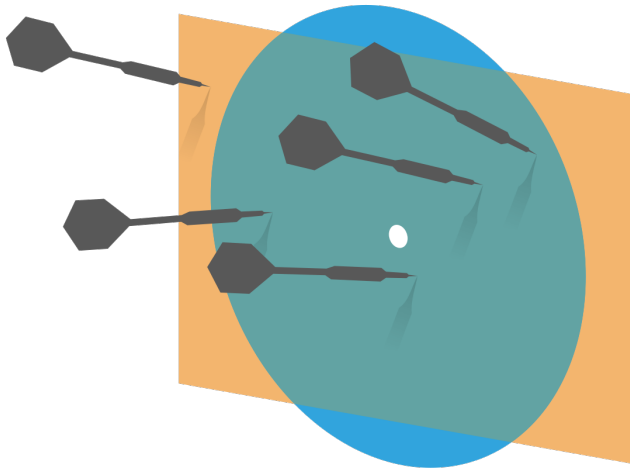


Probability and Statistics for Computer Science



“...many problems are naturally
classification problems” ---Prof.
Forsyth

Credit: wikipedia

Last time

- ✱ Decision tree (II)
- ✱ Random forest
- ✱ Support Vector Machine (I)

Objectives

- ✱ Support Vector Machine (II)
 - ✱ Hinge loss + Regularization
 - ✱ Convex function, Gradient Descent
Stochastic Gradient Descent
 - ✱ Training & Validation
- ✱ Naïve Bayesian Classifier

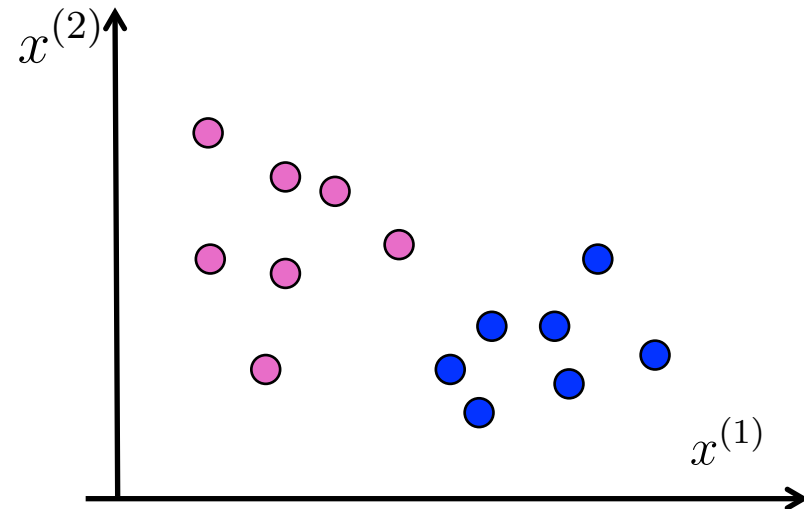
Motivation for Studying Support Vector Machine

- ✱ When solving a classification problem, it is good to try several techniques.
- ✱ Criteria to consider in choosing the classifier include
 - ✱ Accuracy ✓
 - ✱ Training speed
 - ✱ Classification speed ✓
 - ✱ Performance with small training set
 - ✱ Interpretability ✓

Stochastic Gradient Descent

SVM problem formulation

- ✱ At first we assume a binary classification problem
- ✱ The training set consists of N items
 - ✱ Feature vectors x_i of dimension d
 - ✱ Corresponding class labels $y_i \in \{\pm 1\}$
- ✱ We can picture the training data as a d -dimensional scatter plot with colored labels



Decision boundary of SVM

✱ SVM uses a hyperplane as its **decision boundary**

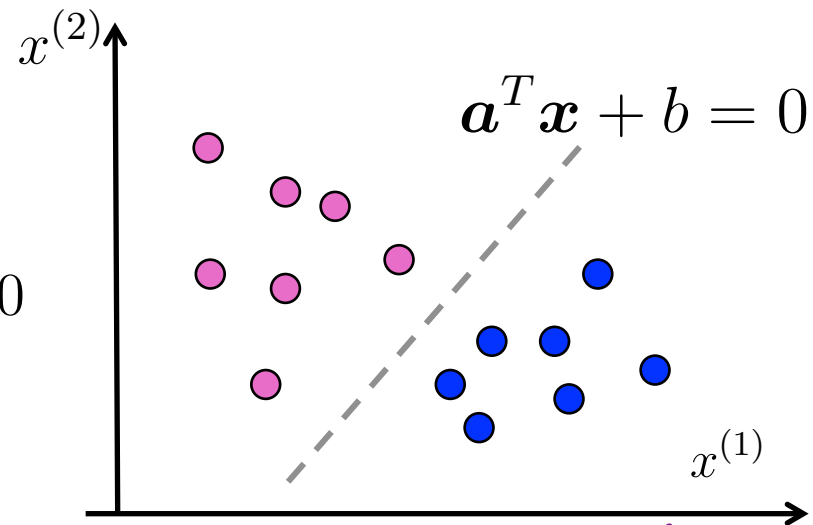
✱ The decision boundary is:

$$a_1x^{(1)} + a_2x^{(2)} + \dots + a_dx^{(d)} + b = 0$$

✱ In vector notation, the hyperplane can be written as:

$$\mathbf{a}^T \mathbf{x} + b = 0$$

$$[\mathbf{a}, b]$$



$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_d \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(d)} \end{bmatrix}$$

Classification function of SVM

- ✱ SVM assigns a class label to a feature vector according to the following rule:

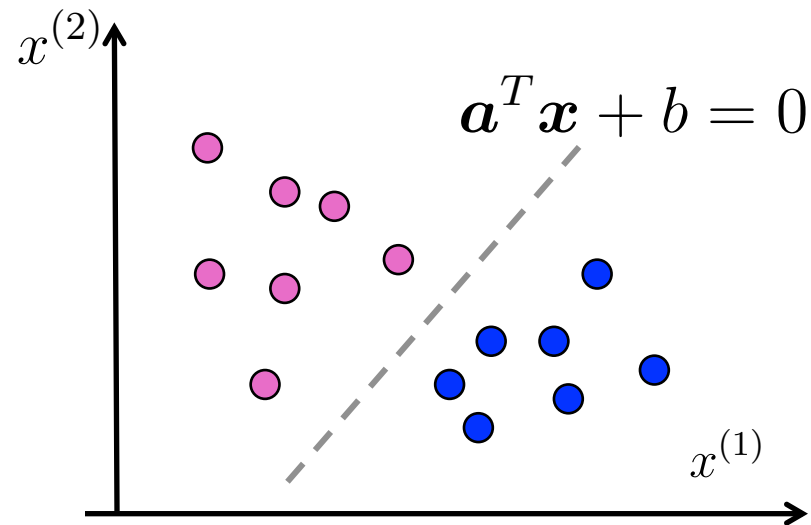
$$+1 \text{ if } \mathbf{a}^T \mathbf{x}_i + b \geq 0$$

$$-1 \text{ if } \mathbf{a}^T \mathbf{x}_i + b < 0$$

- ✱ In other words, the classification function is: $\text{sign}(\mathbf{a}^T \mathbf{x}_i + b)$

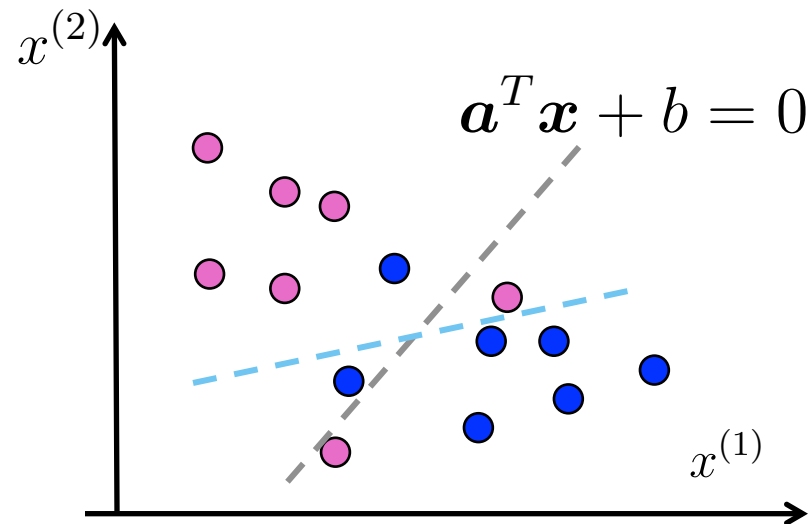
- ✱ Note that

- ✱ If $|\mathbf{a}^T \mathbf{x}_i + b|$ is small, then \mathbf{x}_i was close to the decision boundary
- ✱ If $|\mathbf{a}^T \mathbf{x}_i + b|$ is large, then \mathbf{x}_i was far from the decision boundary



What if there is no clean cut boundary?

- ✱ Some boundaries are better than others for the training data
- ✱ Some boundaries are likely more robust for run-time data
- ✱ We need a quantitative measure to decide about the boundary
- ✱ The **loss function** can help decide if one boundary is better than others



Loss function 1

- * For any given feature vector \mathbf{x}_i with class label $y_i \in \{\pm 1\}$, we want
 - * Zero loss if \mathbf{x}_i is classified correctly $\text{sign}(\mathbf{a}^T \mathbf{x}_i + b) = y_i$
 - * Positive loss if \mathbf{x}_i is misclassified $\text{sign}(\mathbf{a}^T \mathbf{x}_i + b) \neq y_i$
 - * If \mathbf{x}_i is misclassified, more loss is assigned if it's further away from the boundary

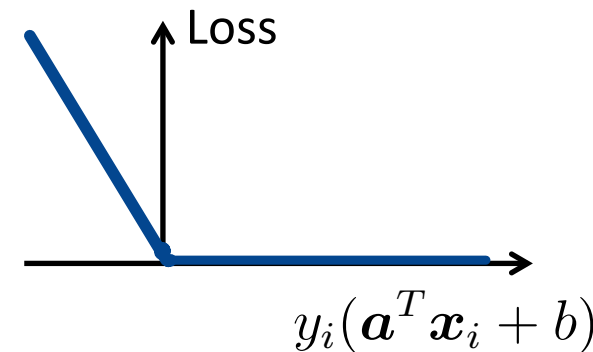
y_i is either -1 or +1

- * This loss function 1 meets the criteria above:

$$\max(0, -y_i(\mathbf{a}^T \mathbf{x}_i + b))$$

- * Training error cost

$$S(\mathbf{a}, b) = \frac{1}{N} \sum_{i=1}^N \max(0, -y_i(\mathbf{a}^T \mathbf{x}_i + b))$$



Q. What's the value of this function ?

$$\max(0, -y_i(\mathbf{a}^T \mathbf{x}_i + b)) \quad \text{if} \quad \text{sign}(\mathbf{a}^T \mathbf{x}_i + b) = y_i$$

- A. 0.
- B. others.

Q. What's the value of this function ?

$$\max(0, -y_i(\mathbf{a}^T \mathbf{x}_i + b)) \quad \text{if } \text{sign}(\mathbf{a}^T \mathbf{x}_i + b) \neq y_i$$

A. 0.

B. A value greater than or equal to 0.

Loss function 1

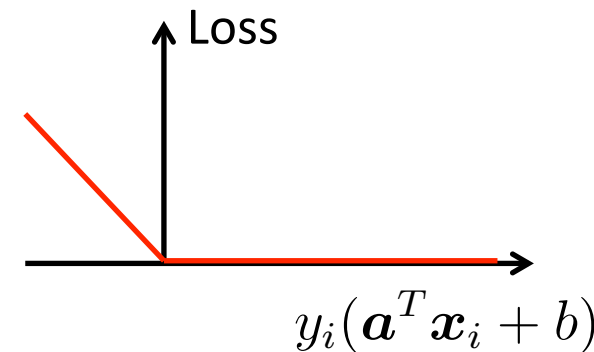
- ✱ For any given feature vector \mathbf{x}_i with class label $y_i \in \{\pm 1\}$, we want
 - ✱ Zero loss if \mathbf{x}_i is classified correctly $\text{sign}(\mathbf{a}^T \mathbf{x}_i + b) = y_i$
 - ✱ Positive loss if \mathbf{x}_i is misclassified $\text{sign}(\mathbf{a}^T \mathbf{x}_i + b) \neq y_i$
 - ✱ If \mathbf{x}_i is misclassified, more loss is assigned if it's further away from the boundary

- ✱ This loss function 1 meets the criteria above:

$$\max(0, -y_i(\mathbf{a}^T \mathbf{x}_i + b))$$

- ✱ Training error cost

$$S(\mathbf{a}, b) = \frac{1}{N} \sum_{i=1}^N \max(0, -y_i(\mathbf{a}^T \mathbf{x}_i + b))$$



The problem with loss function 1

- ✱ Loss function 1 does not distinguish between the following decision boundaries if they both classify \mathbf{x}_i correctly.
 - ✱ One passes the two classes closely
 - ✱ One that passes with a wider margin
- ✱ But leaving a larger margin gives robustness for run-time data- **the large margin principle**

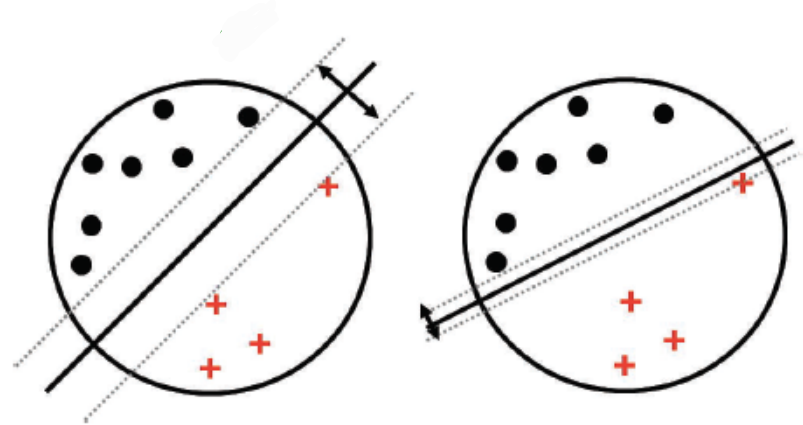


Figure 14.11 Illustration of the large margin principle. Left: a separating hyper-plane with large margin. Right: a separating hyper-plane with small margin.

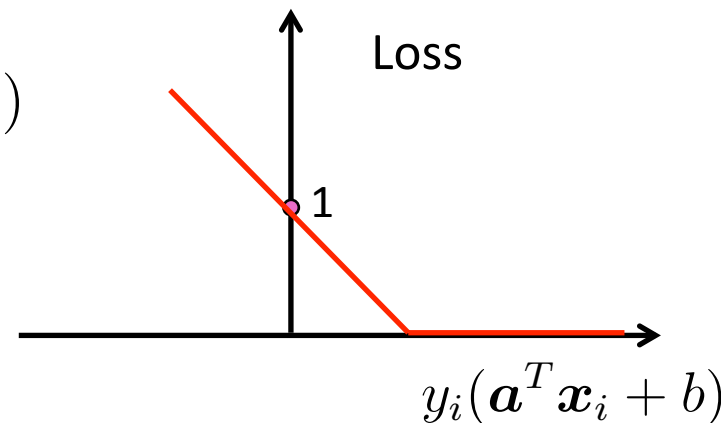
Loss function 2: the hinge loss

- ✱ We want to impose a small positive loss if \mathbf{x}_i is correctly classified but close to the boundary
- ✱ The **hinge loss** function meets the criteria above:

$$\max(0, 1 - y_i(\mathbf{a}^T \mathbf{x}_i + b))$$

- ✱ Training error cost

$$S(\mathbf{a}, b) = \frac{1}{N} \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{a}^T \mathbf{x}_i + b))$$



The problem with loss function 2

- ✱ Loss function 2 favors decision boundaries that have large $\|\mathbf{a}\|$ because increasing $\|\mathbf{a}\|$ can zero out the loss for a correctly classified \mathbf{x}_i near the boundary.
- ✱ But large $\|\mathbf{a}\|$ makes the classification function $\text{sign}(\mathbf{a}^T \mathbf{x}_i + b)$ extremely sensitive to small changes in \mathbf{x}_i and make it less robust to run-time data.
- ✱ So small $\|\mathbf{a}\|$ is better.

Hinge loss with regularization penalty

- ✱ We add a penalty on the square magnitude $\|\mathbf{a}\|^2 = \mathbf{a}^T \mathbf{a}$

\mathbf{x}

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

- ✱ Training error cost

$$S(\mathbf{a}, b) = \left[\frac{1}{N} \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{a}^T \mathbf{x}_i + b)) \right] + \lambda \left(\frac{\mathbf{a}^T \mathbf{a}}{2} \right)$$

- ✱ The **regularization parameter** λ trade off between these two objectives

Q. What does the penalty discourage?

$$S(\mathbf{a}, b) = \left[\frac{1}{N} \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{a}^T \mathbf{x}_i + b)) \right] + \lambda \left(\frac{\mathbf{a}^T \mathbf{a}}{2} \right)$$

- A. Too big a magnitude of the vector \mathbf{a}
- B. Too many data points in the training set

How to compute the decision boundary?

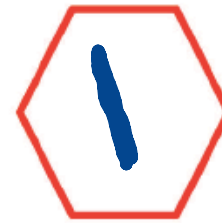
minimize Loss function $S(\vec{a}, b)$

$$(\vec{a}^*, b^*) = \operatorname{argmin}(S(\vec{a}, b))$$

$$\downarrow$$
$$\begin{bmatrix} a_1^* \\ a_2^* \\ \vdots \\ a_d^* \\ b^* \end{bmatrix}$$

Convex set and convex function

- ✱ If a set is convex, any line connecting two points in the set is completely included in the set



(a)

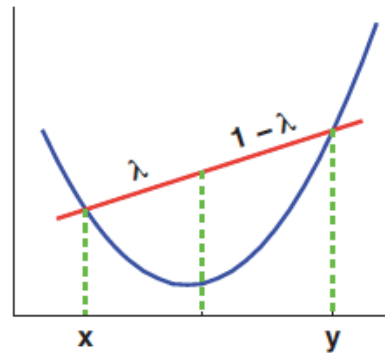


(b)

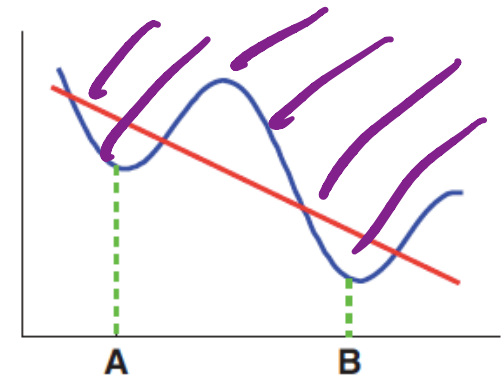
Figure 7.4 (a) Illustration of a convex set. (b) Illustration of a nonconvex set.

- ✱ A convex function: the area above the curve is convex

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y)$$



(a)

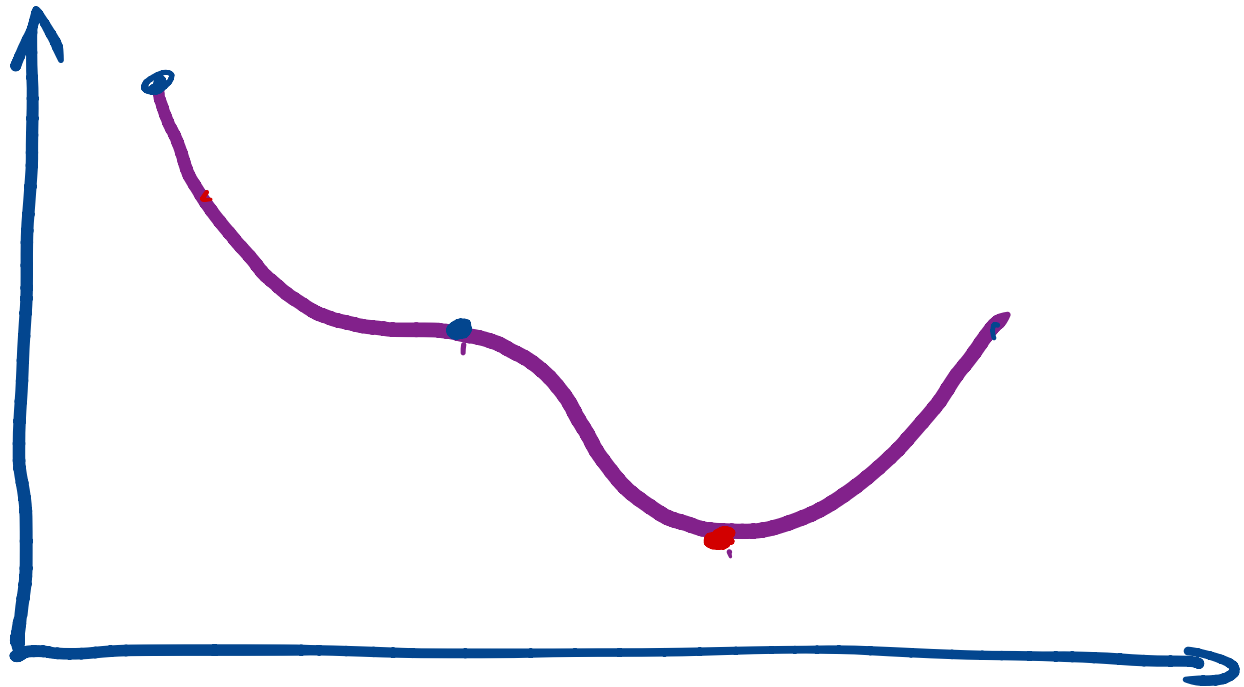


(b)

Q. Is this curve a convex curve?

A. YES

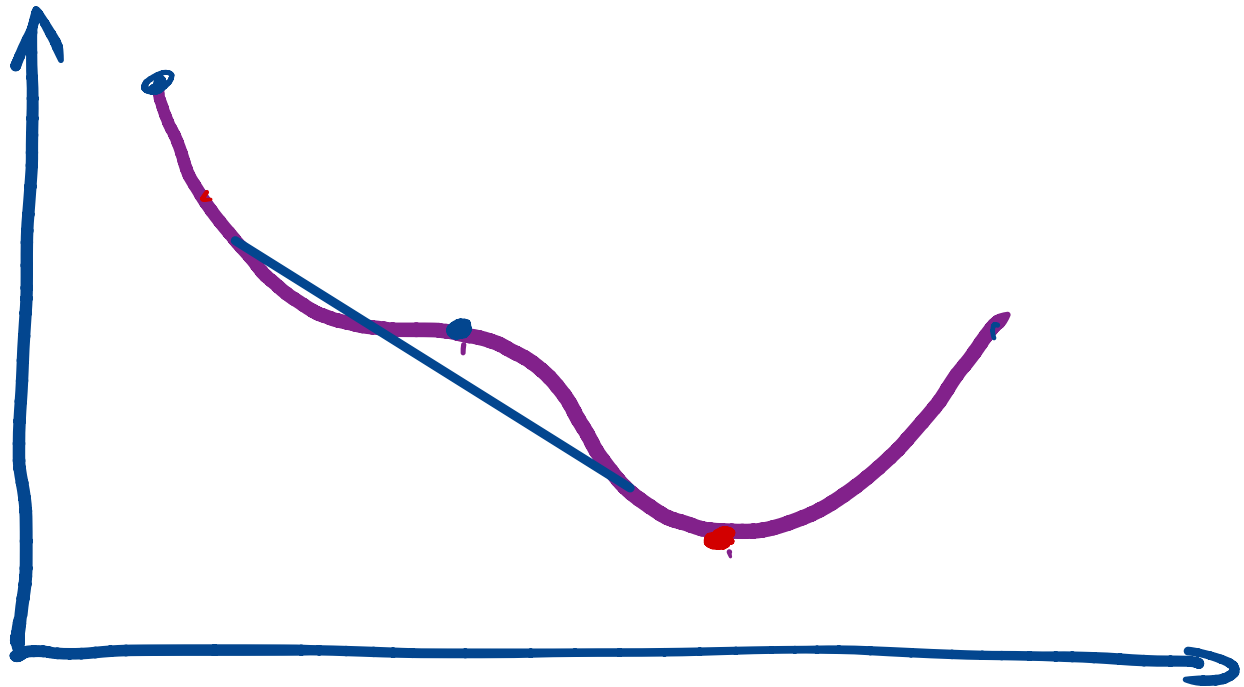
B. NO



Q. Is this curve a convex curve?

A. YES

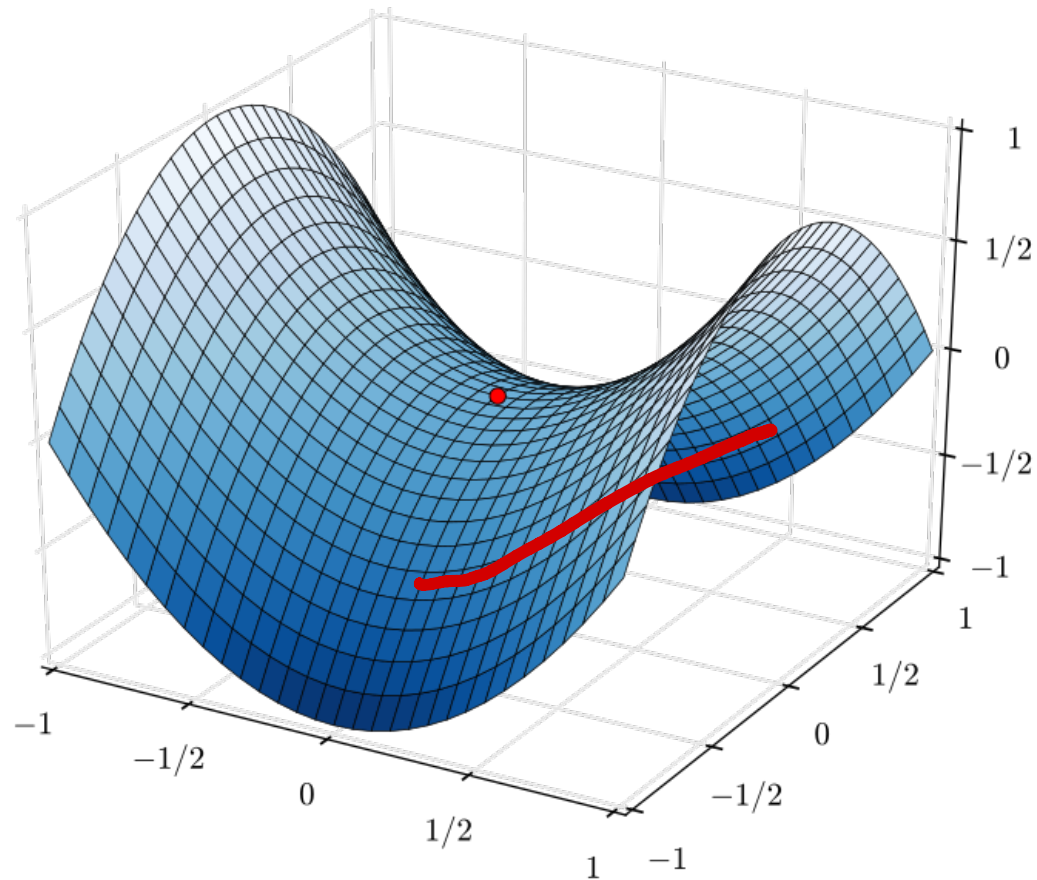
B. NO



Q. Is this surface convex?

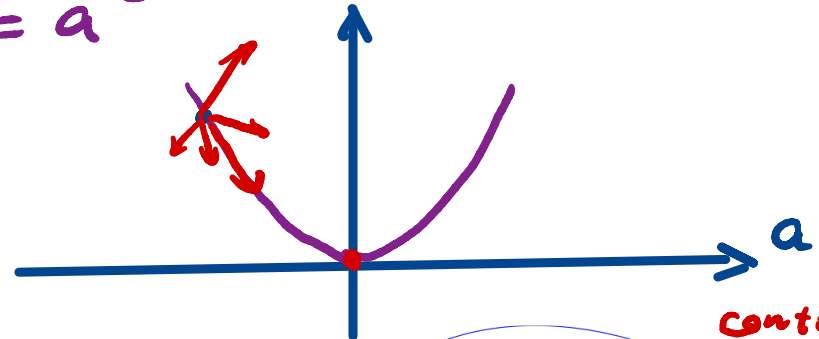
A. YES

B. NO



Iterative minimization by gradient descent

- * For a function such as $f(a) = a^2$



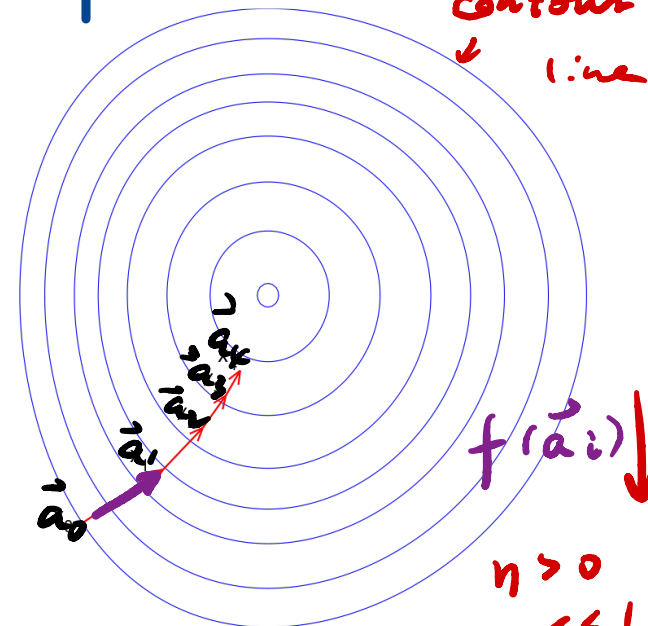
- * A convex surface $f(\vec{a}) = \text{height}(\vec{a})$



$$\vec{a}_i = \vec{a}_0 + \eta_0 \vec{p}$$

$$\vec{p} = -\nabla f$$

f is scalar



contour
line

$\eta > 0$
 $\eta \ll 1$

Gradient Descent

$$\vec{a} = [a_1, a_2, \dots, a_d]^T \quad \text{let's omit } b \text{ for now}$$

$$\text{Loss function: } f = f(\vec{a})$$

$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial a_1} \\ \frac{\partial f}{\partial a_2} \\ \vdots \\ \frac{\partial f}{\partial a_d} \end{pmatrix}$$

\vec{a}^* minimizes $f(\vec{a})$

$$\vec{a}_0, \vec{a}_1, \dots, \vec{a}_n$$

$$f(\vec{a}_0) > \dots > f(\vec{a}_n)$$

$$\vec{a}_{n+1} = \vec{a}_n + \eta \vec{p}_n$$

Taylor expansion

$$f(\vec{a}_{n+1}) = f(\vec{a}_n + \eta_n \vec{p}_n) = f(\vec{a}_n) + \eta_n (\nabla f)^T \vec{p}_n + O(\eta_n^2)$$

$$\text{if } \vec{p}_n = -\nabla f(\vec{a}_n)$$

$$\begin{aligned} f(\vec{a}_{n+1}) &= f(\vec{a}_n) - \eta_n (\nabla f)^T \nabla f \\ &= f(\vec{a}_n) - \eta_n \|\nabla f\|^2 \end{aligned}$$

$$\eta_n > 0$$

$$\eta_n < 1$$

learning
rate
or stepsize

Stochastic gradient descent

$$\text{if } f(\vec{a}) = \frac{1}{k} \sum_{j=1}^k Q(\vec{a}, j)$$

$k \rightarrow$ # of data
in training set

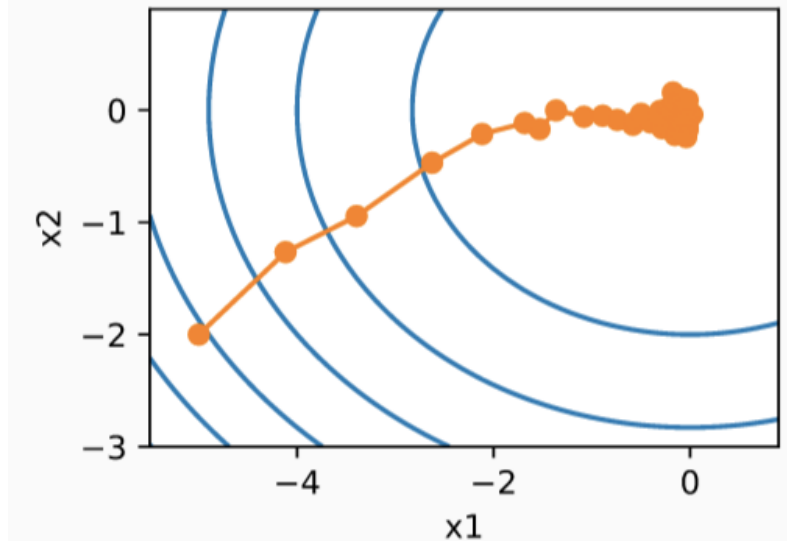
Stochastic gradient descent

$$\text{approximates } f(\vec{a}) \doteq g(\vec{a}) = \frac{1}{m} \sum_{i=1}^m Q(\vec{a}, i)$$

$$x_k \in \{x_i\}$$

i is RV. uniform in $[1, k]$

We often choose $m = 1$



The difference btw GD and SGD

GD

Loss:

$$f(\vec{a}) = \frac{1}{K} \sum_{j=1}^K Q(\vec{a}, j) + \text{penalty}$$

$$Q = \max(0, 1 - y_j(\vec{a}^T x_j + b))$$

$$\vec{a}_{n+1} = \vec{a}_n - \eta \nabla f(\vec{a}_n)$$

$$\lim_{n \rightarrow \infty} \vec{a}_n = \underset{\vec{a}}{\operatorname{argmin}} (f(\vec{a}))$$

if f is convex

SGD

Loss:

$$g(\vec{a}) = \frac{1}{m} \sum_{i=1}^m Q(\vec{a}, i) + \text{penalty}$$

$$g(\vec{a}) = f(\vec{a}) + \epsilon$$

↑
noise

$$\vec{a}_{n+1} = \vec{a}_n - \eta \nabla g(\vec{a}_n)$$

$$\lim_{n \rightarrow \infty} E[(\vec{a}_n - \vec{a}^*)^2] = 0$$

not always

given convex Loss
and other
conditions

Update parameters of the hyperplane during the stochastic gradient descent

- * Since $S_k(\mathbf{a}, b) = \max(0, 1 - y_k(\mathbf{a}^T \mathbf{x}_k + b))$ and $S_0(\mathbf{a}, b) = \lambda \left(\frac{\mathbf{a}^T \mathbf{a}}{2} \right)$

We have the following updating equations:

$$S(\mathbf{a}, b) = S_k + S_0 \quad \text{if } m=1$$

$$\text{if } y_k(\mathbf{a}^T \mathbf{x}_k + b) \geq 1$$

$$\mathbf{a} \leftarrow \mathbf{a} - \eta(\lambda \mathbf{a})$$

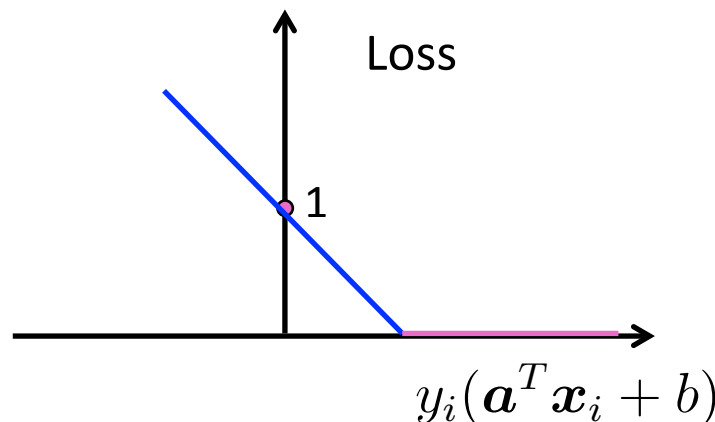
$$b \leftarrow b$$

$$\text{if } y_k(\mathbf{a}^T \mathbf{x}_k + b) < 1$$

$$\mathbf{a} \leftarrow \mathbf{a} - \eta(\lambda \mathbf{a} - y_k \mathbf{x}_k)$$

$$b \leftarrow b - \eta(-y_k)$$

We often set $\eta_n = \frac{1}{n}$



Training procedure-minimizing the cost function

- * The training error cost $S(\mathbf{a}, b)$ is a function of decision boundary parameters (\mathbf{a}, b) , so it can help us find the best decision boundary.
- * Fix λ and set some initial values for (\mathbf{a}, b)
- * Search iteratively for (\mathbf{a}, b)
- * Repeat the previous steps for several values of λ and choose the one that gives the decision boundary with best accuracy on a validation data set.

Validation/testing of SVM model

- * Split the labeled data into **training**, **validation** and **test** sets.
- * For each choice of λ , run stochastic gradient descent to find the best decision boundary parameters (\mathbf{a} , \mathbf{b}) using the training set.
- * Choose the best λ based on accuracy on the validation set.
- * Finally evaluate the SVM's accuracy on the **test** set.
- * This process avoids overfitting the data.

Extension to multiclass classification

* All vs. all

1 2 3

* Train a separate binary classifier for each pair of classes.

1 2

* To classify, run all classifiers and see which class it will be labeled most with.

1 3

2 3

* Computational complexity is quadratic to the number of classes.

$$C \rightarrow \binom{C}{2}$$

* One vs. all

1 vs not 1
2 vs not 2
3 vs not 3

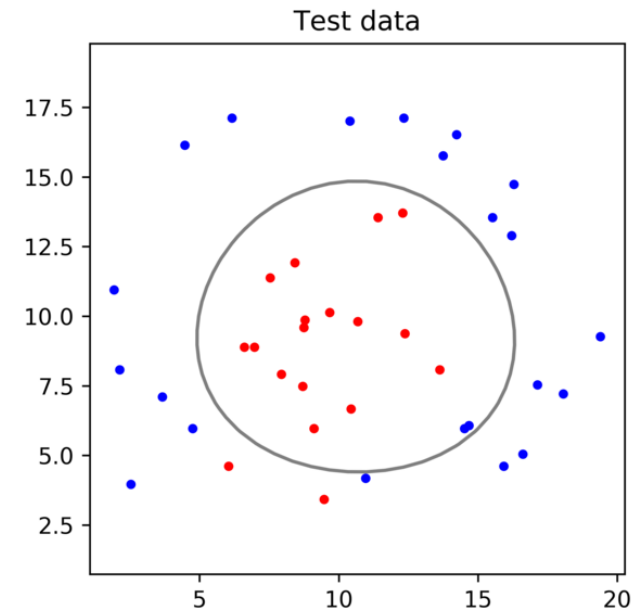
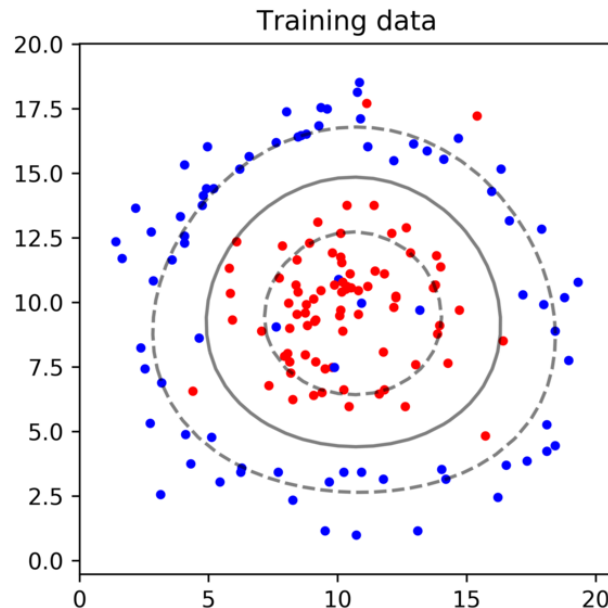
* Train a separate binary classifier for each class against all else.

* To classify, run all classifiers and see which label gets the highest score

* Computational complexity scales linearly.

What if the data is inseparable linearly?

- ✱ There is a chance the data is inseparable
- ✱ Use the non-linear **SVM with kernels!**
- ✱ Decision boundary is curved



Naïve Bayes classifier

* Training

- * Use the training data $\{(x_i, y_i)\}$ to estimate a probability model $P(y|\mathbf{x})$ *

$$\mathbf{x} = [x^{(1)} \ x^{(2)} \ \dots \ x^{(d)}]^T$$

ie. petal length.
petal width ...

- * Assume that the features of $\{\mathbf{x}\}$ are conditionally independent given the class label y

$$P(D|\theta)$$

$$P(\mathbf{x}|y) = \prod_{j=1}^d P(x^{(j)}|y)$$

↑
Data

* Classification

- * Assign the label $\underset{y}{\operatorname{argmax}} P(y|\mathbf{x})$ to a feature vector \mathbf{x}

$$P(x|D \cap Y|D \cap Z|D) = P(x|D) P(Y|D) P(Z|D)$$

Additional References

- ✱ Robert V. Hogg, Elliot A. Tanis and Dale L. Zimmerman. “Probability and Statistical Inference”
- ✱ Kelvin Murphy, “Machine learning, A Probabilistic perspective”

See you next time

*See
You!*

