matrix inversion vs. gradient descent

- ◆ Assume, we have D **features** and N **points**
- ◆ Overall time via **matrix inversion**
 - The closed form solution involves computing:

$$\mathbf{w} = \left(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

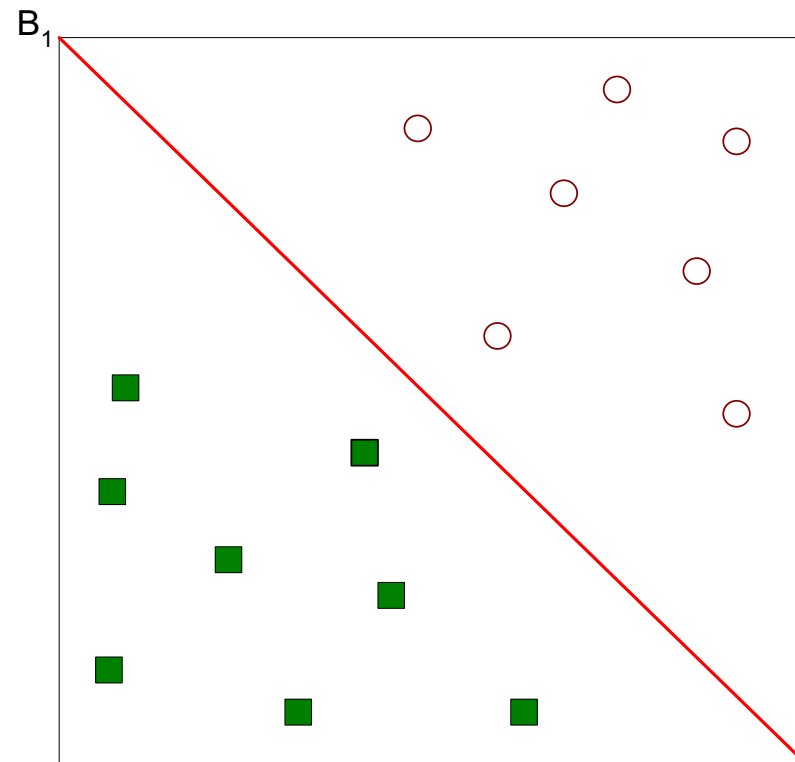
- Total time is $O(D^2N + D^3 + DN)$, assuming $O(D^3)$ matrix inversion
 - If $N > D$, then total time is $O(D^2N)$
- ◆ Overall time via **gradient descent**
 - Gradient: $\frac{d\mathcal{L}}{d\mathbf{w}} = \sum_n -2(y_n - \mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n + \lambda \mathbf{w}$
 - Each iteration: $O(ND)$; T iterations: $O(TND)$
- ◆ Which one is faster?
 - **Small problems** $D < 100$: probably faster to run **matrix inversion**
 - **Large problems** $D > 10,000$: probably faster to run **gradient descent**

Support Vector Machines



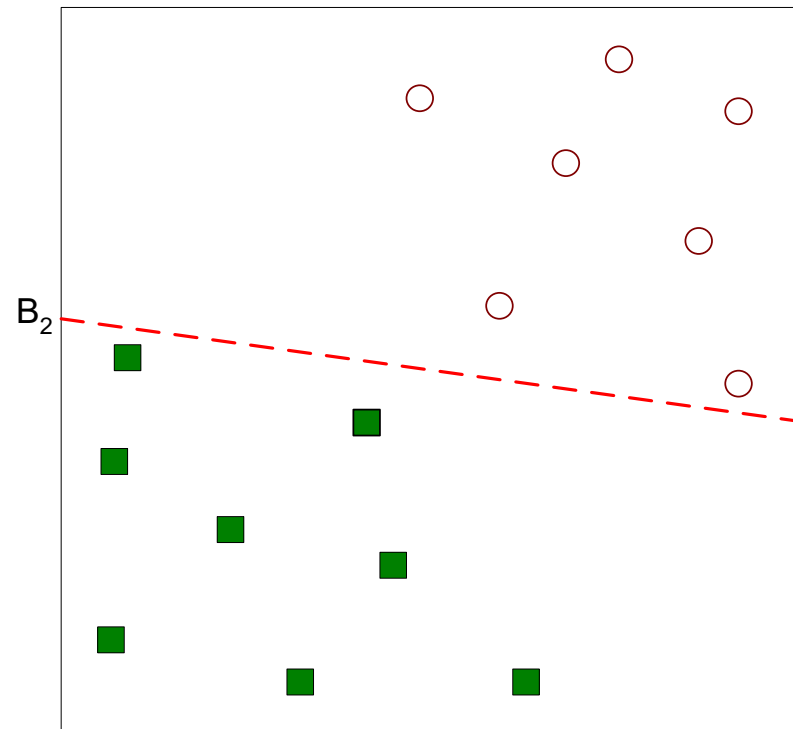
- Find a linear hyperplane (decision boundary) that will separate the data

Support Vector Machines



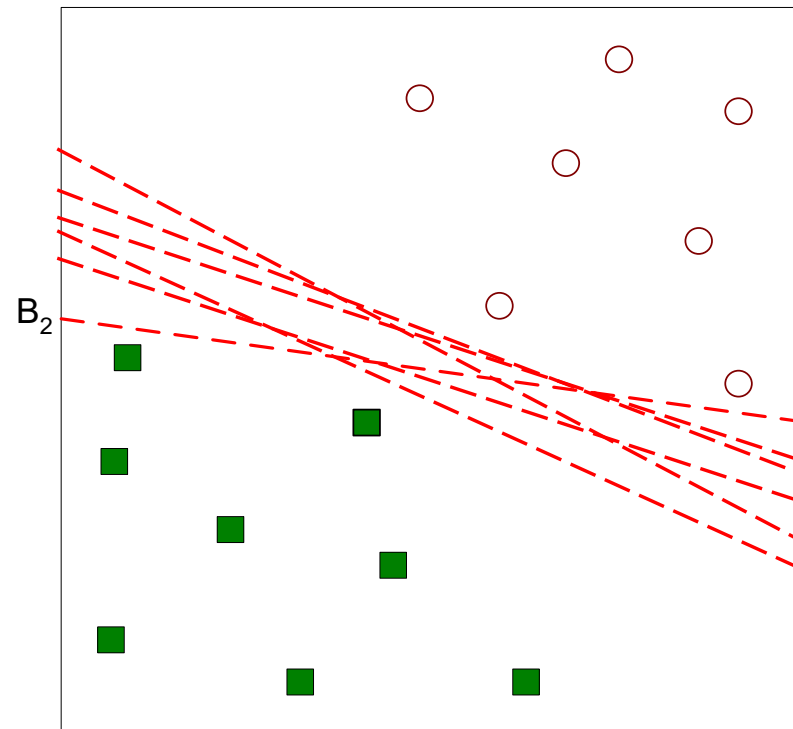
- One Possible Solution

Support Vector Machines



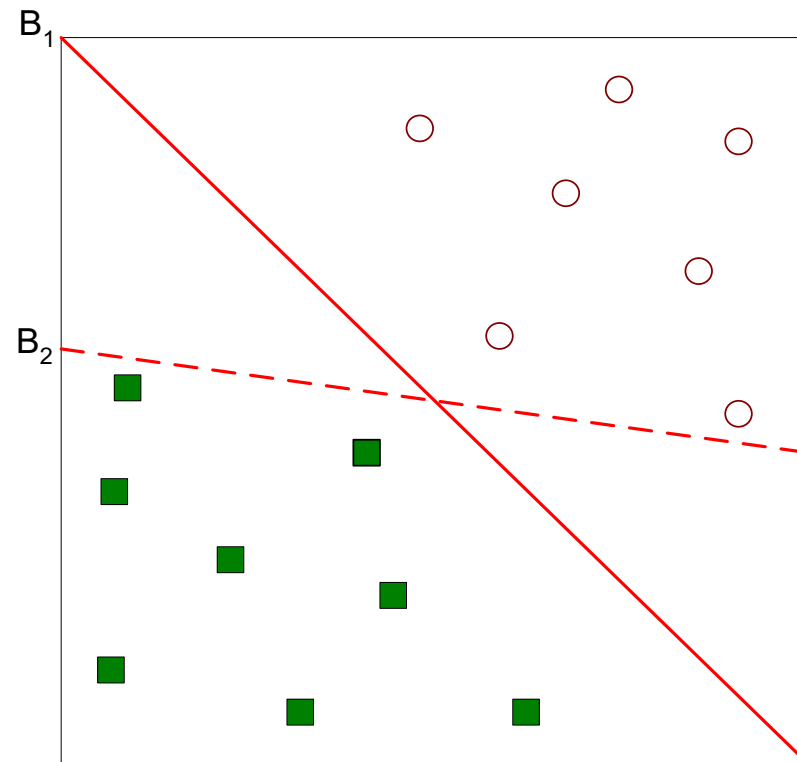
- Another possible solution

Support Vector Machines



- Other possible solutions

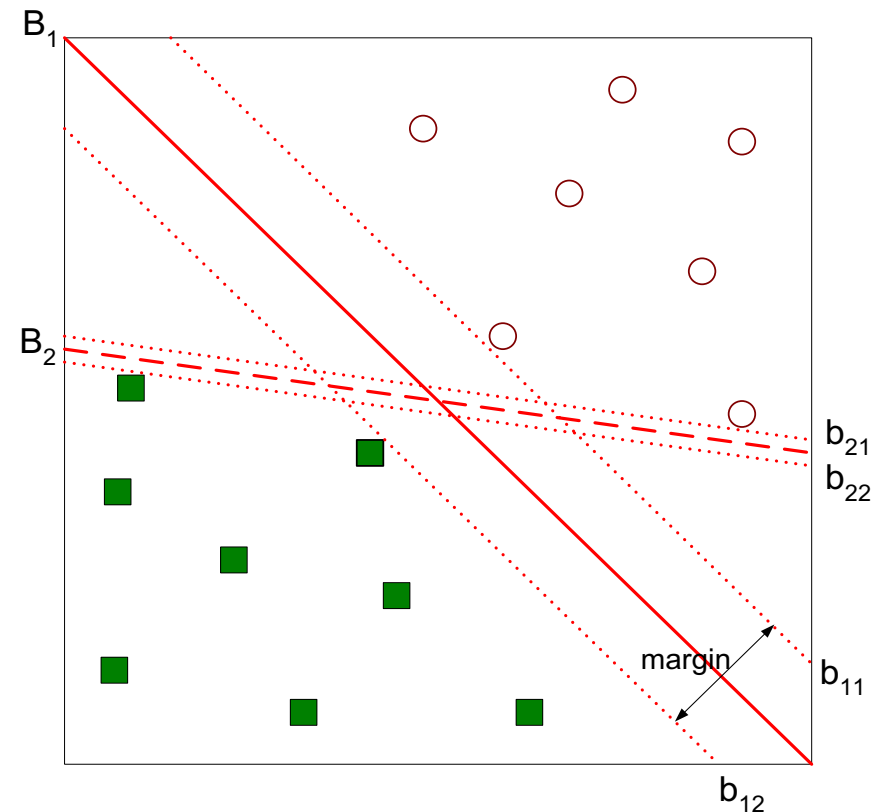
Support Vector Machines



- Which one is better? B_1 or B_2 ?
- How do you define better?

Support Vector Machines

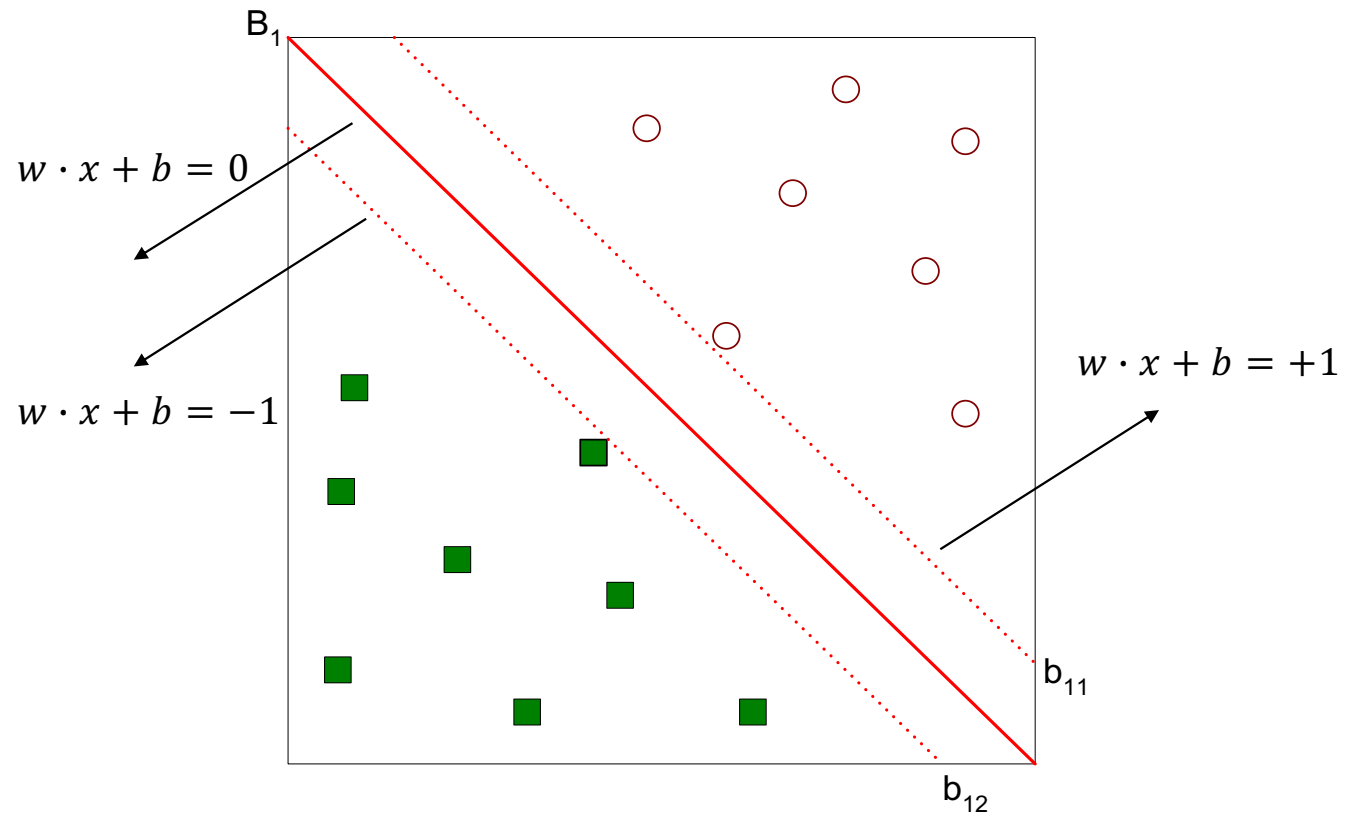
- For every separating hyperplane B_i , there is a pair of margin hyperplanes, b_{i1} and b_{i2} such that they touch the closest instances of classes.
- The margin is the distance between the pair of margin hyperplanes.
- Find the hyperplane that **maximizes** the margin $\Rightarrow B_1$ is better than B_2 .
- A larger margin tends to have a better generalization error.



Support Vector Machines

- The distance of any point x from the hyperplane is
$$D(x) = \frac{|w \cdot x + b|}{\|w\|}$$
- Let the closest points from the hyperplane be k_+ and k_- .
- We have the constraints

$$\frac{w^T x_i + b}{\|w\|} \geq k_+ \quad \text{if } y_i = 1,$$
$$\frac{w^T x_i + b}{\|w\|} \leq -k_- \quad \text{if } y_i = -1.$$

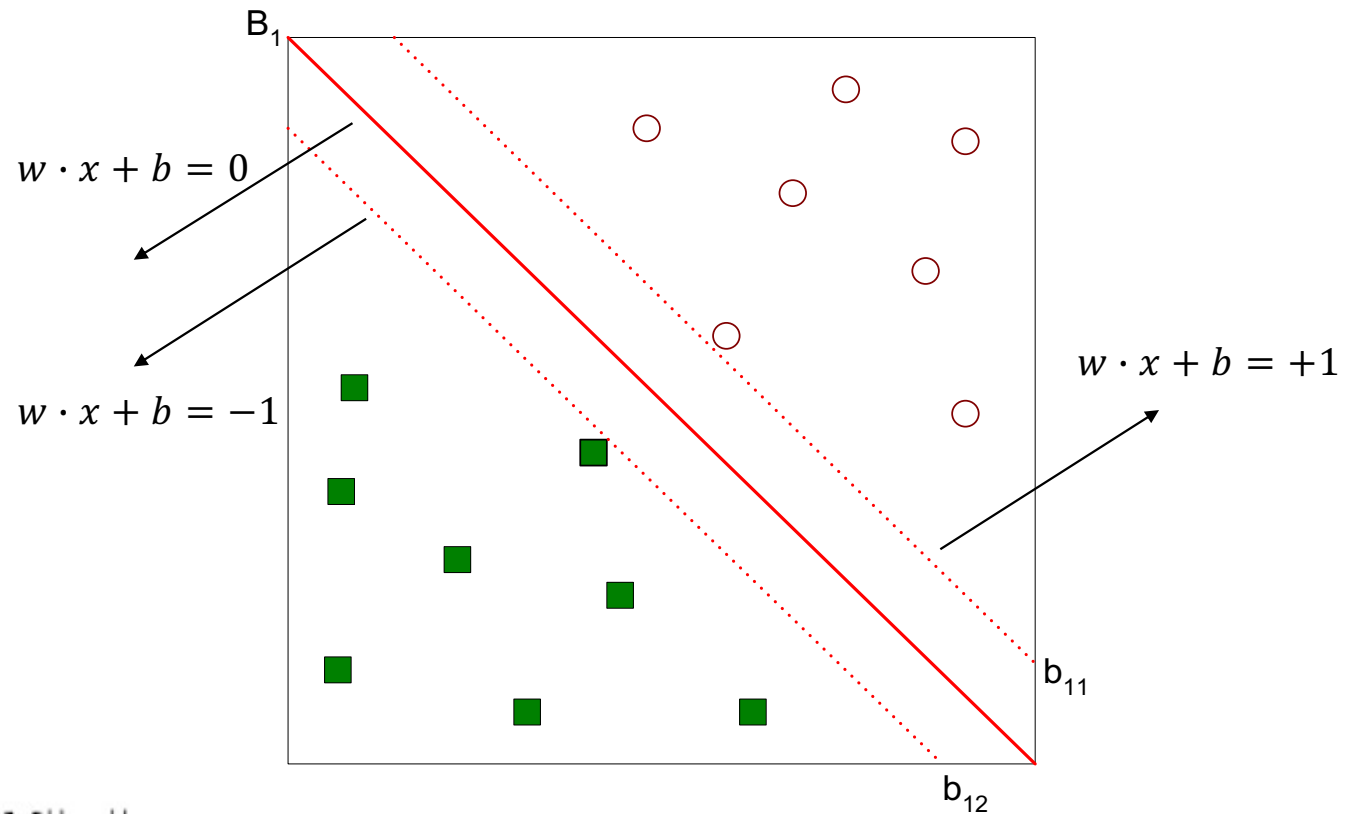


$$\text{Margin} = \frac{2}{\|\vec{w}\|}$$

Support Vector Machines

- Then, we have the product of y_i and $(w^T x_i + b)$ as $y_i(w^T x_i + b) \geq M\|w\|$ where M is a parameter.
- If $k_+ = M$ and $k_- = M$, then margin = $k_+ - k_- = 2M$.
- The optimization problem is:

$$\begin{aligned} & \max_{\mathbf{w}, b} M \\ & \text{subject to } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq M\|\mathbf{w}\|. \end{aligned}$$



$$\text{Margin} = \frac{2}{\|\vec{w}\|}$$

Linear SVM

- We measure the distance as follows, $k_+ = \frac{1}{\|w\|} w \cdot x_+ + b - 1$, and $k_- = \frac{1}{\|w\|} w \cdot x_- + b - 1$.
- $M = \frac{1}{2} [k_+ - k_-] = \frac{1}{2} \left[\frac{1}{\|w\|} w \cdot x_+ + b - 1 - \frac{1}{\|w\|} w \cdot x_- - b + 1 \right] = \frac{1}{2} \left[\frac{1}{\|w\|} w \cdot x_+ - \frac{1}{\|w\|} w \cdot x_- \right] = \frac{1}{2} \left[\frac{1}{\|w\|} (1 - b) - \frac{1}{\|w\|} (-1 - b) \right] = \frac{1}{\|w\|}$
- The size of the margin M is inversely proportional to the norm of the weight vector $\|w\|$.
- Maximize M amounts to minimizing $\|w\|^2$. The SVM is commonly represented as

$$\begin{aligned} & \min_{\mathbf{w}, b} \frac{\|\mathbf{w}\|^2}{2} \\ & \text{subject to} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1. \end{aligned}$$

Learning Model Parameters

$$\begin{aligned} & \min_{\mathbf{w}, b} \frac{\|\mathbf{w}\|^2}{2} \\ & \text{subject to} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1. \end{aligned}$$

- The equation above represents a constrained optimization problem with linear inequalities. And the objective function is convex and quadratic with respect to w .
- We can rewrite the objective function, which is known Lagrangian primal problem $L_P = \frac{1}{2} \|w\|^2 - \sum_i \lambda_i [y_i(w^T x_i + b) - 1]$, where the parameters $\lambda_i \geq 0$ correspond to the constraints and are called the Lagrange multipliers.

SVM Target

- Optimization (Quadratic Programming):

$$\min_{w,b} \frac{\|w\|^2}{2}$$
$$y_i(w^T x_i + b) \geq 1, \forall i$$

- Solved by Lagrange multiplier method:

$$L_P = \frac{1}{2} \|w\|^2 - \sum_i \lambda_i [y_i(w^T x_i + b) - 1]$$

where λ is the Lagrange multiplier

Lagrangian

- Consider optimization problem:

$$\min_w f(w)$$
$$h_i(w) = 0, \forall 1 \leq i \leq l$$

- Lagrangian:

$$\mathcal{L}(w, \boldsymbol{\beta}) = f(w) + \sum_i \beta_i h_i(w)$$

where β_i 's are called Lagrange multipliers

Lagrangian

- Consider optimization problem:

$$\min_w f(w)$$

$$h_i(w) = 0, \forall 1 \leq i \leq l$$

- Solved by setting derivatives of Lagrangian to 0

$$\frac{\partial \mathcal{L}}{\partial w_i} = 0; \quad \frac{\partial \mathcal{L}}{\partial \beta_i} = 0$$

Generalized Lagrangian

- Consider optimization problem:

$$\min_w f(w)$$

$$g_i(w) \leq 0, \forall 1 \leq i \leq k$$

$$h_j(w) = 0, \forall 1 \leq j \leq l$$

- Generalized Lagrangian:

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_i \alpha_i g_i(w) + \sum_j \beta_j h_j(w)$$

where α_i, β_j 's are called Lagrange multipliers

Lagrange Duality

- The primal problem

$$p^* := \min_w f(w) = \min_w \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta)$$

- The dual problem

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

- Always true:

$$d^* \leq p^*$$

Lagrange Duality

- Theorem: under **proper conditions**, there exists (w^*, α^*, β^*) such that

$$d^* = \mathcal{L}(w^*, \alpha^*, \beta^*) = p^*$$

Moreover, (w^*, α^*, β^*) satisfy Karush-Kuhn-Tucker (KKT) conditions:

$$\frac{\partial \mathcal{L}}{\partial w_i} = 0, \quad \alpha_i g_i(w) = 0$$

$$g_i(w) \leq 0, \quad h_j(w) = 0, \quad \alpha_i \geq 0$$

Learning Model Parameters

$$\begin{aligned} & \min_{\mathbf{w}, b} \frac{\|\mathbf{w}\|^2}{2} \\ & \text{subject to} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1. \end{aligned}$$

Given $L_P = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i \lambda_i [y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1]$, where the parameters $\lambda_i \geq 0$ correspond to the constraints and are called the Lagrange multipliers, to minimize the Lagrangian, we have,

$$\frac{\partial L_P}{\partial \mathbf{w}} = 0, \rightarrow \mathbf{w} = \sum_i \lambda_i y_i \mathbf{x}_i$$

$$\frac{\partial L_P}{\partial b} = 0, \rightarrow 0 = \sum_i \lambda_i y_i$$

With the Karash-Kuhn-Tucker (KKT) conditions between (\mathbf{w}, b) and λ_i , we also have $\lambda_i [y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1] = 0$.

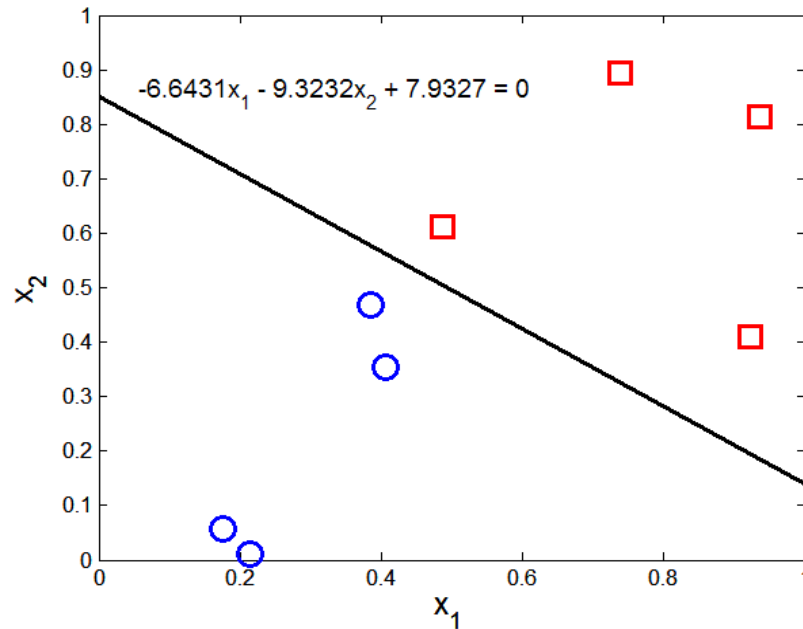
Learning Model Parameters

$$\begin{aligned} & \min_{\mathbf{w}, b} \frac{\|\mathbf{w}\|^2}{2} \\ & \text{subject to} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1. \end{aligned}$$

With all together, we have

$$\begin{aligned} & \max_{\lambda_i} \quad \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ & \text{subject to} \quad \sum_{i=1}^n \lambda_i y_i = 0, \\ & \quad \lambda_i \geq 0. \end{aligned}$$

An Example of Linear SVM

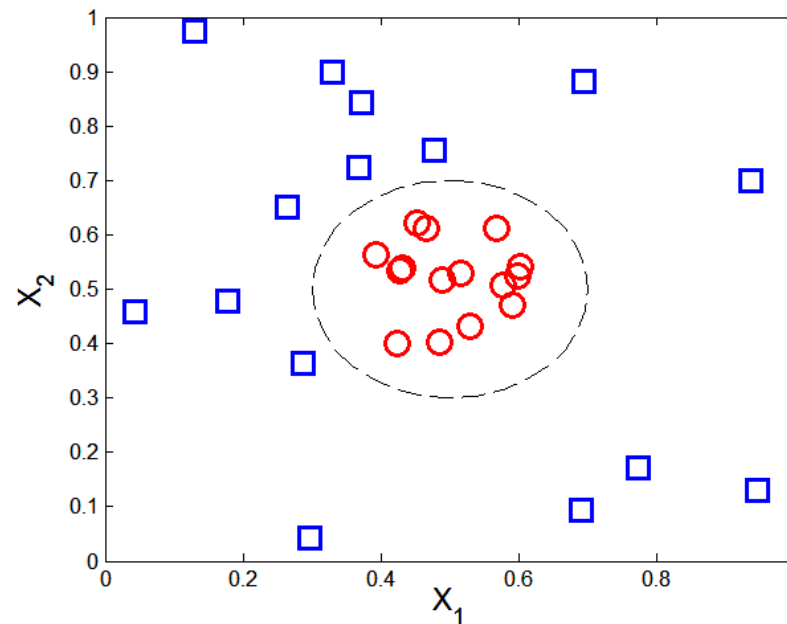


Support vectors

x_1	x_2	y	λ
0.3858	0.4687	1	65.5261
0.4871	0.611	-1	65.5261
0.9218	0.4103	-1	0
0.7382	0.8936	-1	0
0.1763	0.0579	1	0
0.4057	0.3529	1	0
0.9355	0.8132	-1	0
0.2146	0.0099	1	0

Nonlinear Support Vector Machines

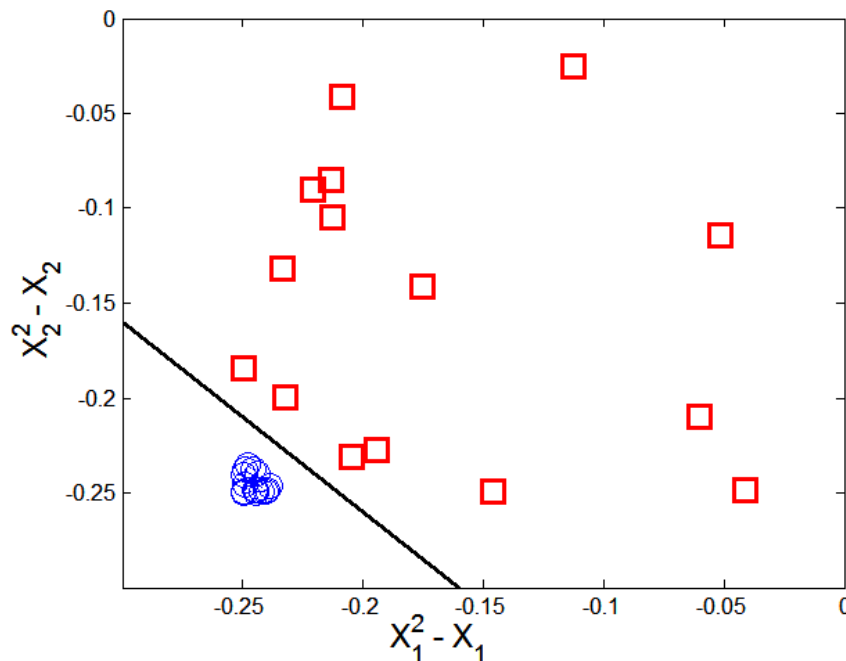
- What if decision boundary is not linear?



$$y(x_1, x_2) = \begin{cases} 1 & \text{if } \sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2} > 0.2 \\ -1 & \text{otherwise} \end{cases}$$

Nonlinear Support Vector Machines

- Transform data into higher dimensional space



$$x_1^2 - x_1 + x_2^2 - x_2 = -0.46.$$

$$\Phi: (x_1, x_2) \rightarrow (x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, 1).$$

$$w_4 x_1^2 + w_3 x_2^2 + w_2 \sqrt{2}x_1 + w_1 \sqrt{2}x_2 + w_0 = 0.$$

Decision boundary:

$$\vec{w} \bullet \Phi(\vec{x}) + b = 0$$

Characteristics of SVM

- The learning problem is formulated as a convex optimization problem
 - Efficient algorithms are available to find the global minima
 - Many of the other methods use greedy approaches and find locally optimal solutions
 - High computational complexity for building the model
- Robust to noise
- Overfitting is handled by maximizing the margin of the decision boundary,
- SVM can handle irrelevant and redundant attributes better than many other techniques
- The user needs to provide the type of kernel function and cost function
- Difficult to handle missing values
- What about categorical variables?

Slides Credit

- [1] Subhransu Maji. Linear model in CMPSCI689.
- [2] Yingyu Liang. SVM II in COS495
- [3] Introduction to Machine Learning 2nd edition, Pang-Ning Tan, Michael Steinbach, Vipin Kumar, and Anuj Karpatne
- [4] A Course in Machine Learning, Hal Daume III.