

SSU

- MINIMALIZACE EMPIRICKÉHO RIZIKA
- ODHAD MAXIMÁLNÍ VĚROHODNOSTI
- EM ALGORITHMUS
- HLUBOKÉ SÍŤE A JEJICH UČENÍ
- STANDARDNÍ A HLUBOKÉ NEURONOVÉ SÍŤE A JEJICH UČENÍ

- OBJECT FEATURES
 - $x \in X$
 - ARE OBSERVABLE

- STATE OF OBJECT
 - $z \in Y$
 - USUALLY HIDDEN

- PREDICTION STRATEGY
 - INFERENCE RULE
 - $h: X \rightarrow Y$
 - IF z IS
 - CATEGORICAL VARIABLE
 - CLASSIFICATION
 - REAL VALUED VARIABLE
 - REGRESSION

- TRAINING EXAMPLES
 - $\{(x_i, z_i) \mid x_i \in X, z_i \in Y\}$

- LOSS FUNCTION
 - $\ell: Y \times Y \rightarrow \mathbb{R}_+$
 - PENALISES WRONG ANSWERS (PREDICTIONS)
 - $\ell(z, h(x))$ IS LOSS FOR PREDICTING $z' = h(x)$ IF TRUE STATE IS z

- GOAL IS TO FIND OPTIMAL PREDICTION STRATEGY R

- THAT MINIMALISES COST

- STATISTICAL MACHINE LEARNING

- X, Y ARE RANDOM VARIABLES

- THEY ARE JOINED BY UNKNOWN PROBABILITY DISTRIBUTION FUNCTION $P(X, Y)$

- WE CAN COLLECT EXAMPLES (x, y) DRAWN FROM $P(X, Y)$

- TYPICAL CONCEPTS

- REGRESSION

$$Y = f(x) + \epsilon$$

- WHERE f IS UNKNOWN AND ϵ IS RANDOM ERROR

- CLASSIFICATION

$$P(X, Y) = P(Y) \cdot P(X | Y)$$

- WHERE $P(Y)$ IS PRIOR CLASS PROBABILITY AND $P(X | Y)$ IS CONDITIONAL FEATURE DISTRIBUTION

- RISK OF INFERENCE RULE

$$R(R) = E \mathcal{L}(Y, R(X)) = \sum_{x \in X} \sum_{y \in Y} P(x, y) \cdot \mathcal{L}(y, R(x))$$

- HOW TO ESTIMATE IT, IF $P(X, Y)$ IS NOT KNOWN

- TESTING SET

- COLLECT I.I.D. TEST SAMPLES $S^{\mathcal{L}} = \{(x^i, y^i) \in X \times Y \mid i = 1 \dots \mathcal{L}\}$ DRAWN FROM $P(X, Y)$

- ESTIMATE BY EMPIRICAL RISK

$$R(R) \approx R_{S^{\mathcal{L}}}(R) = \frac{1}{\mathcal{L}} \sum_{i=1}^{\mathcal{L}} \mathcal{L}(y^i, R(x^i))$$

- CHOOSING OPTIMAL INFERENCE RULE $R(x)$

- IF $P(x, y)$ IS KNOWN

- BEST POSSIBLE RISK IS

$$R^* = \inf_{R \in \mathcal{Y}^X} R(R) = \inf_{R \in \mathcal{Y}^X} \sum_{x \in X} \sum_{y \in Y} P(x, y) \ell(y, R(x)) = \sum_{x \in X} P(x) \inf_{y' \in Y} \sum_{y \in Y} P(y|x) \ell(y, y')$$

- THE BEST INFERENCE RULE IS BAYES RULE

$$R^*(x) = \operatorname{ARGMIN}_{y' \in Y} \sum_{y \in Y} P(y|x) \cdot \ell(y, y')$$

- BUT $P(x, y)$ IS NOT KNOWN

- LEARNING TYPES

- TRAINING DATA

- $T^m = \{(x^i, y^i) \in X \times Y \mid i = 1, \dots, m\}$... SUPERVISED DATA

- $T^m = \{x^i \in X \mid i = 1, \dots, m\}$... UNSUPERVISED DATA

- $T^m = T^m_{\ell} \cup T^m_w$, WHERE T^m_{ℓ} ARE LABELED AND T^m_w ARE UNLABELED \therefore ... SEMI-SUPERVISED LEARNING

- PRIOR KNOWLEDGE ABOUT TASK

- DISCRIMINATIVE LEARNING

- OPTIMAL INFERENCE RULE R^* IS IN SOME CLASS OF RULES H

- REPLACE TRUE RISK BY EMPIRICAL RISK

$$R_T(R) = \frac{1}{|T|} \sum_{(x, y) \in T} \ell(y, R(x))$$

- MINIMIZE WITH RESPECT TO $R \in H$

$$R_T^* = \operatorname{ARGMIN}_{R \in H} R_T(R)$$

- SVM, DNN, STRUCTURED OUTPUT SVM

- GENERATIVE LEARNING

- TRUE $P(x, y)$ IS IN SOME PARAMETRISED FAMILY OF DISTRIBUTIONS

- $P = P_{\theta^*} \in \mathcal{P}$ FOR EXAMPLE

- SPLIT DEPENDENCE ON UNKNOWN $\theta \in \Theta$ AND RANDOM T

$$-1. \theta_T^* = \underset{\theta \in \Theta}{\operatorname{ARGMAX}} P_{\theta}(T)$$

- MAXIMUM LIKELIHOOD ESTIMATOR

-2. SET $R_T^* = R_{\theta_T^*}$ WHERE R_{θ} DENOTES BAYES INFERENCE RULE FOR PROBABILITY DISTRIBUTION P_{θ}

- MIXTURE MODELS, HIDDEN MARKOV MODEL, MARKOV RANDOM FIELD

- EMPIRICAL RISK MINIMIZATION

- WE WANT TO FIND PREDICTION STRATEGY WITH MINIMAL EXPECTED RISK

$$R(R) = \int \sum_{y \in Y} \ell(y, R(x)) P(x, y) dx = E_{(x, y) \sim P} (\ell(y, R(x)))$$

- BUT WE DON'T KNOW $P(x, y)$

- SO WE CAN MINIMIZE EMPIRICAL RISK

$$R_{S^2}(R) = \frac{1}{2} \sum_{i=1}^2 \ell(y^i, R(x^i))$$

- IS IT GOOD APPROXIMATION?

- LAW OF LARGE NUMBERS

- ARITHMETIC MEAN OF RESULTS OF RANDOM TRIALS GETS CLOSER TO THE EXPECTED VALUE AS MORE TRIALS ARE PERFORMED

- Hoeffding Inequality

- LET $\{z^1, \dots, z^2\} \in [a, b]^2$ BE REALIZATIONS OF INDEPENDENT RANDOM VARIABLES WITH SAME EXPECTED VALUE μ . THEN FOR ANY $\epsilon > 0$ IT HOLDS THAT:

$$P\left(\left|\frac{1}{2} \sum_{i=1}^2 z^i - \mu\right| \geq \epsilon\right) \leq 2e^{-\frac{2\epsilon^2}{(b-a)^2}}$$

"PROBABILITY OF RESULT BEING DIFFERENT THAN ϵ IS SMALLER THAN..."

- FOR WHICH ϵ IS $\hat{\mu}_\epsilon$ IN INTERVAL $(\hat{\mu}_\epsilon - \epsilon, \hat{\mu}_\epsilon + \epsilon)$ WITH PROBABILITY AT LEAST ν ?

$$P(|\hat{\mu}_\epsilon - \mu| < \epsilon) = 1 - P(|\hat{\mu}_\epsilon - \mu| \geq \epsilon) \geq 1 - 2e^{-\frac{2\epsilon^2}{(b-a)^2}} = \nu$$

- FROM THIS

$$\epsilon = |b-a| \sqrt{\frac{\log(2) - \log(1-\nu)}{2\epsilon^2}}$$

- MINIMAL NUMBER OF SAMPLES TO $\hat{\mu}$ BE IN $(\hat{\mu}_\epsilon - \epsilon, \hat{\mu}_\epsilon + \epsilon)$ WITH PROBABILITY AT LEAST ν

$$l = \frac{\log(2) - \log(1-\nu)}{2\epsilon^2} (b-a)^2$$

- PROBABILITY OF BEING BAD TEST SET CAN BE BOUNDED BY

- l_{\min}, l_{\max} IS MAXIMAL AND MINIMAL VALUE OF LOSS FUNCTION

$$P(|R_{S_2}(h) - R(h)| \geq \epsilon) \leq 2e^{-\frac{2\epsilon^2}{(l_{\min} - l_{\max})^2}}$$

- LEARNING BY EMPIRICAL RISK MINIMALIZATION

-ERM

-EMPIRICAL RISK

$$R_{T^m}(h) = \frac{1}{m} \sum_{i=1}^m \ell(y^i, h(x^i))$$

- ERM LEARNING ALGORITHM RETURNS h_m

$$h_m \in \underset{h \in H}{\text{ARGMIN}} R_{T^m}(h)$$

- DEPENDING ON H AND ℓ WE GET INSTANCES OF

-SVM

-LINEAR REGRESSION

-LOGISTIC REGRESSION

-NN

ADA BOOST

- BUT IF WE DON'T CONSTRAIN H

- WE DON'T HAVE ANY GUARANTEE THAT EMPIRICAL RISK IS GOOD APPROXIMATION OF TRUE RISK

- AND REGARDLESS OF THE NUMBER OF EXAMPLES

- IF WE DON'T THROW AWAY MEMORIZING PREDICTING STRATEGY FOR EXAMPLE

- EXCESS ERROR

- DEVIATION OF LEARNED PREDICTOR FROM THE BEST ONE

$$\underbrace{(R(\mathcal{A}_m) - R^*)}_{\text{EXCESS ERROR}} = \underbrace{(R(\mathcal{A}_m) - R(\mathcal{A}_H))}_{\substack{\text{ESTIMATION ERROR} \\ \text{"ERROR BY NOT PROBABLY SELECTED FROM } H"}} + \underbrace{(R(\mathcal{A}_H) - R^*)}_{\substack{\text{APPROXIMATION ERROR} \\ \text{"ERROR BY BOUNDING SPACE OF } H"}}$$

- STATISTICALLY CONSISTENT LEARNING ALGORITHM

- IF FOR ANY $\epsilon > 0$ AND $\delta > 0$ HOLDS:

$$\lim_{m \rightarrow \infty} P(R(\mathcal{A}_m) - R(\mathcal{A}_H) \geq \epsilon) = 0$$

- "WE CAN MAKE ESTIMATION ERROR ARBITRARILY SMALL IF WE HAVE ENOUGH EXAMPLES"

- COMBINING WITH UNIFORM LAW OF LARGE NUMBERS

- "PROBABILITY OF BEING "BAD TRAINING SET" FOR AT LEAST ONE HYPOTHESIS FROM H CAN BE MADE ARBITRARILY LOW IF WE HAVE ENOUGH EXAMPLES"

- LINEAR CLASSIFIER

$$h(x; w, b) = \text{SIGN}(\langle w, \overset{\text{FEATURES}}{\phi(x)} \rangle + b) = \begin{cases} +1 & \text{IF } \langle w, \phi(x) \rangle + b \geq 0 \\ -1 & \text{IF } \langle w, \phi(x) \rangle + b < 0 \end{cases}$$

- WITH MINIMAL EXPECTED RISK

$$R^{0/1}(h) = \mathbb{E}_{(x, y) \sim P} (l^{0/1}(y, h(x))) \quad \text{WHERE } l^{0/1}(y, y') = \mathbb{I}[y \neq y']$$

- ERM LEADS TO

$$(w^*, b^*) \in \underset{(w, b) \in (\mathbb{R}^n \times \mathbb{R})}{\text{ARGMIN}} R_{T^m}^{0/1}(h(\cdot; w, b))$$

- WHERE EMPIRICAL RISK IS

$$R_{T^m}^{0/1}(h(\cdot; w, b)) = \frac{1}{m} \sum_{j=1}^m \mathbb{I}[y^j \neq h(x^j; w, b)]$$

- VAPNIK - CHERVOMENKIS DIMENSION

- VC

- LET $H \subseteq \{-1, +1\}^X$ AND $\{x^1, \dots, x^m\} \in X^m$ BE A SET OF m INPUT OBSERVATIONS

- THE SET $\{x^1, \dots, x^m\}$ IS SHATTERED BY H IF FOR ALL $\sigma \in \{-1, +1\}^m$ THERE EXISTS $h \in H$ SUCH THAT $h(x^i) = \sigma^i, i \in \{1, \dots, m\}$

- VC DIMENSION OF H IS CARDINALITY OF LARGEST SET OF POINTS FROM X WHICH CAN BE SHATTERED BY H .

- VC DIMENSION OF LINEAR CLASSIFIER OPERATING IN n -DIMENSIONAL SPACE IS $n+1$

or

- EXAMPLES ARE LINEARLY SEPARABLE IF THERE EXISTS $(w, b) \in \mathbb{R}^{n+1}$ SUCH THAT

$$\sigma^i (\langle w, \phi(x^i) \rangle + b) > 0, i \in \{1, \dots, m\}$$

- CAN BE SOLVED BY PERCEPTRON ALGORITHM

- SVM

- PRIMAL PROBLEM

- ERROR OF LINEAR CLASSIFIER

$$(w', b') = \underset{\|w\| \leq r, b \in \mathbb{R}}{\text{ARGMIN}} \left(\frac{1}{m} \sum_{i=1}^m \max \{ 0, 1 - \gamma^i (\langle w, \phi(x^i) \rangle + b) \} \right)$$

- IT CAN BE FORMULATED AS EQUIVALENT QUADRATIC PROGRAM

$$(w^*, b^*, f^*) = \underset{\substack{(w, b) \in \mathbb{R}^{m+1} \\ f \in \mathbb{R}^m}}{\text{ARGMIN}} \left(\frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_{i=1}^m f_i \right)$$

REGULARIZATION CONSTANT

s.t.

$$\gamma^i (\langle w, \phi(x^i) \rangle + b) \geq 1 - f_i \quad i \in \{1, \dots, m\}$$

$$f_i \geq 0$$

- DUAL PROBLEM

- CONVEX QUADRATIC PROGRAM

$$\alpha^* = \underset{\alpha \in \mathbb{R}^m}{\text{ARGMAX}} \left(\sum_{i=1}^m \alpha_i - \frac{1}{2\lambda} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j \gamma^i \gamma^j \langle \phi(x^i), \phi(x^j) \rangle \right)$$

$$\text{s.t.} \quad \sum_{i=1}^m \alpha_i \gamma^i = 0$$

$$0 \leq \alpha_i \leq \frac{1}{m} \quad i \in \{1, \dots, m\}$$

- THE SOLUTION OF PRIMAL SOLUTION (w^*, b^*) IS OBTAINED FROM α^* BY

$$w^* = \sum_{i=1}^m \gamma^i \phi(x^i) \alpha_i^* \quad , \quad b^* = \gamma^i - \langle w^*, \phi(x^i) \rangle \quad \forall i \in I_{SV}^<$$

WHERE $I_{SV}^< = \{ i \in \{1, \dots, m\} \mid 0 < \alpha_i^* < \frac{1}{m} \}$ ARE BOUNDARY SVs

- TO REPRESENT THE CLASSIFIER WE NEED ONLY SUPPORT VECTORS

- TRAINING EXAMPLES WITH INDICES $I_{SV} = \{ i \in \{1, \dots, m\} \mid \alpha_i^* > 0 \}$

- KERNEL SVM

- WE CAN USE KERNEL FUNCTION

$$k(x, x') = \langle \phi(x), \phi(x') \rangle$$

- WE CAN TRAIN SVM ONLY USING KERNEL FUNCTIONS

$$\alpha^* = \underset{\alpha \in \mathbb{R}^m}{\text{ARGMIN}} \left(\sum_{i=1}^m \alpha_i - \frac{1}{2\gamma} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j k(x^i, x^j) \right)$$

↑ INSTEAD OF $\langle \phi(x), \phi(x') \rangle$

$$\text{s.t.} \quad \sum_{i=1}^m \alpha_i y_i = 0$$
$$0 \leq \alpha_i \leq \frac{1}{m} \quad i \in \{1, \dots, m\}$$

- EVALUATION OF PREDICTION RULE

$$h(x; \alpha^*, b^*) = \text{SIGN}(\langle w^*, \phi(x) \rangle + b^*) = \text{SIGN} \left(\sum_{i \in I_{SV}} y_i \alpha_i^* k(x^i, x) + b^* \right)$$

- EXAMPLE OF KERNEL

$$k(x, x') = \langle \phi(x), \phi(x') \rangle = (\langle x, x' \rangle + 1)^2$$

- POLYNOMIAL OF d DEGREE

$$k(x, x') = \langle x, x' \rangle^d$$

- INHOMOGENEOUS POLYNOMIAL OF d DEGREE

$$k(x, x') = (\langle x, x' \rangle + 1)^d$$

- GAUSSIAN KERNEL

$$k(x, x') = \exp(-\sigma \|x - x'\|^2)$$

- HILBERT SPACE

- COMPLETE VECTOR SPACE WITH DOT PRODUCT $\langle \cdot, \cdot \rangle$
 $H \times H \rightarrow \mathbb{R}$

- SATISFIES FOLLOWING REQUIREMENTS

- SYMMETRY

$$\langle f, g \rangle = \langle g, f \rangle \quad \forall f, g \in H$$

- LINEARITY

$$\langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle = \alpha_1 \langle f_1, g \rangle + \alpha_2 \langle f_2, g \rangle \quad \forall f_1, f_2, g \in H$$

- POSITIVE DEFINITENESS

$$\alpha_1, \alpha_2 \in \mathbb{R}$$

$$\langle f, f \rangle \geq 0 \quad \text{WITH EQUALITY IFF } f=0$$

~~2~~

- DOT PRODUCT DEFINES A NORM $\|f\|_H = \sqrt{\langle f, f \rangle}$

- POSITIVE DEFINITE KERNEL

- LET X BE NON-EMPTY SET

- FUNCTION $\lambda: X \times X \rightarrow \mathbb{R}$ IS POSITIVE DEFINITE KERNEL IF IT IS

- SYMMETRIC

- AND FOR ANY FINITE SET OF INPUTS x^1, \dots, x^m THE KERNEL MATRIX $K \in \mathbb{R}^{m \times m}$ WITH ELEMENTS $K_{i,j} = \lambda(x^i, x^j)$ IS POSITIVE SEMI-DEFINITE

$$K = \begin{pmatrix} \lambda(x^1, x^1) & \lambda(x^1, x^2) & \dots & \lambda(x^1, x^m) \\ \lambda(x^2, x^1) & \lambda(x^2, x^2) & \dots & \lambda(x^2, x^m) \\ \vdots & \vdots & \ddots & \vdots \\ \lambda(x^m, x^1) & \lambda(x^m, x^2) & \dots & \lambda(x^m, x^m) \end{pmatrix}$$

- POSITIVE SEMI-DEFINITE MATRIX $K \in \mathbb{R}^{m \times m}$

$$- \forall \alpha \in \mathbb{R}^m: \alpha^T K \alpha \geq 0$$

- NEURAL NETWORKS

$$\hat{y} = f(s) = f\left(\sum_{i=1}^n w_i x_i + b_f\right)$$

- WE CAN ADD BIAS TO w_1 AND ENHANCE EACH x BY

$$(1, x)$$

- ACTIVATION FUNCTIONS

- SIGMOID

$$\sigma(s) = \frac{1}{1 + e^{-s}} = \frac{e^s}{e^s + 1}$$

- TANH

$$\text{TANH}(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}} = 2\sigma(s) - 1$$

- RELU

$$f(s) = \max(0, s)$$

- LINEAR NEURON

- SINGLE NEURON WITH ~~ACTIVATION~~ LINEAR ACTIVATION FUNCTION

- LINEAR REGRESSION

$$y = \langle x, w \rangle$$

- SQUARED ERROR

$$L(w) = \sum_{i=1}^n \underbrace{(y_i - \langle w, x_i \rangle)^2}_{L(y_i, \hat{y}_i)} = (y - Xw)^T (y - Xw)$$

- MINIMIZATION OF $L(w)$

$$\frac{\partial L}{\partial w} = 0 \quad \therefore w^* = (X^T X)^{-1} X^T y$$

- LOGISTIC REGRESSION

- NEURON USING SIGMOID ACTIVATION FUNCTION

$$\hat{y} = \sigma(\langle W, X \rangle)$$

- y IS TARGET CLASS

- \hat{y} IS OUTPUT CLASS

- LIKELIHOOD OF LOGISTIC REGRESSION

$$P(y|W, X) = \prod_{i=1}^m \hat{y}_i^{y_i} (1 - \hat{y}_i)^{1 - y_i}$$

- NEGATIVE LOG-LIKELIHOOD

$$L(W) = \sum_{i=1}^m - [y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)]$$

$\underbrace{\hspace{15em}}_{\mathcal{L}(y_i, \hat{y}_i)}$

- CROSS ENTROPY

- ~~MINIMIZE~~ $\frac{\partial L}{\partial W} = 0$ HAS NO ANALYTICAL SOLUTION

- WE HAVE TO USE ANALYTICAL METHODS

- LINEAR LAYER

$$Y = XW$$

- SOFTMAX LAYER

- MULTINOMIAL CLASSIFICATION

- K MUTUALLY EXCLUSIVE CLASSES

$$\tilde{\sigma}_k(s) = \frac{e^{s_k}}{\sum_{c=1}^K e^{s_c}}$$

- REPRESENTS PROBABILITY DISTRIBUTION

- DESCRIBES CLASS MEMBERSHIP PROPERTIES $P(y = k | s) = \tilde{\sigma}_k(s)$

- MULTINOMIAL LOGISTIC REGRESSION

- LINEAR LAYER + SOFTMAX LAYER

$$- \hat{y}_k = \sigma_k(x^T W)$$

$$- \text{CLASSIFIER } \hat{y}(x, W) = \underset{k}{\text{ARGMAX}} \hat{y}_k$$

- LOSS FUNCTIONS

- BINARY CLASSIFICATION

- CROSS-ENTROPY

$$\left(- \sum_{i=1}^m [y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)] \right)$$

- MULTINOMIAL CLASSIFICATION

- MULTINOMIAL CROSS-ENTROPY

$$\left(- \sum_{i=1}^m \sum_{c=1}^K y_{ic} \log(\hat{y}_{ic}) \right)$$

- REGRESSION

- SQUARED ERROR

$$\sum_{i=1}^m (y_i - \hat{y}_i)^2$$

- MULTI-OUTPUT REGRESSION

- SQUARED ERROR

$$\sum_{i=1}^m \sum_{c=1}^K (y_{ic} - \hat{y}_{ic})^2$$

- BACKPROPAGATION OVERVIEW

- METHOD TO COMPUTE GRADIENT OF LOSS FUNCTION WITH RESPECT TO ITS PARAMETERS

$$\nabla L(w)$$

- GRADIENT IS USED BY OPTIMIZATION METHODS

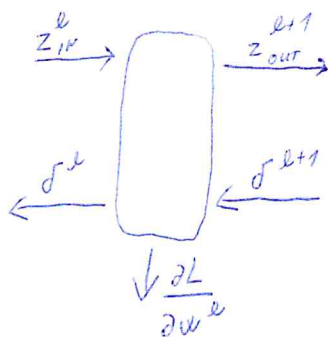
- EXAMPLE OF MULTINOMIAL LOGISTIC REGRESSION

- LOSS IS MULTINOMIAL CROSS-ENTROPY

$$L(w) = - \sum_{i=1}^m \sum_{c=1}^K [\gamma_i = c] \log \left(\frac{\exp(\langle x_i, w_c \rangle)}{\sum_{a=1}^K \exp(\langle x_i, w_a \rangle)} \right)$$

- USAGE OF CHAIN RULE

- DIVIDE AND CONQUER APPROACH



- $\delta^l = \frac{\partial L}{\partial z^l}$ IS SENSITIVITY OF LOSS TO NODE INPUT FOR LAYER L

- THEN:

$$\delta_i^l = \frac{\partial L}{\partial z_i^l} = \sum_j \frac{\partial L}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial z_i^l} = \sum_j \delta_j^{l+1} \frac{\partial z_j^{l+1}}{\partial z_i^l}$$

- WE CAN DO IT SIMILARLY FOR PARAMETERS

$$\frac{\partial L}{\partial w_i^l} = \sum_j \frac{\partial L}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial w_i^l} = \sum_j \delta_j^{l+1} \frac{\partial z_j^{l+1}}{\partial w_i^l}$$

- FOR EACH LAYER, WE NEED TO SPECIFY

- FORWARD PASS: $z^{l+1} = f(z^l)$

- BACKWARD PASS: $\frac{\partial z^{l+1}}{\partial z^l}$

- PARAMETER PASS (OPTIONAL): $\frac{\partial z^{l+1}}{\partial w^l}$

- EXAMPLE

- LINEAR LAYER

- FORWARD

$$z_j^{l+1} = \sum_{i=0}^n w_{ij} z_i^l \quad j = 1, \dots, K$$

- BACKWARD

$$\frac{\partial z_j^{l+1}}{\partial z_i^l} = w_{ij} \quad i = 0, \dots, n \quad j = 1, \dots, K$$

- PARAMETER

$$\frac{\partial z_j^{l+1}}{\partial w_{i,j}} = \mathbb{1}[j=i] z_i^l$$

- SQUARED ERROR

- FORWARD

$$z^{l+1} = \sum_{i=1}^n (y_i - z_i^l)^2$$

- BACKWARD

$$\frac{\partial z^{l+1}}{\partial z_i^l} = -2 (y_i - z_i^l) \quad i \in \{1, \dots, n\}$$

- GRADIENT DESCENT

$$\theta^* = \underset{\theta}{\text{ARGMIN}} L(\theta)$$

- "GOAL IS TO FIND PARAMETERS θ THAT WILL MINIMIZE THE LOSS"

$$\theta^{(t+1)} = \theta^{(t)} - \eta^{(t)} \nabla L(\theta^{(t)})$$

LEARNING RATE

- UPDATING WEIGHTS

- FULL LEARNING

- AFTER ALL TRAINING EXAMPLES ARE USED

- ONLINE LEARNING

- AFTER EACH SAMPLE, BUT IT MAY HAVE CONVERGENCE IN LATER STAGES OF TRAINING

- BATCH LEARNING

- OTHER TYPES

- MOMENTUM

- SIMULATES INERTIA

- ADAGRAD

- ADAPTIVE GRADIENT METHOD

- REDUCE LEARNING RATES FOR PARAMETERS HAVING HIGH VALUES OF GRADIENT

- RMS PROP

- SIMILAR TO ADAGRAD

- REGULARIZATION

- DEALS WITH OVERFITTING

- OTHER POSSIBILITIES

- MORE DATA

- SIMPLER MODEL

- L2

$$- L(W) = (y - Xw)^T (y - Xw) + \lambda w^T w$$

↖ L2 REGULARIZATION

- MINIMIZE THE WEIGHTS IN NETWORK

- L1

- SUM ABSOLUTE VALUES OF WEIGHTS

- EARLY STOPPING

- USE VALIDATION SET

- STOP IS VALIDATION LOSS STARTS TO GROW

- DROPOUT

- AUGMENTING DATASET

- DON'T INITIALIZE WEIGHTS TO ZERO
 - NO GRADIENT
- INITIALIZE BY SMALL NUMBERS
 - GAUSSIAN INITIALIZATION
 - WORKS OK FOR SHALLOW NETWORK
- PROBLEM OF VANISHING GRADIENT
- CONVOLUTIONAL NEURAL NETWORKS
 - USES A SMALL RECEPTIVE FIELD WHICH SHARES WEIGHTS
 - TRANSLATION EQUIVARIANCE
 - EACH RECEPTIVE FIELD USES MULTIPLE FILTERS
- STRIDE
 - HIGHER STRIDE PRODUCES SMALLER OUTPUT
- ZERO PADDING
 - CONVOLUTIONAL LAYER REDUCES SPATIAL SIZE OF OUTPUT
 - THIS IS FIXED BY ZERO PADDING
- NONLINEARITIES
 - SIGMOID
 - TANH
 - RELU
- MAX POOLING
 - REDUCES SPATIAL RESOLUTION
 - ~~RE~~ LESS PARAMETERS
 - IT HELPS WITH OVERFITTING
- LE NET - 5
- ALEX NET
- VGG NET
- GOOGLE NET
- RES NET

- GENERATIVE LEARNING

- MODEL JOINT PROBABILITY DISTRIBUTION $P(x, y)$

- USE INFERENCE RULE

$$L(x) \in \operatorname{ARGMAX}_{y' \in Y} \sum_{y' \in Y} P(y' | x) \ell(y', y)$$

- LEARNING

- IF DISTRIBUTION $P_\theta(x, y)$ IS KNOWN UP TO PARAMETERS $\theta \in \Theta$ ONLY

- THEN ESTIMATE θ^* FROM TRAINING DATA

- GENERATIVE LEARNING MAKES STRONGER ASSUMPTIONS AND IS MORE DATA EFFICIENT WHEN ASSUMPTIONS ARE (NEARLY) CORRECT

- MLE

- MAXIMUM LIKELIHOOD ESTIMATION

$$\theta^*(T^m) = \operatorname{ARGMAX}_{\theta \in \Theta} \frac{1}{m} \sum_{i=1}^m \log P_\theta(x^i) = \operatorname{ARGMAX}_{\theta \in \Theta} \frac{1}{m} \sum_{i=1}^m \log \sum_{y \in Y} P_\theta(x^i, y)$$

- IF θ IS SINGLE PARAMETER OR VECTOR OF HOMOGENEOUS PARAMETERS

- MAXIMIZE LOG-LIKELIHOOD DIRECTLY

- IF θ IS COLLECTION OF HETEROGENEOUS PARAMETERS

- EXPECTATION MAXIMIZATION ALGORITHM (EM)

- INTRODUCE AUXILIARY VARIABLES $\alpha_i(\gamma) \geq 0$ S.T. $\sum_{\gamma \in Y} \alpha_i(\gamma) = 1$ FOR EACH SAMPLE

- CONSTRUCT LOWER BOUND OF LOG-LIKELIHOOD $L(\theta, T^m) \geq L_B(\theta, \alpha, T^m)$

$$L(\theta, T^m) \geq L_B(\theta, \alpha, T^m)$$

$$L(\theta, T^m) = \frac{1}{m} \sum_{i=1}^m \log \sum_{y \in Y} P_\theta(x^i, y) = \frac{1}{m} \sum_{i=1}^m \log \sum_{y \in Y} \frac{\alpha_i(\gamma)}{\alpha_i(\gamma)} P_\theta(x^i, y) \geq$$

$$\geq L_B(\theta, \alpha, T^m) = \frac{1}{m} \sum_{i=1}^m \sum_{y \in Y} \alpha_i(\gamma) \log P_\theta(x^i, y) - \frac{1}{m} \sum_{i=1}^m \sum_{y \in Y} \alpha_i(\gamma) \log \alpha_i(\gamma)$$

- MAXIMIZE LOWER BOUND BLOCK-COORDINATE ASCENT

- INITIALIZE SOME $\theta^{(0)}$

- E-STEP

- FIX CURRENT $\theta^{(t)}$

- MAXIMIZE $L_B(\theta^{(t)}, \alpha, T^m)$ W.R.T. α 's

$$\alpha_i^{(t)}(\gamma) = P_{\theta^{(t)}}(\gamma | x^i)$$

- M-STEP

- FIX CURRENT $\alpha^{(t)}$ ABOVE

- MAXIMIZE $L_B(\theta, \alpha^{(t)}, T^m)$ W.R.T. θ

$$\theta^{(t+1)} = \underset{\theta \in \Theta}{\text{ARGMAX}} \frac{1}{m} \sum_{i=1}^m \sum_{\gamma \in \mathcal{Y}} \alpha_i^{(t)}(\gamma) \log P_{\theta}(x^i, \gamma)$$

- THE SEQUENCE OF LIKELIHOOD VALUES IS INCREASING AND SEQUENCE $\alpha^{(t)}$ IS CONVERGENT

- HIDDEN MARKOV MODELS

- FOR

- TEXT RECOGNITION

- SPEECH RECOGNITION

- ROBOT SELF LOCALISATION

- $S = (s_1, s_2, \dots, s_m)$ IS SEQUENCE OF LENGTH m WITH ELEMENTS FROM FINITE SET K .

- JOINT PROBABILITY DISTRIBUTION

$$P(s_1, s_2, \dots, s_m) = P(s_1) \cdot P(s_2 | s_1) \cdot P(s_3 | s_2, s_1) \dots P(s_m | s_1, \dots, s_{m-1})$$

- JOINT PROBABILITY DISTRIBUTION IS MARKOV MODEL IF

$$P(S) = P(s_1) \cdot P(s_2 | s_1) \cdot P(s_3 | s_2) \dots P(s_m | s_{m-1}) = P(s_1) \cdot \prod_{i=2}^m P(s_i | s_{i-1})$$

- MOST PROBABLE SEQUENCE OF STATES CAN BE COMPUTED BY DYNAMIC PROGRAMMING

- LEARNING MARKOV MODEL

- WE ARE GIVEN TRAINING DATA

- ESTIMATE PARAMETERS OF MARKOV MODEL BY MAXIMUM LIKELIHOOD ESTIMATE

- $\alpha(s_{i-1} = l, s_i = k)$ IS FRACTION OF SEQUENCES IN T_m , FOR WHICH $s_{i-1} = l$ AND $s_i = k$

- ESTIMATES FOR CONDITIONAL PROBABILITIES ARE GIVEN BY

$$P(s_i = k | s_{i-1} = l) = \frac{\alpha(s_{i-1} = l, s_i = k)}{\sum_{k' \in K} \alpha(s_{i-1} = l, s_i = k')}$$

- HIDDEN MARKOV MODELS

- $S = (s_1, s_2, \dots, s_n)$

- SET OF HIDDEN STATES

- SEQUENCE OF CHARACTERS FOR EXAMPLE

- "USED FOR LANGUAGE MODEL"

- $X = (x_1, x_2, \dots, x_m)$

- SEQUENCE OF FEATURES

- SEQUENCE OF IMAGES OF CHARACTERS FOR EXAMPLE

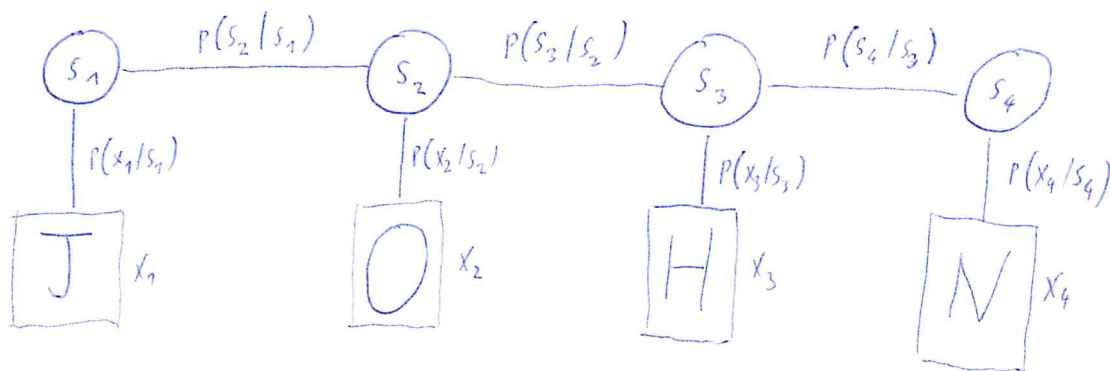
- "USED FOR APPEARANCE MODEL OF CHARACTERS"

- $P(s_i | s_{i-1})$

- LANGUAGE MODEL

- $P(x_i | s_i)$

- APPEARANCE MODEL



- FIND MOST PROBABLE SEQUENCE OF HIDDEN STATES GIVEN SEQUENCE OF FEATURES

$$S^* \in \underset{S \in K^m}{\text{ARGMAX}} P(s_1) \prod_{i=2}^m P(s_i | s_{i-1}) \prod_{i=1}^m P(x_i | s_i)$$

- LEARNING MODEL PARAMETERS FROM TRAINING DATA

- SAME AS FOR MARKOV MODELS

- MARKOV RANDOM FIELDS

- LET (V, E) BE UNDIRECTED GRAPH

- LET $S = \{s_i, i \in V\}$ BE FIELD OF RANDOM VARIABLES INDEXED BY THE NODES OF THE GRAPH AND TAKING VALUES FROM FINITE SET K .

- JOINT PROBABILITY DISTRIBUTION $p(s)$ IS GIBBS RANDOM FIELD ON GRAPH (V, E) IF IT FACTORISES OVER NODES AND EDGES

$$p(s) = \frac{1}{Z(w)} \exp \left[\sum_{i \in V} w_i(s_i) + \sum_{\{i, j\} \in E} w_{ij}(s_i, s_j) \right]$$

- PROBABILITY DISTRIBUTION $p(s)$ IS MARKOV RANDOM FIELD W.R.T. GRAPH (V, E) IF

$$p(s_A, s_B | s_C) = p(s_A | s_C) \cdot p(s_B | s_C)$$

HOLDS FOR ANY SUBSET $A, B \subset V$ AND SEPARATING SET C

- MULTICLASS LINEAR CLASSIFIER

- X

- SET OF OBSERVATIONS

- $Y = \{1, \dots, Y\}$

- SET OF CLASS LABELS

- $h(x; w) \in \operatorname{ARGMAX}_{y \in Y} \langle w_y, \phi(x) \rangle$

- $\phi: X \rightarrow \mathbb{R}^d$ IS FEATURE MAP

- $w = (w_1, \dots, w_Y) \in \mathbb{R}^{d \cdot Y}$ ARE PARAMETERS

- WE CAN REWRITE THE SCORING FUNCTION TO

$$\langle w_y, \phi(x) \rangle = \langle w, \phi(x, y) \rangle$$

- WHERE:

- $\phi: X \times Y \rightarrow \mathbb{R}^{d \cdot Y}$ IS

$$\phi(x, y) = (0, \dots, \underbrace{\phi(x)}_{y\text{-TH SLOT}}, \dots, 0)$$

- EXAMPLE

- OCR FOR SEQUENCES

- $X = (x_1, \dots, x_L) \in I^L$

- SEQUENCE OF IMAGES WITH CHARACTERS

- $Y = (y_1, \dots, y_L) \in A^L$

- SEQUENCE OF CHARACTERS FROM $A = \{A...Z\}$

- $P(x_i | y_i)$

- APPEARANCE MODEL FOR CHARACTERS

- $P(y_i | y_{i-1})$

- LANGUAGE MODEL

$$y \in \underset{y \in A^L}{\text{ARGMAX}} \left(P(y_1) \prod_{i=2}^L P(y_i | y_{i-1}) \prod_{i=1}^L P(x_i | y_i) \right)$$

- MAP

$$y \in \underset{y \in A^L}{\text{ARGMAX}} \left(\log P(y_1) + \sum_{i=2}^L \log P(y_i | y_{i-1}) + \sum_{i=1}^L \log P(x_i | y_i) \right)$$

- LET'S ASSUME FOLLOWING PARAMETRIZATION

$$\log P(y_1) = \langle w, \phi(y_1) \rangle$$

$$\log P(y_i | y_{i-1}) = \langle w, \phi(y_{i-1}, y_i) \rangle$$

$$\log P(x_i | y_i) = \langle w, \phi(x_i, y_i) \rangle$$

- MAP THEN BECOMES LINEAR CLASSIFIER

$$y = \underset{(y_1, \dots, y_L) \in A^L}{\text{ARGMAX}} \left\langle w, \phi(y_1) + \sum_{i=2}^L \phi(y_{i-1}, y_i) + \sum_{i=1}^L \phi(x_i, y_i) \right\rangle$$

$$\phi(x, y)$$

- TRAIN BY ERM

- CORRECTLY CLASSIFIED EXAMPLE

$$\langle \phi(x^i, y^i), w \rangle > \langle \phi(x^i, y), w \rangle \quad \forall y \in Y \setminus \{y^i\}$$

- PERCEPTRON ALGORITHM

- WE WANT TO FIND w SO THAT

$$\langle w, a^i \rangle > 0 \quad \forall i \in \{1, 2, \dots, k\}$$

- 1. $w \leftarrow 0$

- 2. FIND VIOLATING $\langle w, a^i \rangle \leq 0$

- 3. IF THERE IS NO VIOLATING

- RETURN w

OTHERWISE

$$w \leftarrow w + a^i$$

GOTO 2

- STRUCTURED PERCEPTRON

- LEARNING $\hat{y}(x; w) \in \text{ARGMAX}_{y \in Y} \langle w, \phi(x, y) \rangle$

- PAIR EXAMPLES $T_m = \{(x^i, y^i) \in (X \times Y) \mid i=1, \dots, m\}$

- IT LEADS TO SOLVING

$$\langle \phi(x^i, y^i), w \rangle - \langle \phi(x^i, y), w \rangle > 0 \quad \forall i \in \{1, \dots, m\}, y \in Y \setminus \{y^i\}$$

- 1. $w \leftarrow 0$

- 2. FIND MISCLASSIFIED EXAMPLE $(x^i, y^i) \in T_m$ SUCH THAT

$$y^i \neq \hat{y}^i \in \text{ARGMAX}_{y \in Y} \langle w, \phi(x^i, y) \rangle \quad \text{"PREDICTION PROBLEM"}$$

- 3. IF THERE IS NO MISCLASSIFIED EXAMPLE RETURN w

OTHERWISE

$$w \leftarrow w + \phi(x^i, y^i) - \phi(x^i, \hat{y}^i)$$

GOTO 2.

- STRUCTURED OUTPUT SVM

- APPROXIMATES EMPIRICAL RISK MINIMIZATION BY CONVEX PROBLEM

$$w^* \in \underset{w \in W_r}{\text{ARGMIN}} R^\Psi(w) \quad \text{WHERE} \quad R_{T^m}(w) = \frac{1}{m} \sum_{i=1}^m \frac{\ell(\gamma_i, \phi(x_i^i, w))}{\Psi(x_i^i, \gamma_i^i, w)}$$

- WHERE

- $W_r \subseteq \mathbb{R}^n$ - CONVEX FEASIBLE SET, E.G. $W_r = \{w \in \mathbb{R}^n \mid \|w\| \leq r\}$

- $\Psi: X \times Y \times \mathbb{R}^n \rightarrow \mathbb{R}$ - CONVEX PROX + APPROXIMATING

THE TRUE LOSS $\ell(\gamma^i, \phi)$

$$\langle w, \phi(x_i^i, \gamma^i) \rangle \geq \langle w, \phi(x_i^i, \phi) \rangle + \ell(\gamma^i, \phi) \quad \forall \phi \in Y \setminus \{\gamma^i\}$$

- MARGIN RESCALING LOSS

$$\Psi(x_i^i, \gamma_i^i, w) = \max \left\{ 0, \max_{\phi \in Y \setminus \{\gamma^i\}} \{ \ell(\gamma^i, \phi) + \langle w, \phi(x_i^i, \phi) \rangle - \langle w, \phi(x_i^i, \gamma^i) \rangle \} \right\}$$

- SO-SVM

- CONVEX CONSTRAINED OPTIMIZATION PROBLEM

$$w^* \in \underset{w \in W_r}{\text{ARGMIN}} R^\Psi(w)$$

$$R^\Psi(w) = \frac{1}{m} \sum_{i=1}^m \max_{\phi \in Y} \{ \ell_i(\phi) + \langle w, \phi_i(\phi) \rangle \}$$

- $W_r = \{w \in \mathbb{R}^n \mid \|w\| \leq r\}$ IS CONVEX SET

- $\ell_i(\phi) = \ell(\gamma_i^i, \phi)$, $\phi_i(\phi) = \phi(x_i^i, \phi) - \phi(x_i^i, \gamma_i^i)$

- SO-SVM UNCONSTRAINED PROBLEM

$$w^* \in \underset{w \in \mathbb{R}^n}{\text{ARGMIN}} \left(\frac{\lambda}{2} \|w\|^2 + R^\Psi(w) \right)$$

- WITH SLACK VARIABLES

$$w^* = \underset{w \in \mathbb{R}^n, f \in \mathbb{R}^m}{\text{ARGMIN}} \left(\frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_{i=1}^m f_i \right)$$

$$\text{s.t. } f_i \geq \ell_i(\gamma) + \langle w, \phi_i(\gamma) \rangle \quad \forall i \in \{1, \dots, m\} \forall \gamma \in Y$$

- CUTTING PLANE ALGORITHM

- 1. $w_1 \in W$, $\epsilon \leftarrow 1$

- 2. COMPUTE NEW CUTTING PLANE AND OBJECTIVE VALUE

$$\alpha_\epsilon = \frac{1}{m} \sum_{i=1}^m \phi_i(x^i)$$

$$b_\epsilon = \frac{1}{m} \sum_{i=1}^m l_i(x^i)$$

$$R^\Psi(w_\epsilon) = b_\epsilon + \langle w_\epsilon, \alpha_\epsilon \rangle$$

WHERE x^i IS SOLUTION OF LOSS AUGMENTED PREDICTION PROBLEM:

$$x^i = \underset{x \in X}{\operatorname{ARGMAX}} (l_i(x) + \langle w_i, \phi_i(x) \rangle) = \underset{x \in X}{\operatorname{ARGMAX}} (l_i(x, y) + \langle w_i, \phi(x, y) \rangle)$$

- 3. SOLVE REDUCED PROBLEM

$$w_{\epsilon+1} = \underset{w \in W}{\operatorname{ARGMIN}} R^\Psi(w), \quad \text{WHERE } R^\Psi(w) = \max_{i=1, \dots, \epsilon} (b_i + \langle w, \alpha_i \rangle)$$

- 4. IF $\min_{i=1, \dots, \epsilon+1} R(w_\epsilon) - R^\Psi(w_{\epsilon+1}) \leq \epsilon$

- EXIT

ELSE

$\epsilon \leftarrow \epsilon + 1$ AND GO TO 2

- EMSE MODELING

- "WISDOM OF CROWD"
- AVERAGING OR TAKING MAJORITY VOTE
- CANCELING EFFECT OF NOISE OF INDIVIDUAL OPINIONS

- $g_m(x) = E_{T_m} (h_m(x))$

- PREDICTOR AVERAGED OVER MULTIPLE DATASETS

- ERROR = BIAS² + VARIANCE + NOISE

- APPROACHES "OFF-CENTER" "SPREAD"

- BAGGING

- BOOTSTRAP AGGREGATING

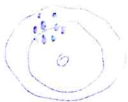
- SAMPLE DIFFERENT TRAINING SETS FROM ORIGINAL TRAINING DATA

- TRAIN PREDICTORS ON THESE SETS AND AVERAGE THEM

- HIGH VARIANCE LOW BIAS  "IN THE CENTER, BUT UNACCURATE"

- BOOSTING

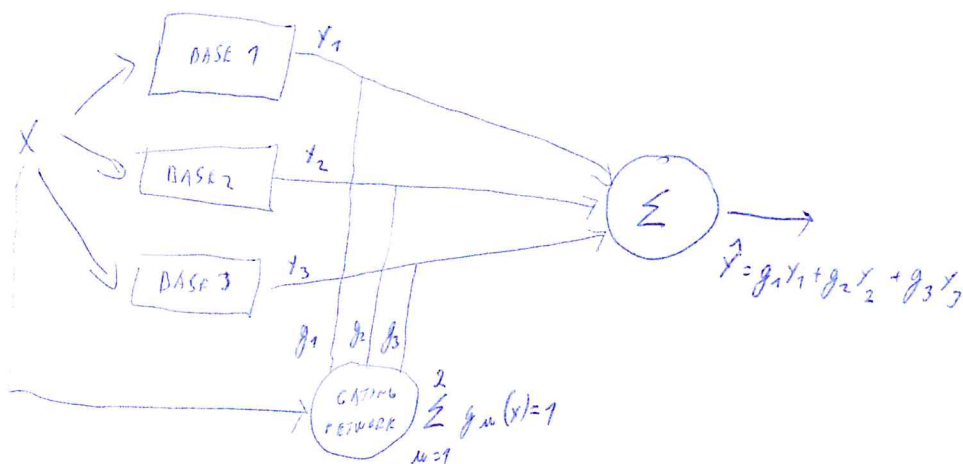
- SEQUENTIALLY TRAIN LOW VARIANCE HIGH BIAS PREDICTORS

 - "ACCURATE BUT OFF-CENTER"

- SUBSEQUENT PREDICTORS LEARN FROM MISTAKES OF PREVIOUS ONE

- INCREASE WEIGHTS FOR MISCLASSIFIED SAMPLES

- COMBINATION OF BASE AND META LEARNERS



- DECISION / REGRESSION TREES

- TRAINING SET

$$T^m = \{(x_i, y_i) \mid i = 1 \dots m\} \quad x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$$

- INPUT SPACE SPLIT INTO REGIONS DEFINED IN LEAVES

$$R_r, r \in \{1, \dots, M\}$$

- WE CAN MODEL REGION RESPONSES BY CONSTANTS $c_r, r \in \{1, \dots, M\}$

- PREDICTION

$$h(x) = \sum_{r=1}^M c_r \mathbb{I}[x \in R_r]$$

- FOR SUM OF SQUARES LOSS FUNCTION $\sum_{i=1}^m (y_i - h(x_i))^2$

- WE SET RESPONSES TO BE AVERAGES OVER REGIONS

$$c_r = \frac{1}{|S_r|} \sum_{x_i \in R_r} y_i$$

$$S_r = \{(x_i, y_i) : (x_i, y_i) \in T^m \wedge x_i \in R_r\}$$

- BUT NP-COMPLEX LEARNING

- WE NEED HEURISTICS

- GREEDY APPROACH

- CHOOSE RECURSIVELY THE MOST IMPORTANT ATTRIBUTE

- IMPORTANT IS THE ATTRIBUTE THAT REDUCE LOSS THE MOST

- THEN SPLIT ATTRIBUTE j IN SPLIT POINT s INTO TWO PARTS (RECURSIVE REGION)

- ORDINAL ("CONTINUOUS") ATTRIBUTE:

$$R_L(j, s) = \{x \mid x \in R \wedge x_j \leq s\} \quad R_R(j, s) = \{x \mid x \in R \wedge x_j > s\}$$

- NOMINAL ("DISCRETE") ATTRIBUTE:

$$R_L(j, s) = \{x \mid x \in R \wedge x_j = s\} \quad R_R(j, s) = \{x \mid x \in R \wedge x_j \neq s\}$$

- HOW TO CHOOSE ATTRIBUTE j AND SPLIT POINT s

$$\text{ARGMIN}_{(j,s)} \left(\text{MIN}_{c_L} \left(\sum_{x_i \in R_L(j,s)} (y_i - c_L)^2 \right) + \text{MIN}_{c_R} \left(\sum_{x_i \in R_R(j,s)} (y_i - c_R)^2 \right) \right)$$

- IMPER OPTIMIZATION IN REGION

- BY AVERAGING DESIRED OUTPUT

$$\hat{c}_L = \frac{1}{|S_L(j,s)|} \sum_{x_i \in R_L(j,s)} y_i$$

$$\hat{c}_R = \frac{1}{|S_R(j,s)|} \sum_{x_i \in R_R(j,s)} y_i$$

- BIAS AND VARIANCE

- SMALL CHANGE IN DATA MAY LEAD TO BIG DIFFERENCE IN FINAL TREES

- GROWTH DECISION TREES HAS USUALLY

- LOW BIAS

- HIGH VARIANCE

- WHICH LEADS TO OVERFITTING

- IDEA IS TO AVERAGE MULTIPLE MODELS

- BOOTSTRAPING

- WAY HOW TO MAKE SEVERAL DATASETS FROM SINGLE DATASET

- NEW DATASETS HAS THE SAME SIZE AS ORIGINAL

- SAMPLING FROM T^m WITH REPLACEMENT

- NEW DATASETS T_i^m HAS $\approx 63,2\%$ ~~number~~ OF UNIQUE SAMPLES OF ORIGINAL DATASET

- BAGGING

- BOOTSTRAP AGGREGATING

- USE BOOTSTRAPING TO MAKE SEVERAL NEW DATASETS

- TRAIN MODELS ON DATASETS

- AVERAGE MODELS

- IN CASE OF DECISION TREES

- RANDOM FOREST

- GROWS TREES TO MAXIMA DEPTH TO PRODUCE BIAS

- TREES ARE DECORRELATED BY

- TRAINING EACH ON DIFFERENT BOOTSTRAP DATASET

- RANDOMIZATION OF SPLIT ATTRIBUTE SELECTION

- SELECT RANDOM ~~ATTRIBUTE~~ SET OF
ATTRIBUTES ~~TO~~ CONSIDERED FOR SPLITTING

- AVERAGING FOR REGRESSION

- MAJORITY VOTE FOR CLASSIFICATION

- OUT OF BAG ERROR

- HOW TO TEST ERRORS

- MEASURE ERRORS ONLY OF TREES, WHICH WERE NOT
TRAINED ON SAMPLE (x_i, y_i)

- "K-FOLD CROSS-VALIDATION"

- BOOSTING

- SEQUENTIALLY TRAIN WEAK LEARNERS

- LOW VARIANCE HIGH BIAS

- SUBSEQUENT PREDICTORS FIX MISTAKES OF PREVIOUS ONES REDUCING BIAS

- ADA-BOOST

- BINARY CLASSIFIER

- 1. INITIALIZE WEIGHTS OF TRAINING SAMPLES

$$w_i = \frac{1}{m} \quad i=1, 2, \dots, m$$

- 2. FOR $k=1$ TO K :

- a) FIT CLASSIFIER $f_k(x; \theta_k)$ TO TRAINING DATA USING LOSS WEIGHTED BY w_i :

$$\theta_k = \underset{\theta}{\text{ARGMIN}} \sum_{i=1}^m w_i \mathbb{I}[y_i \neq f_k(x_i; \theta)]$$

- b) COMPUTE WEIGHTED ERROR RATE

$$\epsilon_k = \frac{\sum_{i=1}^m w_i \mathbb{I}[y_i \neq f_k(x_i; \theta_k)]}{\sum_{i=1}^m w_i}$$

- c) COMPUTE WEIGHT

$$\alpha_k = \log\left(\frac{1 - \epsilon_k}{\epsilon_k}\right)$$

- d) SET $w_i \leftarrow w_i \cdot \exp(\alpha_k \cdot \mathbb{I}[y_i \neq f_k(x_i; \theta_k)])$ FOR $i=1, 2, \dots, m$

- 3. RETURN $h_m(x) = \text{SIGN} \left[\sum_{k=1}^K \alpha_k f_k(x; \theta_k) \right]$

- FORWARD STAGEWISE ADDITIVE MODELING

- FSAM

-1. INITIALIZE $f_0(x) = 0$

-2. FOR $k=1$ TO K :

-a) FIND

$$(\beta_k, \theta_k) = \underset{\beta, \theta}{\text{ARGMIN}} \sum_{i=1}^m \ell(y_i, f_{k-1}(x_i) + \beta b(x_i; \theta))$$

WHERE $b(x_i, \theta_k)$ IS BASIS FUNCTION AND β_k IS CORRESPONDING COEFFICIENT

-b) SET $f_k(x) = f_{k-1}(x) + \beta_k b(x; \theta_k)$

-3. RETURN $h_m(x) = f_k(x)$

- FSAM UPDATE IS SIMILAR TO GRADIENT DESCENT

$$f_k(x) = f_{k-1}(x) + \beta_k b(x; \theta_k)$$

↑ STEP SIZE ↑ NEGATIVE OF GRADIENT

- GRADIENT BOOSTING MACHINE

-1. INITIALIZE $f_0(x) = 0$ OR $f_0(x) = \underset{\gamma}{\text{ARGMIN}} \sum_{i=1}^m \ell(y_i, \gamma)$

-2. FOR $k=1$ TO K :

-a) COMPUTE

$$g_k = \left[\frac{\partial \ell(y_i, f_{k-1}(x_i))}{\partial f_{k-1}(x_i)} \right]_{i=1}^m$$

-b) FIT REGRESSION MODEL $b(\cdot; \theta)$ TO $-g_k$

$$\theta_k = \underset{\theta}{\text{ARGMIN}} \sum_{i=1}^m [(-g_k)_i - b(x_i; \theta)]^2$$

-c) CHOOSE FIXED STEP SIZE $\beta_k = \beta \in (0, 1]$ OR USE LINE SEARCH

$$\beta_k = \underset{\beta > 0}{\text{ARGMIN}} \sum_{i=1}^m \ell(y_i, f_{k-1}(x_i) + \beta b(x_i; \theta_k))$$

-d) SET $f_k(x) = f_{k-1}(x) + \beta_k b(x; \theta_k)$

-3. RETURN $h_m(x) = f_k(x)$

- GRADIENT BOOSTED TREES

- ALL WEAK LEARNERS ARE DECISION OR REGRESSION TREES
- LIMIT SIZE OF TREE
- META PARAMETERS (K, β) HAVE TO BE FOUND USING CROSS VALIDATION
(NUMBER OF TREES, LEARNING RATE)
- MODEL IS BUILT SEQUENTIALLY
- UNLIKE RANDOM FORESTS