

Perceptrons

Vibhav Gogate

The University of Texas at Dallas

Several Slides courtesy of Carlos Guestrin, Luke Zettlemoyer, Vincent Ng, Pedro Domingos, and Dan Weld.

Linear Classifiers

- Inputs are **feature values** (x_1, x_2, \dots, x_n)
- Each feature has a **weight** (w_1, w_2, \dots, w_n)
- Sum is the **activation** $\sum_{i=0}^n w_i x_i$

Bias w_0 , Input $x_0=1$ for all examples

- If the activation is:
 - Positive, output *class* = +ve
 - Negative, output *class* = -ve

Example: Spam

- Imagine 3 features (spam is “positive” class):
 - free (number of occurrences of “free”)
 - money (occurrences of “money”)
 - BIAS (intercept, always has value 1)

$$\vec{w} \cdot \vec{x}$$

$$\sum_{i=0}^n w_i x_i$$

| | x | w | |
|--------------|-----------|-----------|-----------|
| “free money” | BIAS : 1 | BIAS : -3 | (1)(-3) + |
| | free : 1 | free : 4 | (1)(4) + |
| | money : 1 | money : 2 | (1)(2) + |
| | ... | ... | ... |
| | | | = 3 |

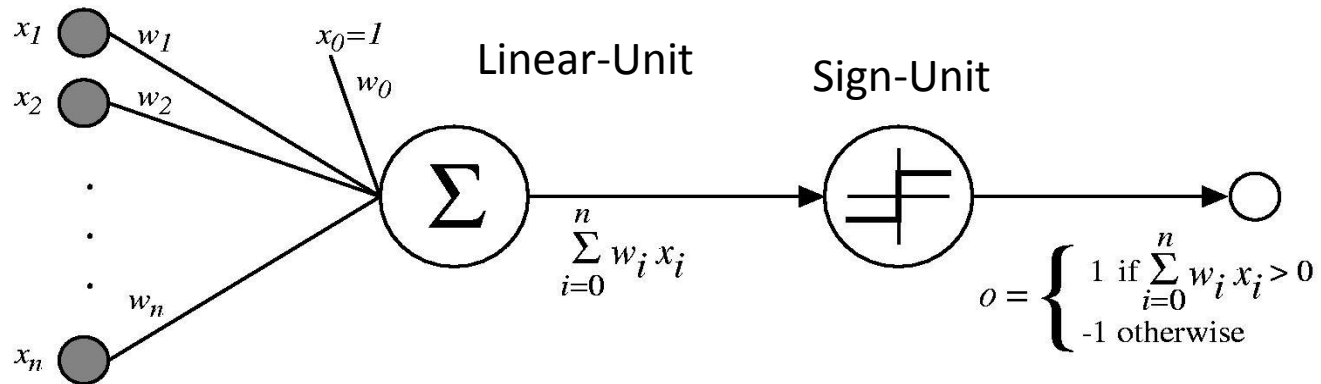
$$\vec{w} \cdot \vec{x} > 0 \Rightarrow \text{SPAM!!!}$$

Who needs probabilities?

- Previously: model data with distributions
- Joint: $P(X,Y)$
 - e.g. Naïve Bayes
- Conditional: $P(Y|X)$
 - e.g. Logistic Regression
- But wait, why probabilities?
- Lets try to be error-driven!

| mpg | cylinders | displacemen | horsepower | weight | acceleration | modelyear | maker |
|------|-----------|-------------|------------|--------|--------------|-----------|---------|
| good | 4 | 97 | 75 | 2265 | 18.2 | 77 | asia |
| bad | 6 | 199 | 90 | 2648 | 15 | 70 | america |
| bad | 4 | 121 | 110 | 2600 | 12.8 | 77 | europa |
| bad | 8 | 350 | 175 | 4100 | 13 | 73 | america |
| bad | 6 | 198 | 95 | 3102 | 16.5 | 74 | america |
| bad | 4 | 108 | 94 | 2379 | 16.5 | 73 | asia |
| bad | 4 | 113 | 95 | 2228 | 14 | 71 | asia |
| bad | 8 | 302 | 139 | 3570 | 12.8 | 78 | america |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| good | 4 | 120 | 79 | 2625 | 18.6 | 82 | america |
| bad | 8 | 455 | 225 | 4425 | 10 | 70 | america |
| good | 4 | 107 | 86 | 2464 | 15.5 | 76 | europa |
| bad | 5 | 131 | 103 | 2830 | 15.9 | 78 | europa |
| | | | | | | | |

Perceptron

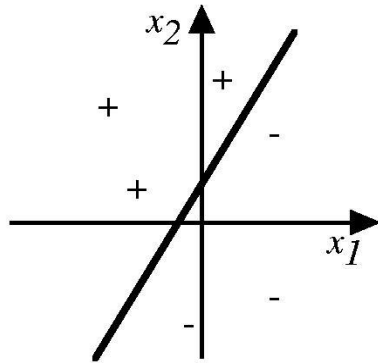


$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise.} \end{cases}$$

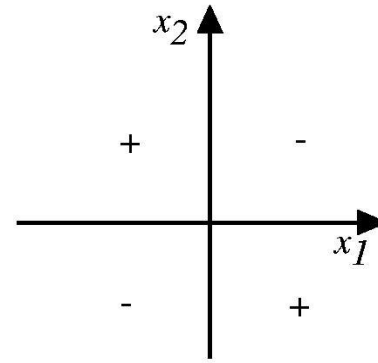
Sometimes we'll use simpler vector notation:

$$o(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w} \cdot \vec{x} > 0 \\ -1 & \text{otherwise.} \end{cases}$$

Decision Surface of a Perceptron



(a)



(b)

Represents some useful functions

- What weights represent $g(x_1, x_2) = AND(x_1, x_2)$?

But some functions not representable

- All not linearly separable
- Therefore, we'll want networks of these...

Perceptron Training Rule

$$w_i \leftarrow w_i + \Delta w_i$$

where

$$\Delta w_i = \eta(t - o)x_i$$

Where:

- $t = c(\vec{x})$ is target value
- o is perceptron output
- η is small constant (e.g., 0.1) called *learning rate*

Perceptron Training Rule

- Converges if the data is linearly separable
 - Provided the learning rate is sufficiently small
 - Proof on the class website (a bit involved)
- Convergence is not assured if data is not linearly separable
 - In fact in many cases, it will not converge
- Can we use some other algorithm to guarantee convergence?
 - YES!! – Gradient Descent
 - Gradient Descent yields a new rule for learning called **the Delta rule**

Gradient Descent

To understand, consider simpler *linear unit*, where

$$o = w_0 + w_1x_1 + \cdots + w_nx_n$$

Let's learn w_i 's that minimize the squared error

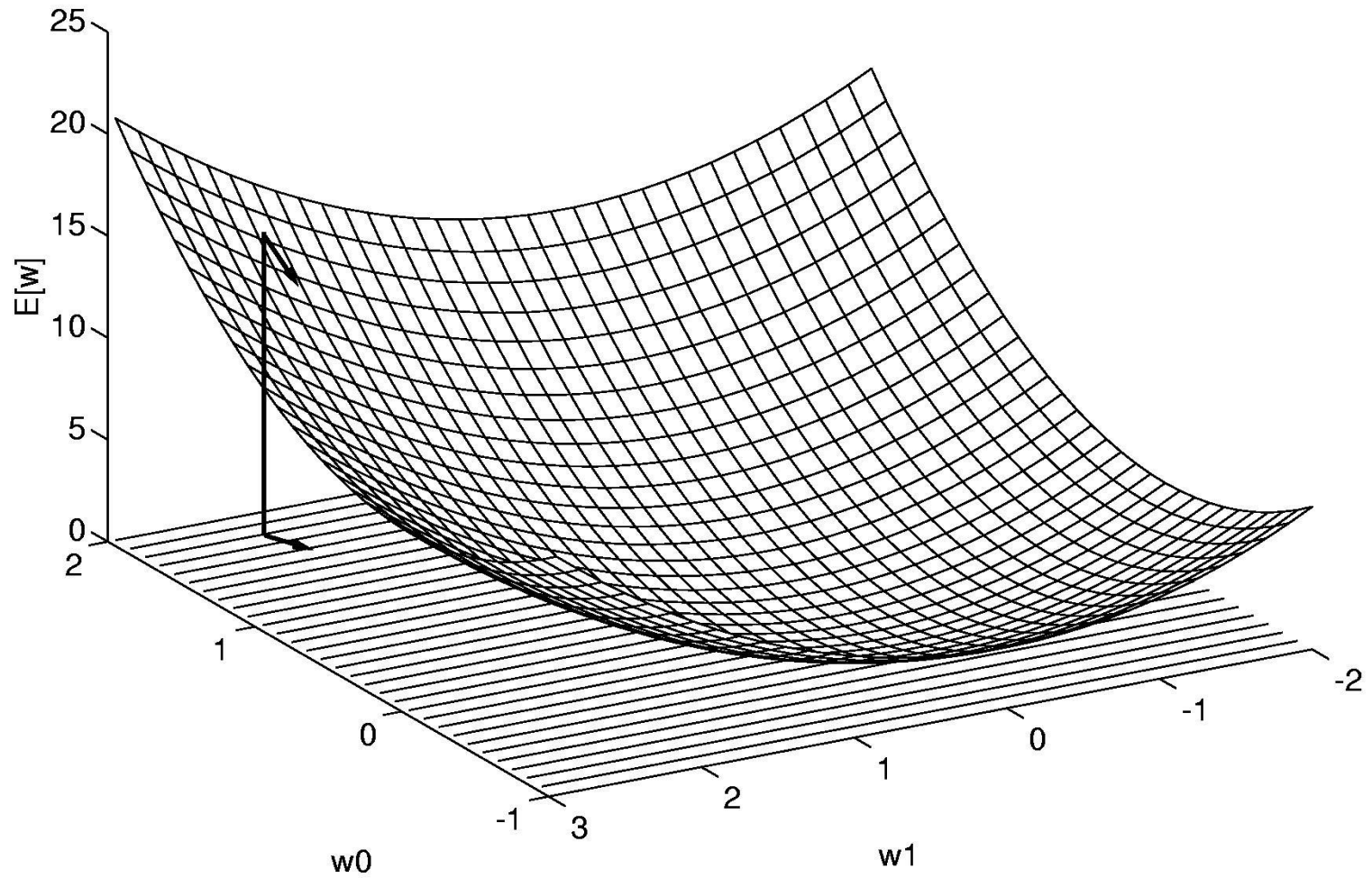
$$E[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

Where D is set of training examples

Note that the delta rule ($E[\vec{w}]$ given above) uses a different o_d as compared with the perceptron rule



Gradient Descent



Gradient:

$$\nabla E[\vec{w}] \equiv \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

Training rule:

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

I.e.:

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

Gradient Descent

$$\begin{aligned}\frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2 \\&= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\&= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\&= \sum_d (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d) \\ \frac{\partial E}{\partial w_i} &= \sum_d (t_d - o_d) (-x_{i,d})\end{aligned}$$

Gradient Descent

GRADIENT-DESCENT(*training_examples*, η)

Initialize each w_i to some small random value

Until the termination condition is met, Do

- Initialize each Δw_i to zero.
- For each $\langle \vec{x}, t \rangle$ in *training_examples*, Do
 - Input instance \vec{x} to unit and compute output o
 - For each linear unit weight w_i , Do

$$\Delta w_i \leftarrow \Delta w_i + \eta(t - o)x_i$$

- For each linear unit weight w_i , Do

$$w_i \leftarrow w_i + \Delta w_i$$

Summary

Perceptron training rule guaranteed to succeed if

- Training examples are linearly separable
- Sufficiently small learning rate η

Linear unit training rule uses gradient descent

- Guaranteed to converge to hypothesis with minimum squared error
- Given sufficiently small learning rate η
- Even when training data contains noise
- Even when training data not separable by H

Batch vs. Incremental Gradient Descent

Batch Mode Gradient Descent:

Do until convergence

1. Compute the gradient $\nabla E_D[\vec{w}]$
2. $\vec{w} \leftarrow \vec{w} - \eta \nabla E_D[\vec{w}]$

Incremental Mode Gradient Descent:

Do until convergence

For each training example d in D

1. Compute the gradient $\nabla E_d[\vec{w}]$
2. $\vec{w} \leftarrow \vec{w} - \eta \nabla E_d[\vec{w}]$

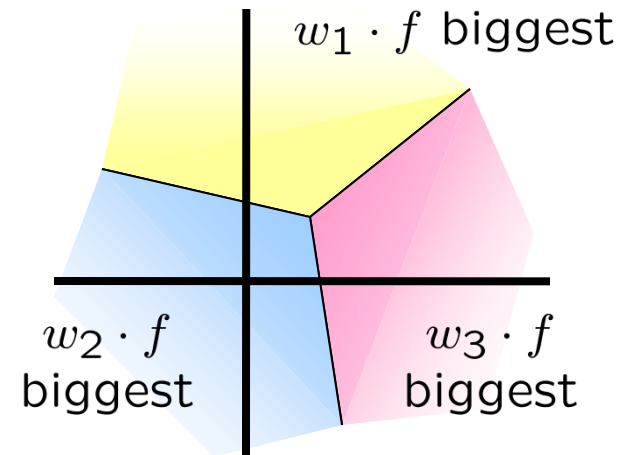
$$E_D[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$E_d[\vec{w}] \equiv \frac{1}{2} (t_d - o_d)^2$$

Incremental Gradient Descent can approximate *Batch Gradient Descent* arbitrarily closely if η made small enough

Multiclass Decision Rule

- If we have more than two classes:
 - Have a weight vector for each class: w_y
 - Calculate an activation for each class

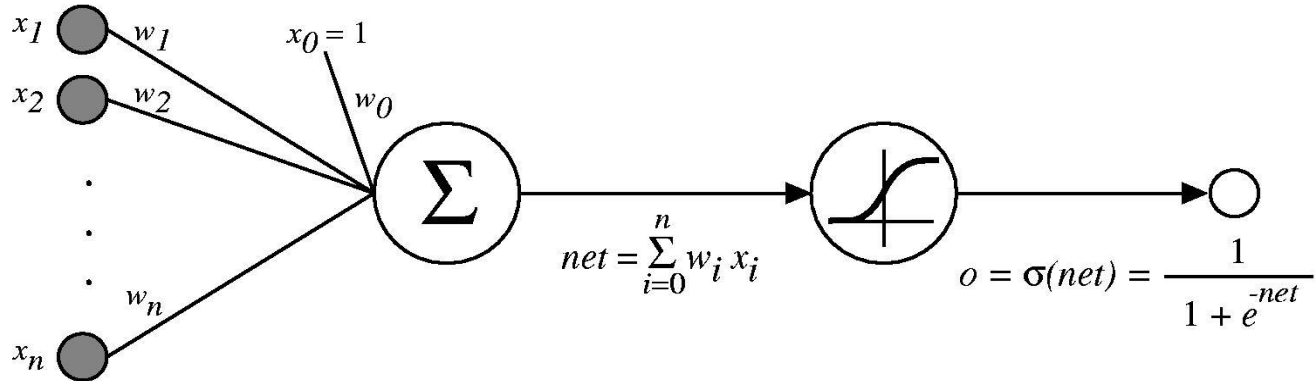


$$\text{activation}_w(x, y) = w_y \cdot f(x)$$

- Highest activation wins

$$y = \arg \max_y (\text{activation}_w(x, y))$$

Sigmoid Unit



$\sigma(x)$ is the sigmoid function

$$\frac{1}{1 + e^{-x}}$$

Nice property: $\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$

Error Gradient for a Sigmoid Unit

$$\begin{aligned}\frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum_d (t_d - o_d) \left(-\frac{\partial o_d}{\partial w_i} \right) \\ &= - \sum_d (t_d - o_d) \frac{\partial o_d}{\partial net_d} \frac{\partial net_d}{\partial w_i}\end{aligned}$$

But we know:

$$\frac{\partial o_d}{\partial net_d} = \frac{\partial \sigma(net_d)}{\partial net_d} = o_d(1 - o_d)$$

$$\frac{\partial net_d}{\partial w_i} = \frac{\partial(\vec{w} \cdot \vec{x}_d)}{\partial w_i} = x_{i,d}$$

So:

$$\frac{\partial E}{\partial w_i} = - \sum_{d \in D} (t_d - o_d) o_d (1 - o_d) x_{i,d}$$

Let: $\delta_k = - \frac{\partial E}{\partial net_k}$

Comparison: Linear vs Sigmoid Unit

- Training rule for linear unit: $o = \sum_{i=0}^n w_i x_i$

$$w_i = w_i + \Delta w_i$$

where

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} = \eta(t - o)x_i$$

- Training rule for Sigmoid unit $o = \text{Sig}(\sum_{i=0}^n w_i x_i)$

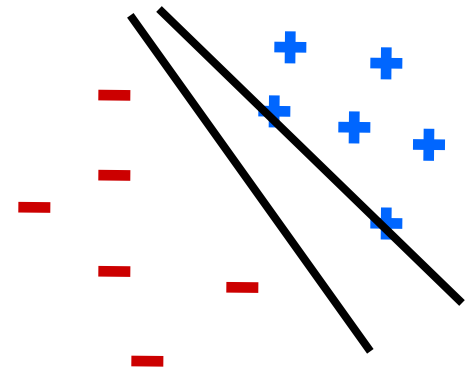
$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} = \eta(t - o)o(1 - o)x_i$$

- Training rule for unit of your choice!!! (e.g., *tanh*)
 - Same idea (set up the error function)
 - Use Gradient Descent (compute derivatives)

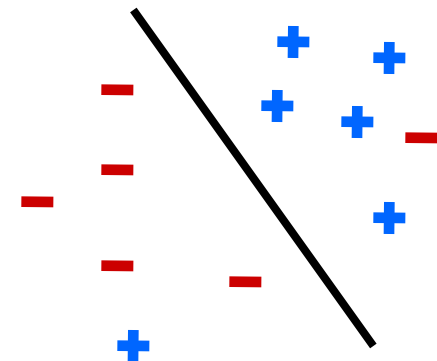
Properties of Perceptrons

- **Separability:** some parameters get the training set perfectly correct
- **Convergence:** if the training is separable, perceptron will eventually converge (binary case)

Separable

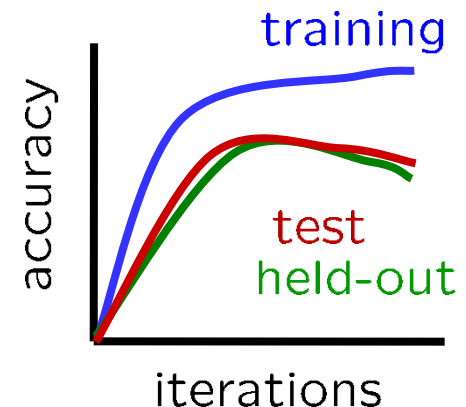
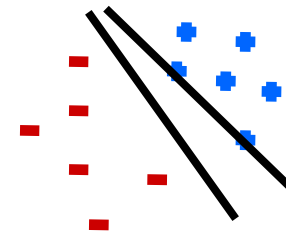
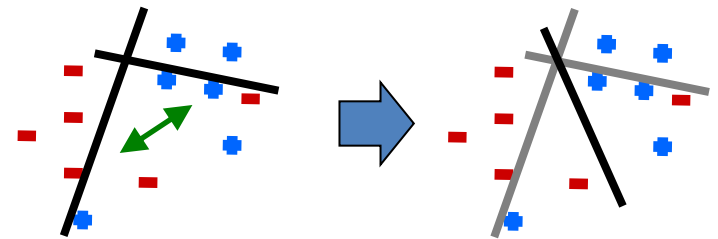


Non-Separable



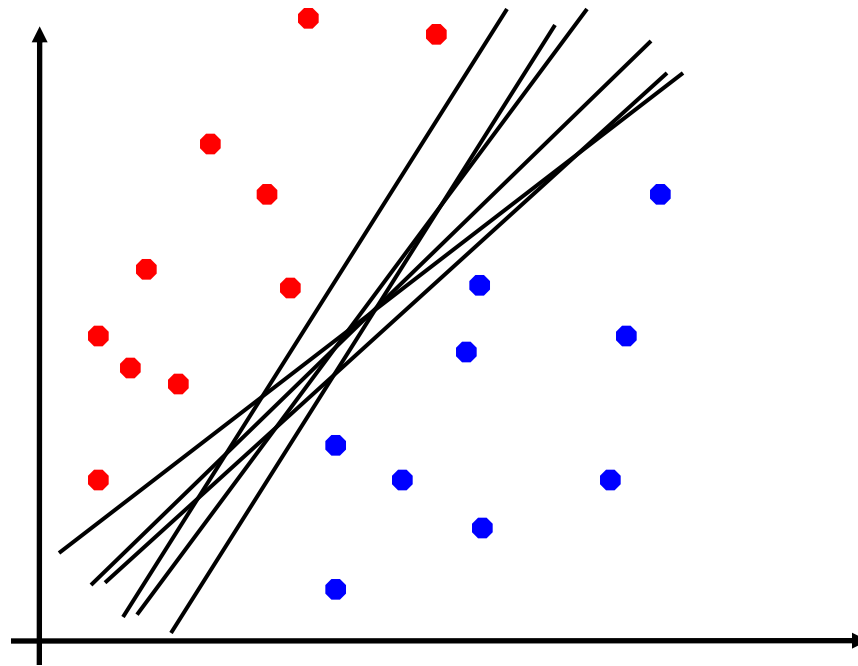
Problems with the Perceptron

- Noise: if the data isn't separable, weights might thrash
 - Averaging weight vectors over time can help (averaged perceptron)
- Mediocre generalization: finds a “barely” separating solution
- Overtraining: test / validation accuracy usually rises, then falls
 - Overtraining is a kind of overfitting



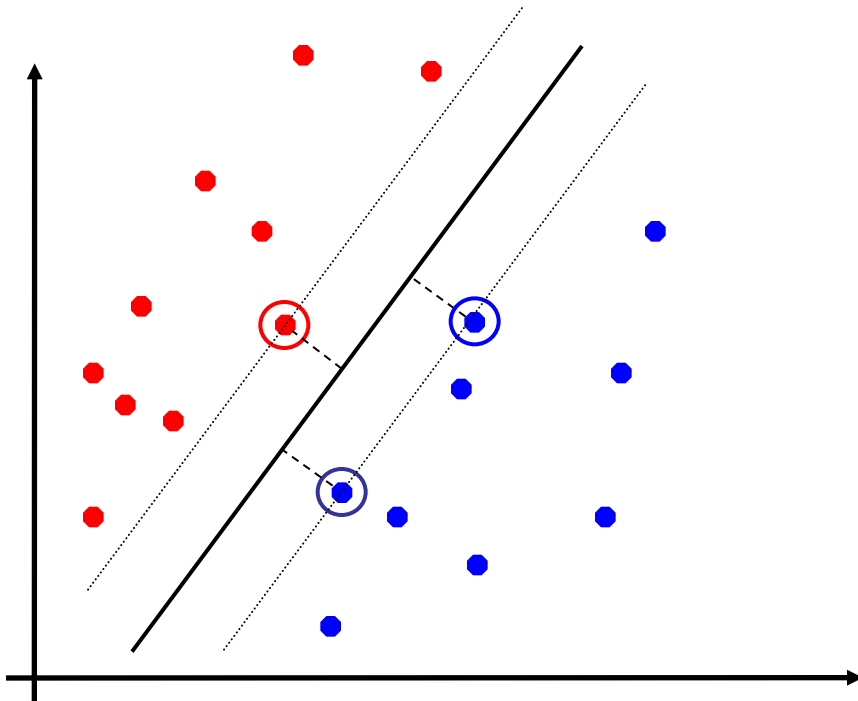
Linear Separators

- Which of these linear separators is optimal?



Support Vector Machines

- **Maximizing the margin:** good according to intuition, theory, practice
- SVMs find the separator with max margin

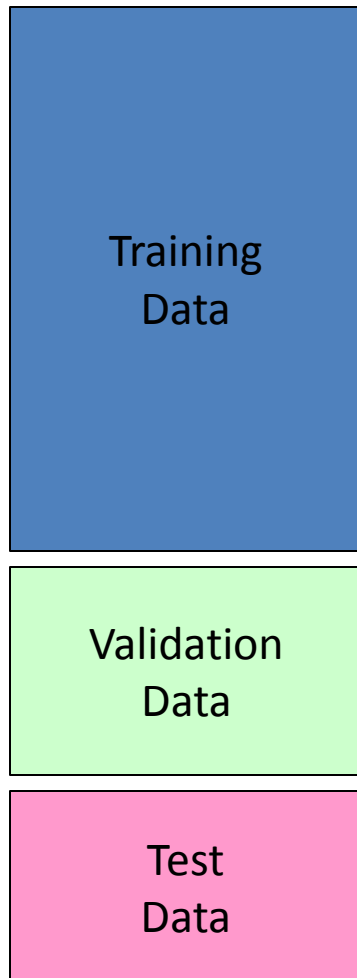


SVM

$$\min_w \frac{1}{2} \|w\|^2$$

$$\forall i, y \quad w_{y^*} \cdot f(x_i) \geq w_y \cdot f(x_i) + 1$$

Three Views of Classification



- Naïve Bayes:
 - Parameters from data statistics
 - Parameters: probabilistic interpretation
 - Training: one pass through the data
- Logistic Regression:
 - Parameters from gradient ascent
 - Parameters: linear, probabilistic model, and discriminative
 - Training: one pass through the data per gradient step; regularization essential
- The Perceptron:
 - Parameters from reactions to mistakes
 - Parameters: discriminative interpretation
 - Training: go through the data until accuracy on validation set maxes out