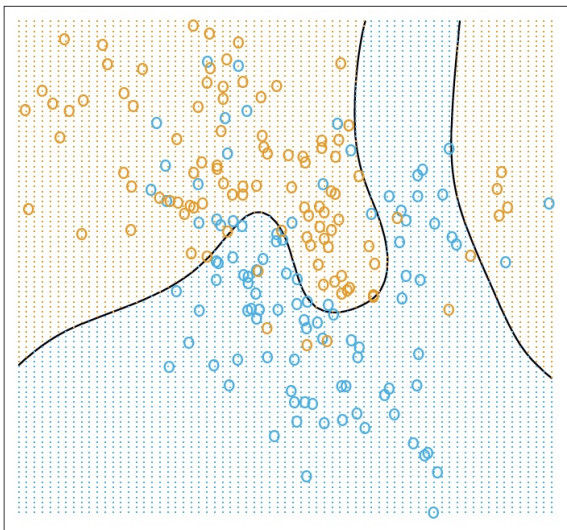


Topic 3: INTRO TO CLASSIFICATION

CMSC 35400/STAT 37710 Machine Learning
Risi Kondor, The University of Chicago



[Hastie, Tibshirani & Friedman]

Classification

Two-class classification in \mathbb{R}^n is the prototypical supervised ML problem.

- Input space: \mathcal{X} (in the simplest case, $\mathcal{X} = \mathbb{R}^n$)
- Output space: $\mathcal{Y} = \{-1, +1\}$ (in k -class case $\mathcal{Y} = \{1, \dots, k\}$)
- Training set: $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$
- Goal: find a good hypothesis $f: \mathcal{X} \rightarrow \mathcal{Y}$.

We will look at three very simple classifiers:

1. Model based classification
2. k -nearest neighbors
3. The perceptron

1. Model based classification

Model based classification

1. Assume that x and y come from a joint distribution $p(x, y)$ (similar to mixture models for clustering).
2. Define a marginal distribution for y (Bernoulli):

$$p_0(y = +1) = \pi \quad p_0(y = -1) = 1 - \pi.$$

3. Define the conditionals for x given y :

$$p_+(x) = p(x|y = +1)$$

$$p_-(x) = p(x|y = -1).$$

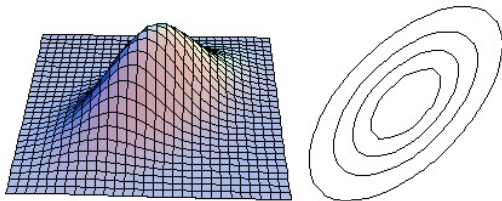
4. Estimate the parameters of p_+ , p_- and p_0 from the training data.
5. Classify future examples according to

$$\hat{y} = \begin{cases} +1 & \text{if } p(y = +1|x) \geq p(y = -1|x) \\ -1 & \text{if } p(y = +1|x) < p(y = -1|x). \end{cases}$$

This is called the **Bayes optimal classifier**.

The Gaussian model

Assume that $\mathcal{X} = \mathbb{R}^n$ and $p_+(\mathbf{x})$ and $p_-(\mathbf{x})$ are both multivariate normal (Gaussian) distributions:

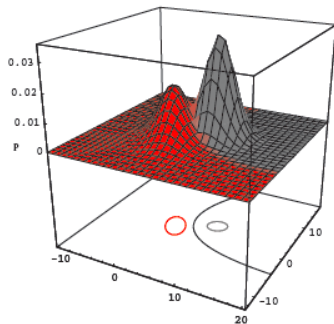
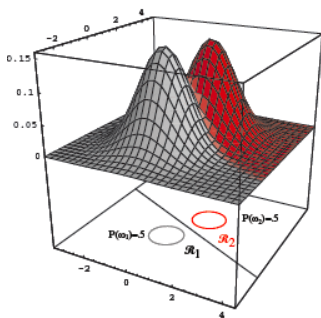


$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})/2}$$

where $\boldsymbol{\mu} \in \mathbb{R}^n$, and $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$ is a symmetric positive definite matrix.

$$\mathbb{E}(\mathbf{x}) = \boldsymbol{\mu} \quad \text{Cov}(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\Sigma}_{i,j}$$

The Gaussian model



How do we learn (estimate) the parameters μ_+ , Σ_+ , μ_- , Σ_- , π ? Maximum likelihood (easier than clustering because the y are not latent, so no need for EM).

MLE for two-Gaussians model

- MLE for π :

$$\hat{\pi} = \frac{n_+}{n} = \frac{\text{number of training examples with } y_i = 1}{\text{total number of training examples}}$$

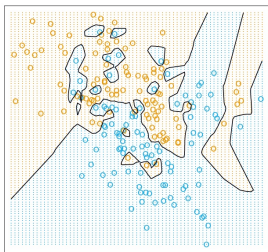
- MLE for μ and Σ :

$$\mu_+ = \frac{1}{n_+} \sum_{i: y_i=+1} \mathbf{x}_i \quad \Sigma_+ = \frac{1}{n_+} \sum_{i: y_i=+1} (\mathbf{x}_i - \mu_+)(\mathbf{x}_i - \mu_+)^{\top}$$

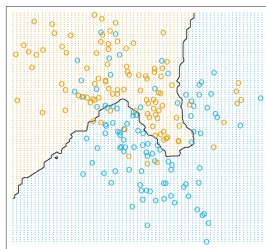
$$\mu_- = \frac{1}{n_-} \sum_{i: y_i=-1} \mathbf{x}_i \quad \Sigma_- = \frac{1}{n_-} \sum_{i: y_i=-1} (\mathbf{x}_i - \mu_-)(\mathbf{x}_i - \mu_-)^{\top}$$

2. k -nearest neighbors

k -nearest neighbors



1-nearest neighbor



15-nearest neighbor

The hypothesis returned by the k -NN classifier is

$$\hat{f}(\mathbf{x}) = \text{sgn}\left(\sum_{i \in \text{NN}_k(\mathbf{x})} y_i\right),$$

where $\text{NN}_k(\mathbf{x}) = \{i_1, \dots, i_k\}$ are the indices of the k points in S closest to \mathbf{x} . (When k is even and $\sum_{i \in \text{NN}_k(\mathbf{x})} y_i = 0$, break ties arbitrarily.)

What effect does k have on the behavior of k -NN?

Model based vs. k -NN classification

The Two-Gaussians classifier is a **generative** learning algorithm:

- Given $\mu_+, \Sigma_+, \mu_-, \Sigma_-, \pi$ we could generate more data.
- Nice, but from an ML point of view, often an overkill.

k -NN is a **discriminative** learning algorithm:

- It just focuses on learning the hypothesis $f: \mathbf{x} \mapsto \{+1, -1\}$. Doesn't try to learn the underlying p .
- This is more in the ML spirit. However some branches of ML (e.g., graphical models) are fundamentally generative.

3. The Perceptron

Linear classifiers

The Perceptron belongs to the wide class of linear classifiers.

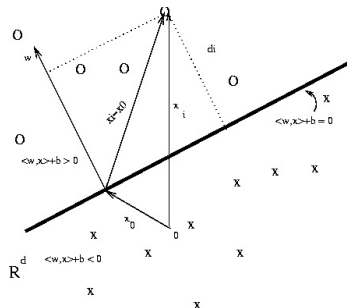
Input space: $\mathcal{X} = \mathbb{R}^n$

Output space: $\mathcal{Y} = \{-1, +1\}$

Hypothesis (affine hyperplane):

$$f(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$$

- \mathbf{w} is the normal to the **separating hyperplane**
- b is the **bias**

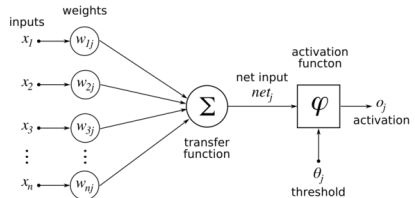
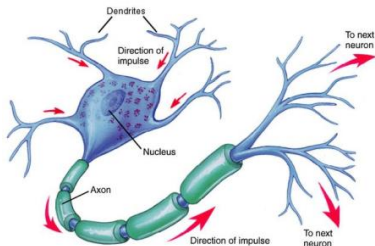


$$\text{Empirical error : } \mathcal{E}_{\text{train}}(f) = \frac{1}{m} \sum_{i=1}^m \ell_{0/1}(f(\mathbf{x}_i), y_i)$$

where $\ell_{0/1}(\hat{y}, y) = 0$ if $\hat{y} = y$, else $\ell_{0/1}(\hat{y}, y) = 1$.

Of all possible hyperplanes that separate the data which one is best?

Inspiration: Artificial Neural Networks



$$o_j = \varphi\left(\theta_j + \sum_i w_{i,j} x_i\right),$$

An old idea going back to [McCulloch & Pitts, 1943].

The perceptron (Rosenblatt, 1958)

Consider the following a simplified scenario:

- $\|\mathbf{x}_i\| = 1$ for all i .
- Assume that all datapoints satisfy the ground truth concept

$$h_{\text{true}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{v} \cdot \mathbf{x} \geq 0 \\ -1 & \text{otherwise} \end{cases},$$

where \mathbf{v} is some fixed vector and w.l.o.g. $\|\mathbf{v}\| = 1$.

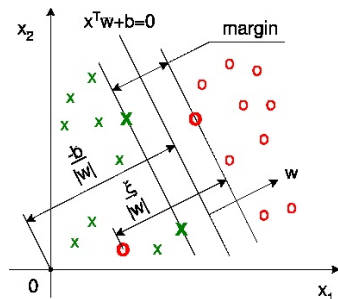
- Since h_{true} goes through the origin, set $b = 0$ for all h .
- We want to solve the problem in its online form.

The perceptron (Rosenblatt, 1958)

```
 $\mathbf{w} \leftarrow 0 ;$   
 $t \leftarrow 1 ;$   
while(true){  
    if  $\mathbf{w} \cdot \mathbf{x}_t \geq 0$  predict  $\hat{y}_t = 1$  ; else predict  $\hat{y}_t = -1$  ;  
    if  $((\hat{y}_t = -1) \text{ and } (y_t = 1))$  let  $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}_t$  ;  
    if  $((\hat{y}_t = 1) \text{ and } (y_t = -1))$  let  $\mathbf{w} \leftarrow \mathbf{w} - \mathbf{x}_t$  ;  
     $t \leftarrow t + 1$  ;  
}
```

Note: In this form, the Perceptron is an **online** learning algorithm.

The margin



The projection of \mathbf{x} onto the line orthogonal to decision boundary $\mathbf{w} \cdot \mathbf{x} + b = 0$ is

$$P\mathbf{x} = \mathbf{w} \cdot \mathbf{x} / \|\mathbf{w}\|.$$

On the boundary

$$\mathbf{w} \cdot \mathbf{x} = -b \quad \Rightarrow \quad P\mathbf{x} = \frac{-b}{\|\mathbf{w}\|}.$$

The margin of a correctly classified positive resp. negative example are

$$\delta_{(\mathbf{x}, +1)} = P\mathbf{x} - \frac{-b}{\|\mathbf{w}\|}, \quad \delta_{(\mathbf{x}, -1)} = \frac{-b}{\|\mathbf{w}\|} - P\mathbf{x},$$

so the general formula is $\delta_{(\mathbf{x}, y)} = y (\mathbf{w} \cdot \mathbf{x} + b) / \|\mathbf{w}\|$.

The perceptron convergence thm

Suppose that the perceptron is run on a sequence of examples with $\|\mathbf{x}_i\| = 1$ and $y_i(\mathbf{v} \cdot \mathbf{x}_i) \geq 0$ for some \mathbf{v} with $\|\mathbf{v}\| = 1$. In addition, assume that the data obeys $|\mathbf{v} \cdot \mathbf{x}_i| \geq \delta$ for all \mathbf{x}_t (**margin**).

Claim 1: After M mistakes $\mathbf{w} \cdot \mathbf{v} \geq \delta M$.

Initially, $\mathbf{w} \cdot \mathbf{v} = 0$. After a false negative, $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}_t$, and $(\mathbf{w} + \mathbf{x}_t) \cdot \mathbf{v} \geq (\mathbf{w} \cdot \mathbf{v}) + \delta$, so $\mathbf{w} \cdot \mathbf{v}$ increases by at least δ . Similarly for a false positive.

The perceptron convergence thm

Claim 2: After M mistakes, $\|\mathbf{w}\|^2 \leq M$

Initially $\|\mathbf{w}\| = 0$. On a false negative, $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}_t$, and

$$\|\mathbf{w} + \mathbf{x}\|^2 = (\mathbf{w} + \mathbf{x}_t) \cdot (\mathbf{w} + \mathbf{x}_t) = \|\mathbf{w}\|^2 + 2\mathbf{w} \cdot \mathbf{x}_t + \|\mathbf{x}_t\|^2.$$

The third term on the RHS is 1, while the second term must be negative, since we predicted 0. Therefore, $\|\mathbf{w}\|^2$ can increase by at most 1. Similarly for false positives.

The perceptron convergence thm

Theorem. If the perceptron is run a sequence of examples with $|\mathbf{v} \cdot \mathbf{x}_t| \geq \delta$ (plus our other assumptions), then it will make at most $1/\delta^2$ mistakes.

Proof.

$$\delta M \leq \mathbf{w} \cdot \mathbf{v} \leq \|\mathbf{w}\| \leq \sqrt{M}$$



Affine hyperplanes

Question: What if we don't want the separator to go through the origin?

Add extra feature with $[\mathbf{x}_t]_{n+1} = 1$ for all t .

$$x_1w_1 + x_2w_2 + \dots x_nw_n \geq \theta$$

$$\Longleftrightarrow$$

$$x_1w_1 + x_2w_2 + \dots x_nw_n + x_{n+1}w_{n+1} \geq 0$$

with $w_{n+1} = -\theta$

Question: How can we make linear threshold functions more expressive?

Towards multiclass

Rewrite perceptron as competition between \mathbf{w}^+ and \mathbf{w}^-

```
 $\mathbf{w}^+ \leftarrow (0, 0, \dots, 0) ;$   
 $\mathbf{w}^- \leftarrow (0, 0, \dots, 0) ;$   
 $t \leftarrow 1 ;$   
while(1){  
  if  $\mathbf{w}^+ \cdot \mathbf{x}_t \geq \mathbf{w}^- \cdot \mathbf{x}_t$  predict  $\hat{y}_t = 1$  ; else predict  $\hat{y}_t = 0$  ;  
  if  $((\hat{y}_t = 0) \text{ and } (y_t = 1))$  let  $\mathbf{w}^+ \leftarrow \mathbf{w}^+ + \mathbf{x}_t$  ;  
  if  $((\hat{y}_t = 1) \text{ and } (y_t = 0))$  let  $\mathbf{w}^- \leftarrow \mathbf{w}^- + \mathbf{x}_t$  ;  
   $t \leftarrow t + 1$  ;  
}
```

Towards multiclass

If $y_t \in \{1, 2, \dots, \kappa\}$, maintain κ different weight vectors $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_\kappa$.

On erroneously classifying example \mathbf{x}_t as class i instead of class j :

- $\mathbf{w}_i \leftarrow \mathbf{w}_i - \mathbf{x}_t/2$
- $\mathbf{w}_j \leftarrow \mathbf{w}_j + \mathbf{x}_t/2$

[Collins & Duffy, 2002]

General concepts of supervised learning

The general framework

- There is an *unknown* distribution p on $\mathcal{X} \times \mathcal{Y}$.
- The training set $S = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$, is an IID sample from p .
- The hypothesis is a function $\hat{f}: \mathcal{X} \mapsto \mathcal{Y}$.
- The learning algorithm is a function $\mathcal{A}: S \mapsto \hat{f}$.
- The hypothesis space \mathcal{F} is the space of functions accessible to \mathcal{A} .

We want f to be good on future examples drawn from p , where “good” is defined in terms of a loss function $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$.

How can we ever be sure that we can find a good f , especially when \mathcal{F} is huge??? This is what **Statistical Learning Theory** is about.

Three notions of error

- The **training error** or **empirical error**:

$$\mathcal{E}_{\text{emp}}[\hat{f}] = \frac{1}{m} \sum_{i=1}^m \ell(\hat{f}(\mathbf{x}_i), y_i)$$

This is what we can actually measure, and what many learning algorithms try to minimize. → **Empirical Risk Minimization (ERM)**

- The **testing error**:

$$\mathcal{E}_{\text{test}}[\hat{f}] = \frac{1}{m'} \sum_{i=1}^{m'} \ell(\hat{f}(\mathbf{x}'_i), y'_i)$$

This is what the algorithm is going to be evaluated on.

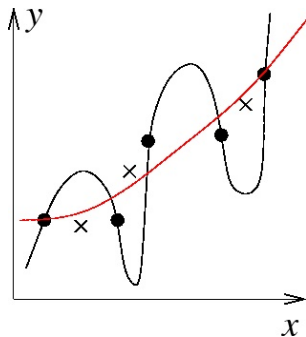
- The **true error**:

$$\mathcal{E}_{\text{true}}[\hat{f}] = \mathbb{E}_{(x,y) \sim p} \ell(\hat{f}(\mathbf{x}), y).$$

This is what an ideal algorithm would minimize, but it can't be measured since p is unknown.

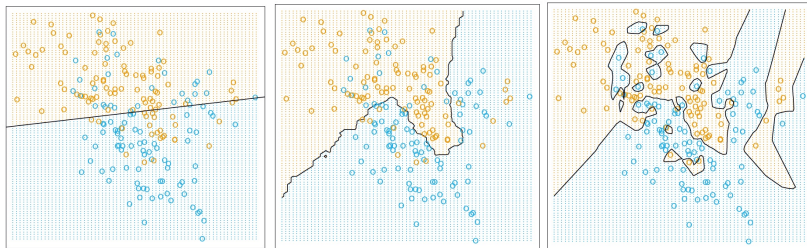
Overfitting

In general, because \hat{f} is explicitly chosen to minimize $\mathcal{E}_{\text{emp}}[\hat{f}]$, it will tend to be an overoptimistic estimate of the true error, i.e., $\mathcal{E}_{\text{true}}[\hat{f}] > \mathcal{E}_{\text{train}}[f]$ and $\mathcal{E}_{\text{test}}[\hat{f}] > \mathcal{E}_{\text{train}}[f]$. This is called **overfitting**.



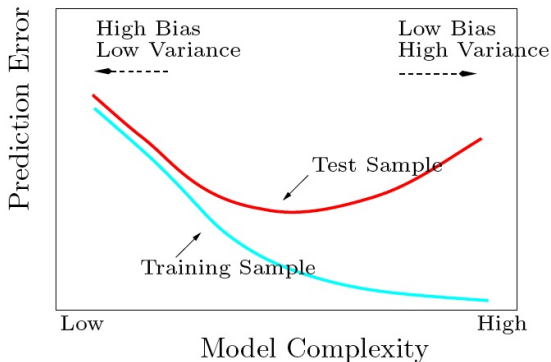
Question: How can we avoid overfitting? Restrict hypothesis space!! But how much?

Overfitting and underfitting



This is the biggest thing to be paranoid about in *all* branches of ML.

The overfitting/underfitting tradeoff



The hypothesis space should be big enough but not too big. Alternatively, we need to add a regularizer (c.f. inverse problems in applied math).

Regularized Risk Minimization (RRM)

The most general framework for finding the right balance between overfitting and underfitting in supervised learning is RRM:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \left[\underbrace{\frac{1}{m} \sum_{i=1}^m \ell(f(x_i), y_i)}_{\text{training error}} + \underbrace{\lambda \Omega(f)}_{\text{regularizer}} \right]$$

Terminology:

1. \mathcal{F} : hypothesis space
2. $\ell(\hat{y}, y)$: loss function
3. $\Omega: \mathcal{F} \rightarrow \mathbb{R}^+$: regularization functional

The right compromise is found by tuning the regularization parameter λ .
[Tykhonov regularization]

Statistical Learning Theory

Concentration of measure

The reason that we have any hope of learning a good \hat{f} at all is that for any **fixed** $f \in \mathcal{F}$, over the choice of training/testing sets,

$$\mathbb{E}(\mathcal{E}_{\text{emp}}[f]) = \mathbb{E}(\mathcal{E}_{\text{test}}[f]) = \mathcal{E}_{\text{true}}[f].$$

Moreover, as the size of training/testing sets increases, these quantities are more and more **concentrated** around $\mathcal{E}_{\text{true}}[f]$

However, a very unlucky choice of training set can always mess us up.

Generalization Bounds

Statistical Learning Theory proves bounds of the form

$$\mathbb{P} [\mathcal{E}_{\text{true}}[\hat{f}] > \mathcal{E}_{\text{emp}}[\hat{f}] + \epsilon] < \delta$$

where ϵ is a complicated function depending on δ , the size of the function class \mathcal{F} , and the nature of the regularization. \rightarrow **Probably Approximately Correct (PAC)** bounds [Valiant, 1984]

Example: if the Vapnik-Chervonenkis dimension of \mathcal{F} is d , then

$$\mathbb{P} \left[\mathcal{E}_{\text{true}}[\hat{f}] > \mathcal{E}_{\text{emp}}[\hat{f}] + \sqrt{\frac{d(\log \frac{2m}{d} + 1) + \log(4/\delta)}{m}} \right] \leq \delta.$$

In reality even the best such bounds are ridiculously loose.

Cross validation

Holdout set

The generalization bounds from statistical learning theory are framed entirely in terms of the training set, therefore they reveal fundamental properties of the learning algorithm. Unfortunately, they are almost always incredibly loose.

The more practical way to estimate $\mathcal{E}_{\text{true}}$ is to split the training data into:

- The true training set, used to train the algorithm
- A **holdout set** $\{(x_1, y_1), (x_2, y_2), \dots, (x_p, y_p)\}$ which is only used to compute the holdout error

$$\mathcal{E}_{\text{h.o.}}[\hat{f}] = \frac{1}{p} \sum_{i=1}^p \ell(\hat{f}(x_i), y_i)$$

Holdout error is an unbiased estimator of $\mathcal{E}_{\text{true}}[\hat{f}]$, i.e., $\mathbb{E}[\mathcal{E}_{\text{h.o.}}[\hat{f}]] = \mathcal{E}_{\text{true}}[\hat{f}]$.

Cross Validation

k -fold cross validation uses the holdout idea to set internal parameters θ of learning algorithms (e.g., k in k -NN):

- Set θ to some value.
- Split the training set into k roughly equal parts S_1, S_2, \dots, S_k .
- For each i , train the algorithm on $S_1, \dots, S_{i-1}, S_{i+1}, \dots, S_p$ to get \hat{f}_i .
- Compute the holdout error of each \hat{f}_i on the corresponding S_i .
- Average to get an estimator of $\mathcal{E}_{\text{true}}(\hat{f}_\theta)$.
- Iterate over a range of θ values and finally set θ to the value that minimizes the cross-validation error.