

# Lecture 06 - Support Vector Machines

NIPS

j = 677

## (1) Naive Bayes

- Laplace Smoothing
- Event Models

## (2) Comments on apply ML

## (3) SVM intros

### (1) Recap

$$x = \begin{bmatrix} 1 \\ 0 \\ \dots \\ 1 \\ 0 \\ \dots \end{bmatrix} \begin{array}{l} \text{a} \\ \text{adv} \\ \dots \\ \text{buy} \\ : \end{array}$$

$x_j = 1$  { whether j appears in email }

→ Generative Model (Gaussian or Bernoulli)

$$p(x|y) = \prod_{j=1}^n p(x_j|y) p(y)$$

i-th → training examples

j-th → index in features

→ Maximum Likelihood

$$\rightarrow \text{parameters } p(y=1) = \phi_y$$

$$p(x_j = 1 | y=0) = \phi_j | y=0$$

$$p(x_j = 1 | y=1) = \phi_j | y=1$$

At prediction time:

$$P(y=1|x) = \frac{P(x|y=1) \cdot P(y=1)}{P(x|y=1) \cdot P(y=1) + P(x|y=0) \cdot P(y=0)}$$

→ Laplace smoothing

More Generally:

$$x \in \{1, \dots, k\}$$

$$\Rightarrow \text{Estimate } p(x=j) = \frac{\sum_{i=1}^n 1\{x_i^{(i)}=j\} + 1}{M+k}$$

$$\phi_j | y=0 = \frac{\sum_{i=1}^n 1\{x_i^{(i)}=j, y^{(i)}=0\} + 1}{M+k}$$

$$\frac{\sum_{i=1}^n 1\{y^{(i)}=0\} + 2}{M+k}$$

$$x_i \in \{1, \dots, k\}$$

Size	< 400	400 - 800	800 - 1200	> 1200
$x_i$	1	2	3	4

$$P(x|y) = \prod_{j=1}^n P(x_j|y) \xrightarrow{\text{multinomial}}$$

- New representation

$$x \in \begin{bmatrix} 1600 \\ 800 \\ 1600 \\ 6200 \end{bmatrix} \in \mathbb{R}^w$$

$$x = \begin{bmatrix} 0 \\ 1 \\ i \\ 1 \\ 1 \end{bmatrix} \begin{array}{l} \xrightarrow{a} \\ \xrightarrow{\text{ad...}} \\ \xrightarrow{\text{buy}} \\ \xrightarrow{\text{atm}} \\ \xrightarrow{\text{new}} \end{array} \begin{array}{l} 1 \\ c \\ 800 \\ 1600 \\ 6200 \end{array}$$

$$x_j \in \{1 \dots 1000\}$$

$R_i$  = length of email  $i$   $\hookrightarrow$  Multinomial event model {new represent'}

$$P(x, y) = P(x|y) \cdot P(y)$$

assume

$$\prod_{j=1}^R P(x_j|y) \cdot P(y)$$

• Parameters:  $\phi_y = P(y=1)$

$$\phi_{k|y=0} = P(x_j=k|y=0)$$

↳ chance of word  $j$  being  $k$   
if  $y = 0$

"doesn't depend on  $j$ th"

- MLE:

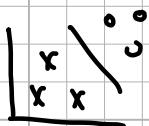
$$\phi_{k|y=0} = \frac{\sum_{i=1}^n 1 \{ y^{(i)} = 0 \} \sum_{j=0}^{w_i} 1 \{ x_j^{(i)} = k \} + 1}{\sum_{i=1}^n 1 \{ y^{(i)} = 0 \} w_i + 1000}$$

word embedding  
 $\hookrightarrow C(230)$

→ Support Vector Machine

- non-linear boundaries

- optimal margin classifier (separable case)



- kernels

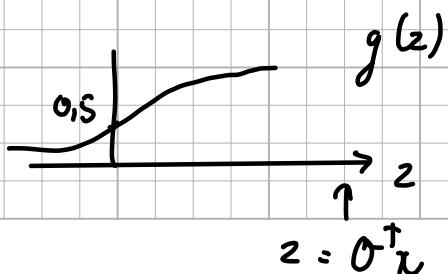
$$x \rightarrow \phi(x) \rightarrow \mathbb{R}^m$$

$$x \rightarrow \phi(x)$$

- inseparable case

functional margin  $\rightarrow$  how confidently

$$h_\phi(x) = g(\phi^\top x)$$



Predict "1" if  $\Omega^T x > 0$  ( $h_\Omega(x) = g(\Omega^T x) \geq 0.5$ )

"0" otherwise

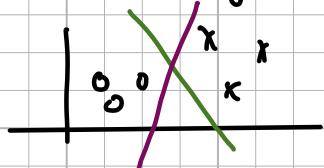
if  $y^{(i)} = 1$ , hope that  $\Omega^T x^{(i)} > 0$

if  $y^{(i)} = 0$ , hope that  $\Omega^T x^{(i)} \ll 0$

Geometric margin

$\rightarrow$  optimal margin classifier

to max geometric margin



Notation:

Labels  $y \in \{-1, +1\}$

$\rightarrow$  have  $h$  output  $\{-1, +1\}$

$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

$$\left[ \begin{array}{c} o_0 \\ o_1 \\ o_2 \\ \vdots \end{array} \right] \left\{ \begin{array}{l} b \\ w \end{array} \right.$$

$$h_\Omega(x) = g(\Omega^T x)$$

$\uparrow R^{n+1}, x_0 = 1$

$$h_{w,b}(x) = g(w^T x + b)$$

$\uparrow R^n \quad \uparrow R$

Drop  $x_0 = 1$   
Convention

functional margin of hyperplane defined by  $(w, b)$  wrt  $(x^{(i)}, y^{(i)})$

$$\hat{y}^{(i)} = y^{(i)} (w^T x^{(i)} + b) \rightarrow \text{very large} \gg 0$$

if  $y^{(i)} = 1$  want  $w^T x^{(i)} + b \geq 0$

if  $y^{(i)} = -1$  want  $w^T x^{(i)} + b \ll 0$

$\rightarrow$  if  $\hat{y}^{(i)} \gg 0$  that means  $h_i(x^{(i)}) = \hat{y}^{(i)}$

$$\text{and } \hat{y} = \min_i \hat{y}^{(i)}$$

$$\Rightarrow f^{(i)} = \frac{(w^T x^{(i)} + b) \cdot y^{(i)}}{\|w\|}$$

$$\text{geometric margin} = \frac{\hat{f}^{(i)}}{\|w\|} \quad \text{functional margin}$$

$$\|w\|=1 \quad , \quad (w, b) \rightarrow \left( \frac{w}{\|w\|}, \frac{b}{\|w\|} \right) \quad \text{Normalize}$$

$\rightarrow$  optimal margin classifier

$$\left\{ \begin{array}{l} \text{Choose } w, b \text{ to max } g \\ \max_{f, w, b} f \text{ s.t. } \frac{y^{(i)}(w^T x^{(i)} + b)}{\|w\|} \geq 1 \quad \forall i = 1 \dots n \end{array} \right.$$

$$\rightarrow \min_{w, b} \|w\|^2$$

$$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1$$

## Lecture 07

### SVM

- Optimization problem
- Representer theorem
- Kernels
- Examples of Kernels

#### (i) Optimal margin classifier

$$\min_{w, b} \frac{1}{2} \|w\|^2 \quad \text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1$$

$$g^{(i)} = \frac{y^{(i)}(w^T x^{(i)} + b)}{\|w\|} \quad (\text{geometric margin})$$

$$g = \min_{i=1 \dots n} g^{(i)}$$

$$y = g(w^T x + b)$$

$$x^{(i)} \in \mathbb{R}^{100}$$

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

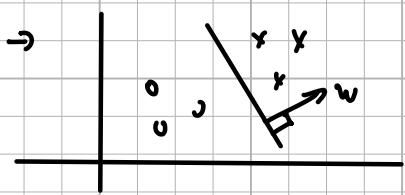
$$\begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_m \end{bmatrix} \rightarrow \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_m \end{bmatrix}$$

(Intuition #1)

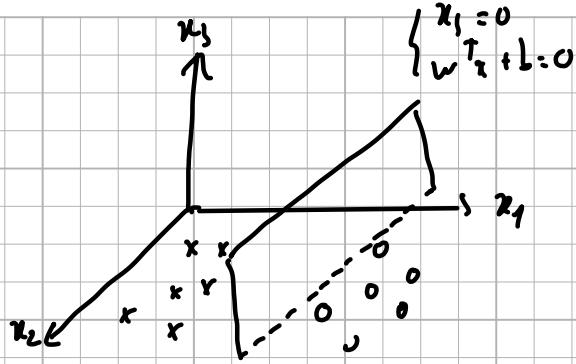
$\rightarrow$  Gradient Descent:  $\alpha_i := 0$

$$\alpha := \alpha - \alpha \sum_{i=1}^m (h_\alpha(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

(Intuition #2)



$$g(w^T u + b)$$



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

$$\text{Min}_{w,b} \frac{1}{2} \|w\|^2 \rightarrow w^T w$$

$$\text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq 1 \quad i=1 \dots m$$

$\langle x, z \rangle = z^T x$  is the inner products

$$\rightarrow \text{min}_{w,b} \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} \underbrace{x^{(i)} T x^{(j)}}_{\langle x^{(i)}, x^{(j)} \rangle}$$

$$\text{s.t. } y^{(i)} \left( \left( \sum_j \alpha_j y^{(j)} \underbrace{x^{(j)}}_{\langle x^{(j)}, x^{(i)} \rangle} \right)^T x^{(i)} + b \right) \geq 1$$

"Dual optimization prob"

$$h_{w,b}(u) = g(w^T u + b)$$

$$= g\left(\sum_i \alpha_i y^{(i)} \langle x^{(i)}, u \rangle + b\right)$$

④ Kernel trick

$$\begin{matrix} x & z \\ \downarrow & \downarrow \\ x^{(i)} & x^{(j)} \end{matrix}$$

1) Write algo in terms of  $\langle x^{(i)}, x^{(j)} \rangle$

2) Let there be mapping  $u \mapsto \phi(u)$

3) Find way to compute  $K(x, z) = \phi(x)^T \phi(z)$

4) Replace  $\langle x, z \rangle$  in algo with  $K(x, z)$

$$\begin{cases} x_1 = 0 \\ w_1^T x_i + b = 0 \end{cases}$$

$$x \in \mathbb{R}^n \rightarrow \phi(x) \in \mathbb{R}^{n^2}$$

$$K(x, z) : \underbrace{\phi(x)^T \phi(z)}_{\text{product}} = \underbrace{(x^T z)^2}_{\frac{1}{2} \|x\|^2} = \left( \sum_{i=1}^n x_i z_i \right) \left( \sum_{j=1}^n x_j z_j \right)$$

$$= \sum_{i=1}^n \sum_{j=1}^n (x_i z_j)(x_j z_i) \hookrightarrow \mathbb{R}^n$$

$$K(x, z) = (x^T z + c)^d$$

$$\hookrightarrow x = \begin{bmatrix} x_1 x_1 \\ \vdots \dots \\ \vdots \dots \\ x_1 x_n \\ \sqrt{2c} x_1 \\ \sqrt{2c} x_2 \\ \sqrt{2c} x_3 \\ c \end{bmatrix} \quad z = \begin{bmatrix} z_1 z_1 \\ \vdots \dots \\ \vdots \dots \\ z_1 z_n \\ \sqrt{2c} z_1 \\ \sqrt{2c} z_2 \\ \sqrt{2c} z_3 \\ c \end{bmatrix}$$

$$K(x, z) = (x^T z + c)^d \rightarrow O(n) \text{ time}$$

$\phi^d$  has all  $\binom{n+d}{d}$  feature of monomial up to order  $d$

$\rightarrow$  SVM = optimal margin classifier + Kernel trick

- how to make kernels

if  $x, z$  are "similar"  $K(x, z) = \phi(x)^T \phi(z)$  is large and otherwise

$$K(x, z) = \exp \left( - \frac{\|x - z\|^2}{2\sigma^2} \right) \rightarrow \text{Gaussian Kernel } \phi(x) \in \mathbb{R}^{\infty}$$

$\hookrightarrow$  When to apply?

Does there exist  $\phi$  s.t  $K(x, z) = \phi(x)^T \phi(z)$

$$\text{and } K(x, x) = \phi(x)^T \phi(x) \geq 0$$

Theorem (Mercer) :  $K$  is a valid Kernel func (i.e  $\exists \phi$  s.t  $K(x, z) = \phi(x)^T \phi(z)$ ) if and only if for any  $d$  points, the corresponding matrix  $K \geq 0$

Let  $\{x^{(1)}, \dots, x^{(n)}\}$  be d points  
 $K \in \mathbb{R}^{d \times d}$  "Kernel matrix"

$$K_{ij} = K(x^{(i)}, x^{(j)})$$

→ Linear Kernel  $K(x_i, x_j) = x_i^T x_j$ ,  $f(x) = x$

→ L1 norm Soft margin from

$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \quad \text{where } \xi_i \geq 0$$

$$s.t. y^{(i)} (w^T x^{(i)} + b) \geq 1 - \xi_i$$

→ Positive sequence classifier

$$\alpha_i \text{ in Lagrangian}$$

$$\max_w w(w) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$s.t. 0 \leq \alpha_i \leq C, i=1 \dots n$$

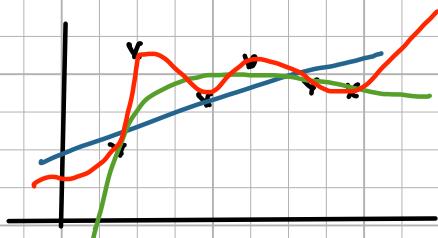
$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

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

## Lecture 08 Learning Theory

- Bias / Variance
- Regularization
- Train / dev (test splits)
- Model selection / Cross-Validation

### ① Bias / Variance



- underfit  $\rightarrow$  high bias
  - "just right"
  - overfit  $\rightarrow$  high variance
- ↑ don't learn  
↑ prejudice  
↓ totally different res

## (2) Regularization

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^m \|y^{(i)} - \theta^T x^{(i)}\|^2 + \underbrace{\frac{\lambda}{2} \|\theta\|^2}_{\text{regularization term}}$$

$$\arg \max_{\theta} \sum_{i=1}^n \log p(y^{(i)} | x^{(i)}; \theta) - \lambda \|\theta\|^2 \approx \|\theta\|^2$$

### • Text classifier

$$m = 100 \quad x = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \text{ad.}$$

$$n = 10,000$$

$$S = \{(x^{(i)}, y^{(i)})\}_{i=1}^m$$

$$P(\theta|S) = \frac{P(S|\theta) P(\theta)}{P(S)}$$

$$\arg \max_{\theta} P(\theta|S) = \arg \max_{\theta} P(S|\theta) P(\theta)$$

$$\text{Assume: } \theta \sim N(0, \tau^2 I)$$

$$= \arg \max_{\theta} \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta) \cdot P(\theta)$$

$$P(\theta) = \frac{1}{\sqrt{2\pi} (\tau^2 I)^{1/2}} \exp\left(-\frac{1}{2} \theta^T (\tau^2 I)^{-1} \theta\right)$$

Logistic regression

Frequentist  $\rightarrow \arg \max P(S|\theta)$  - MLE

Bayesian  $\rightarrow$  prior  $P(\theta) \rightarrow \arg \max_{\theta} P(\theta|S)$   
distribution

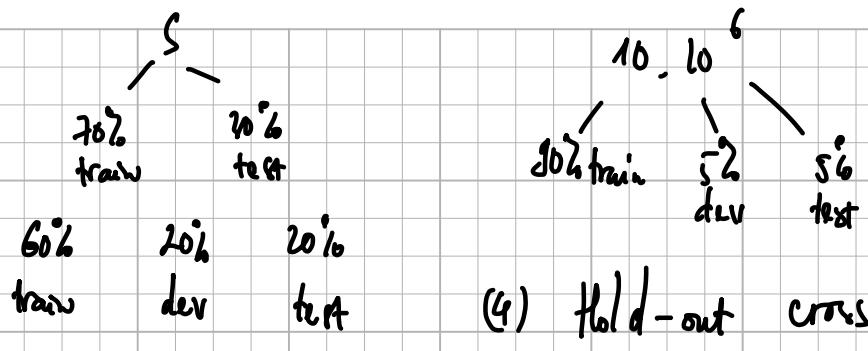
- MAP

maximum a posterior

## (3) Train / dev / test sets

- Train each model i (option for degree of polynomial)  
on train. Get some hypothesis  $h_i$
- Measure error  $S_{\text{dev}}$   $\rightarrow$  pick lowest  $S_{\text{dev}}$
- Evaluate on separate  $S_{\text{test}}$  (optional)

CIFAR



"development set" = "cross validation set" for  $d = 1 \dots 5$   
 for  $i = 1 \dots k$

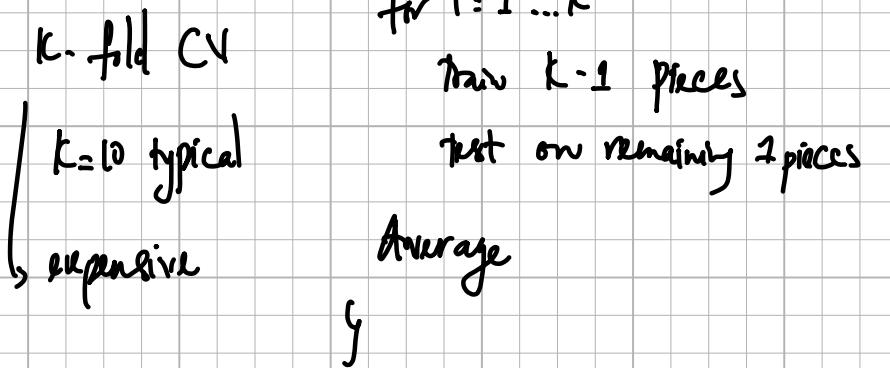
- Small data sets

$$m = 100$$

$$70\%_{\text{train}}, 30\%_{\text{dev}}$$

$$m \in [20, 50]$$

$$\hookrightarrow K = m$$



- feature selection

Start with  $f = \emptyset$

Repeat {

- 1) Try add each feature  $i \in f$ ,  
and see which single-feature addition  
most improves dev set performance.

- 2) Add that feature to  $f$

}

CS230  
ml.yearnly.org

## Lecture 09 - Approx / Estimation Error & ERM

- Setup / Assumptions
- Bias / Variance
- Approx Esti
- Empirical Risk Minimizer
- Uniform Convergence
- VC dimension

Assumptions:

1. Data distribution D

$$(x_i, y) \sim D \quad \begin{matrix} \text{- train} \\ \text{- Test} \end{matrix}$$

2. Independent Samples

$x^{(i)}, y^{(i)}$
:
$x^{(m)}, y^{(m)}$

[Estimator]

$$\sim D \Rightarrow$$



$\hat{h}$  or  $\hat{\theta} \sim$  Sampling distr

Random var

Deterministic func

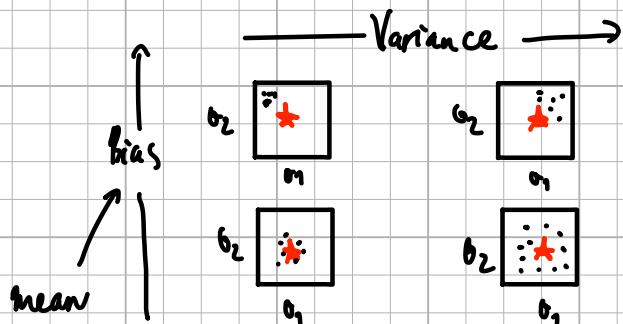
Random Var

$\theta^*$  or  $h^*$  "True" parameter - Not random, Const distribution

Bias / Variance

→ Under fit      "just right"      Overfit

→ Parameter view



$m \rightarrow +\infty$

$$\text{Var}[\hat{\theta}] \rightarrow 0$$

"stats efficiency"

$\hat{\theta} \rightarrow \theta^*$  as  $m \rightarrow \infty$ : Consistent

$$E(\hat{\theta}) = \theta^* \text{ for all } m$$

- fighting variance

(i)  $m \rightarrow \infty$

Small Bz

(ii) Regularization  $L_1, L_2 \rightarrow$  Low Var

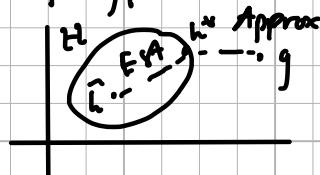
where

$g$  - Best possible Hypothesis

$h^*$  - Best in class  $\mathcal{H}$

$\hat{h}$  - learn from finite data

Space of Hypotheses



$E(g)$  = Bayes Error (Irreducible Error)

$E(h)$  - Risk / Generalization Error

$E(h^*) - E(g)$ : Approx Error } Class

$$= E_{(x,y) \sim D} [1 \{ h(x) \neq y \}]$$

$E(\hat{h}) - E(h^*)$ : Estimation Error } Data

$\hat{E}_s(h)$  : Empirical Risk

$$= \frac{1}{m} \sum_{i=1}^m 1 \{ h(x^{(i)}) \neq y^{(i)} \}$$

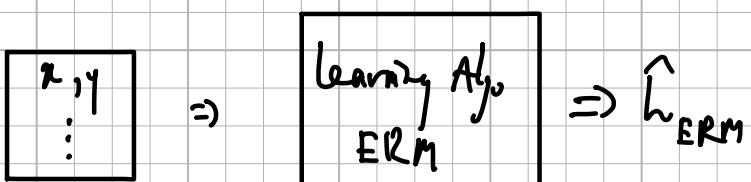
$$E(\hat{h}) = \text{Estimation} + \text{Approx} + \text{Irreducible}$$

Error      Error      Error

Est Var    Est Bias    +    Approx Err    +    Irreducible Err  
 Var                  Bias

- fight high bias
  - make  $\hat{h}$  bigger  $\rightarrow$  Bias / Var trade-off

Empirical Risk minimizer (ERM)



$$\hat{h}_{\text{ERM}} = \underset{h \in \mathcal{H}}{\operatorname{arg\,min}} \frac{1}{m} \sum_{i=1}^m \ell(h(x^{(i)}) + y^{(i)})$$

Uniform Convergence  $\rightarrow$  tools

①  $\hat{E}(h)$  vs  $E(h)$

i/ Union Bound

$A_1, A_2 \dots A_K$  (need not be independent)

$$P(A_1 \cup A_2 \dots \cup A_K) \leq P(A_1) + P(A_2) + \dots + P(A_K)$$

ii/ Hoeffding inequality

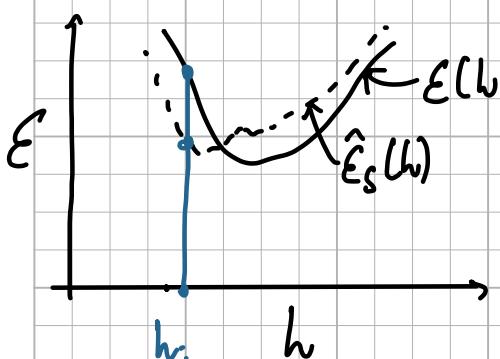
let  $z_1, z_2 \dots z_m \sim \text{Bernoulli}(\phi)$

$$\hat{\phi} = \frac{1}{m} \sum_{i=1}^m z_i$$

let  $\gamma > 0$  [margin]

$$\Pr \left[ \left| \hat{\phi} - \phi \right| > r \right] \leq 2 \exp(-2r^2 m)$$

deviation error margin



$$E[\hat{E}(h_i)] = E(h_i)$$

• Finite Hypotheses

$$|\mathcal{H}| = K$$

$$\Pr \left[ \exists h \in \mathcal{H} : |\hat{E}(h) - E(h)| > \gamma \right] \leq K \cdot 2 \exp(-2\gamma^2 n)$$

$$\Rightarrow \Pr \left[ \forall h \in \mathcal{H} : |\hat{E}_s(h) - E(h)| < \gamma \right] \geq 1 - \underbrace{2K \exp(-2\gamma^2 n)}_{\delta}$$

$$\text{let } \delta = 2K \exp(-\frac{1}{2} \gamma^2 m)$$

$\delta$ : prob of error

$\Gamma$ : margin of error

$m$ : Sample size

fixed  $\delta, \Gamma > 0$

$$m \geq \frac{1}{\gamma^2} \log \frac{2K}{\delta} \Rightarrow \text{Sample Complexity}$$

$$\hat{\epsilon}(h) \leq \hat{\epsilon}(h^*) + \gamma \leq \hat{\epsilon}(h^*) + \gamma \leq \hat{\epsilon}(h^*) + 2\gamma$$

$\Rightarrow$  with prob  $1 - \delta$ , training size  $m$

$$\epsilon(h) \leq \epsilon(h^*) + 2 \sqrt{\frac{1}{m} + \log \frac{2K}{\delta}}$$

- VC dimension

$$VC(T_2) = K$$

$$\epsilon(h) \leq \epsilon(h^*) + O\left(\sqrt{\frac{VC(T_2)}{m}} + \log\left[\frac{m}{VC(T_2)} + \frac{1}{m}\log\frac{1}{\delta}\right]\right)$$

## Lecture 10

## Decision Trees and Ensemble Methods

- Decision Trees

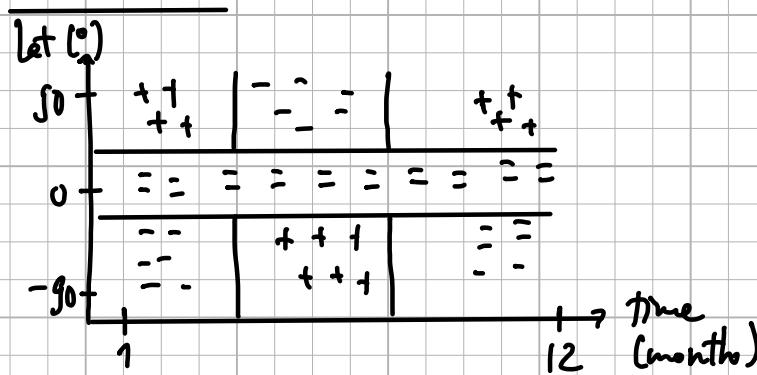
- Ensemble Methods

- Bagging

- Random Forests

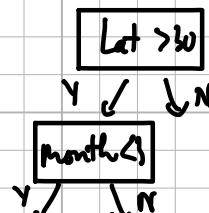
- Boosting

Decision Trees - Non linear Model



Greedy, Top-Down, Recursive

Partitioning  
↓



- Region  $R_p$
- looking for a split  $S_p$
- $S_p(j, t)$
- $= \{(x | x_j < t, x \in R_p)\}$
- $\{x | x_j > t, x \in R_p\}$
- $R_L$

## How to choose Splits?

$$\text{Define } L(R) : \text{loss on } R \Rightarrow L_{\text{missclass}} = 1 - \max_c \hat{P}_c$$

Given  $C$  classes, define

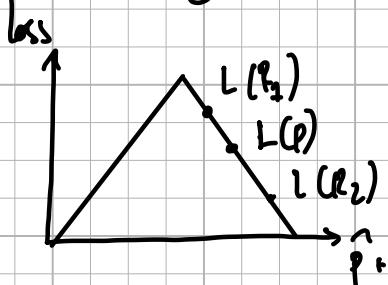
$\hat{P}_c$  to be proportion of examples  
in  $R$  that are of class  $c$

$$\max_j L(\hat{P}_j) - (\underbrace{L(R_1) + L(R_2)}_{\text{parent loss}})$$

→ Missclassification loss has issues

Instead, define cross-entropy loss

$$L_{\text{cross}} = -\sum_c \hat{P}_c \log_2 \hat{P}_c \rightarrow \text{proportion right}$$



→ Missclassification

Reg trees  $R_m$

$$\rightarrow \text{Predict } \hat{y}_m = \frac{\sum_{i \in R_m} y_i}{|R_m|}$$

$$\Rightarrow L_{\text{squared}} = \frac{\sum_{i \in R_m} (y_i - \hat{y}_m)^2}{|R_m|}$$

Categorical Vars

$$k \in \{N\}$$

$k$  categories  $\Rightarrow 2^k$  possible splits

Regularization of DTs

1) min leaf size

2) max depth

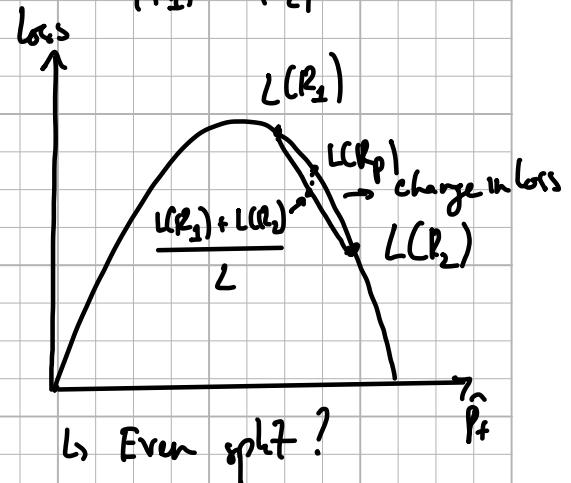
3) max number of nodes

4) minimum decrease in loss

5) Pruning

(missclassification with val set)

$$L(R_p) = \frac{|R_1| L(R_1) + |R_2| L(R_2)}{|R_1| + |R_2|}$$



→ Cross-entropy

$$\sum \hat{P}_c (1 - \hat{P}_c) \rightarrow \text{same shape}$$

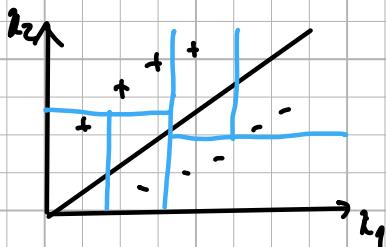
## Runtime

n examples  
f features  
d depth of tree

test time: Train time:

$O(d)$  Each point is part of  $O(d)$  nodes  
 $d < \log_2 n$  Cost of print at each node is  $O(f)$   
So total cost  $O(nfd)$   
Data matrix is of size  $nf$

## No additive structure



## Recap

- ⊕ Easy to explain
- ⊖ High variance
- ⊕ Interpretable
- ⊖ Bad at additive
- ⊕ Categorical vars
- ⊖ Low predictive acc
- ⊕ Fast

## Ensemble → Groups to one

Take  $X_i$ 's which are random variables (RV)  
that are independent identically distributed (iid)

$$\text{Var}(X_i) = \sigma^2 \quad \text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum X_i\right) = \frac{\sigma^2}{n}$$

Drop independence assumption

so now  $X_i$ 's are i.d.,  $X_i$ 's correlated by  $\rho$  (rho)

$$\text{Var}(\bar{X}) = \rho \sigma^2 + \frac{1-\rho}{n} \sigma^2 \quad \downarrow$$

## Ways to ensemble

1/ different algs

3/ Bagging (Random Forest)

2/ different training sets

4/ Boosting (Adaboost, xgboost)

## bagging - bootstrap Aggregation

↳ to measure uncertainty of your est

Have a true Population P

Training Set  $S \sim P$

Assume  $P = S$

Bootstrap samples  $Z \sim S$   
 $\downarrow$   
 $z_1, \dots, z_M$

Train model  $G_m$  on  $Z_m$   
$$G(x) = \frac{\sum_{m=1}^M G_m(x)}{M}$$

## Bias / Variance Analysis

$\text{Var}(\bar{x}) = p\sigma^2 + \frac{1-p}{M}\sigma^2$ , has  
Bootstrapping is driving down  $p$ , lower bound  
More  $M \rightarrow$  less variance  
Bias slightly increase  
because of random subsampling

## DTs + Bagging

DT are high variance, low bias

Ideal fit for bagging

## Random forests

At each split, consider only a fraction of your total features

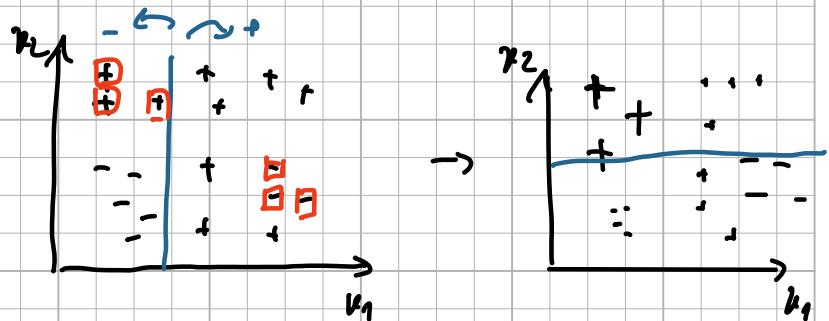
Decrease  $p$

Decorrelate Models

## Boosting

Decrease bias

Additive



$$G(x) = \sum_m \omega_m G_m$$

Each  $G$  trained on reweighted training set

Determine for classifier  $G_m$  a weight  $\omega_m$  proportional to  $\log\left(\frac{1 - \text{err}_m}{\text{err}_m}\right)$

Adaboost + Xgboost (Lecture Notes)