

ADAPTIVE BASIS FUNCTION MODELS

KERNEL METHODS RELY TOO MUCH ON HAVING A GOOD KERNEL. LET'S LEARN RELEVANT FEATURES DIRECTLY FROM DATA.

$f(x) = w_0 + \sum_{m=1}^M w_m \phi_m(x)$, ϕ IS BASIS FUNCTION LEARNED FROM DATA. $\phi_m(x) = \phi(x, v_m)$ ARE PARAMETRIC. MODELS NON LINEAR ANYMORE.
ONLY LOCAL ESTIMATES OF MEAN AND MAP FOR θ WHILE PARAMETERS $\{w, m\}$

CLASSIFICATION AND REGRESSION TREES

RECURSIVELY PARTITION INPUT SPACE, DEFINING LOCAL MODEL IN EACH REGION. EACH LEAF IS A REGION. AXIS PARALLEL SPLITS. PIECEWISE CONSTANT SURFACES.

• GENERAL ALGORITHM

OPTIMAL PARTITIONING IS NP-COMPLEX. \rightarrow GREEDY ALGOS. CART, C4.5, ID3. SPLIT FUNCTION CHOOSES BEST FEATURE AND BEST VALUE FOR FEATURE TO SPLIT ON.

• STOPPING HEURISTICS \rightarrow IS RESIDUAL COST TOO SMALL? $>$ MAX DEPTH? RESPONSE DISTRIBUTIONS SUFFICIENTLY HOMOGENEOUS? N CASES IN RESPONSE TOO SMALL?

• REGRESSION COST! $COST(D) = \sum (y_i - \bar{y})^2$. $\bar{y} = \frac{1}{|D|} \sum y_i$ MEAN OF RESPONSE VARIABLE. ELSE FIT REGRESSION MODEL AT EACH LEAF WITH INDICES
ROOT \rightarrow LEAF AND MEASURE RESIDUAL ERROR.

• CLASSIFICATION COST: FIT A MULTICLASS BY ESTIMATING CLASS-CONDITIONAL PROBABILITIES $\hat{\pi}_c = \frac{1}{|D|} \sum 1(y_i = c)$

ERROR MEASURES:

• MISCLASSIFICATION RATE $\hat{y}_i = \text{argmax}_c \hat{\pi}_c \rightarrow \frac{1}{|D|} \sum 1(y_i \neq \hat{y}) = 1 - \hat{\pi}_{\hat{y}}$

• ENTROPY / DEVIANCE: $H(\hat{\pi}) = - \sum_{\hat{\pi}_c} \hat{\pi}_c \log \hat{\pi}_c$ MIN ENTROPY \leftrightarrow MAX INFORMATION GAIN $H(y) - H(y|x, \hat{c})$
 $\hat{\pi}_c$ IS MLE

• GINI INDEX: $\sum \hat{\pi}_c (1 - \hat{\pi}_c) = 1 - \sum \hat{\pi}_c^2$ EXPECTED ERROR RATE. BEHAVES SIMILARLY TO X-ENTROPY
LEAD TO MORE PURE NODES

• PRUNING EARLY STOPPING IS MYOIC. GROW A FULL TREE THEN PRUNE. PRUNE BRANCHES GIVING LEAST INCREASE IN ERROR.
PICK SMALLEST TREE WHOSE CV ERROR IS WITHIN 1 STDDEV OF MINIMUM.

• PROS:

- EASY TO INTERPRET
- HANDLE MIXED TYPES OF INPUTS
- AUTO VAR. SELECTION
- ROBUST TO OUTLIERS
- SCALE WELL
- CAN MISSING INPUTS

• CONS:

- NOT VERY ACCURATE
- UNSTABLE BECAUSE BUILT HIERARCHICALLY
 \rightarrow ERRORS PROPAGATE TOP TO BOTTOM

RANDOM FORESTS

LET'S REDUCE VARIANCE BY AVERAGING MANY ESTIMATES

BAGGING: TRAIN M TREES ON DIFFERENT DATA SUBSETS RANDOM W/ REPLACEMENT. $f(x) = \sum_{m=1}^M f_m(x)$. BOOTSTRAP AGGREGATING \rightarrow HIGH-CORR PREDICTORS

RANDOM FOREST: TREES ARE TRAINED ON RANDOM SUBSETS OF VARS AND DATA CASES. AWESOME.

• BAYESIANLY WE CAN PERFORM APPROX INFERENCE OVER SPACE OF TREES (STRUCTURE + PARAMS) OR ENSEMBLE OF TREES **BART**

• HIERARCHICAL MIXTURE OF EXPERTS

CAN BE ALTERNATIVE TO TREES. DIFFERENT EXPERT ON EACH PARTITION. ANY NESTED LINEAR DECISION BOUNDARIES. NOT JUST AXIS PARALLEL.
GLOBAL PREDICTION IS AVG OF ALL EXPERTS. CAN USE EM BECAUSE IT'S SMOOTH CONTINUOUS PROBLEM.

GENERALIZED ADDITIVE MODELS

$f(x) = \alpha + f_1(x) + \dots + f_p(x)$. f IS SCATTERPLOT SMOOTHER. MAPS TO $p(y|x)$ WITH LINK FUNCTIONS AS IN GLM.

BASIS FUNCTIONS ARE SMOOTHING/REGRESSION SPLINES. BACKFITTING ALGORITHM. IF X FULL RANK OBJECTIVE IS CONVEX

$O(NDT)$ • MARS: ALLOWS FOR INTERACTION EFFECTS, ANOVA MODEL. $f(x) = \beta_0 + \sum f_1(x_1) + \sum f_{12}(x_1, x_2) + \sum f_{123}(x_1, x_2, x_3)$

→ SIMILAR TO CART BUT WITH PIECEWISE LINEAR BASIS FUNCTION VS STEPS FUNCTIONS.

BACKFITTING: ITERATIVELY UPDATE f_j USING RESIDUALS $(-f_j)$ AS TARGET VECTOR
NORMALIZE AT 0 MEAN WITH EACH STEP "WEIGHTING IRLS"

MARS • GAM + INTERACTION + BASIS FCN + CART + SPLINES

BOOSTING

GREEDY APPROACH FOR FINDING ADAPTIVE BASIS FUNCTION MODELS. ϕ_m GENERATED BY WEAK/BASE LEARNER. m WLS APPLIED SEQUENTIALLY TO WEIGHTED DATA. EARLIER MISCLASSIFIED DATA → MORE WEIGHT. CAN DO ON ANY REG/CLASS ALSO BUT USUALLY ON CLASS. IT'S BESTEST. BOOSTING WORKS AS LONG AS WL DOES BETTER THAN CHANCE. VERY RESILIENT TO OVERFITTING. CAN BE SEEN AS GRADIENT DESCENT IN FUNCTION SPACE. CAN BE EXTENDED TO MANY LOSS FCNS.

BOOSTING PROBLEM: $\min_f \sum L(y, f(x_i))$: L IS LOSS, f IS ADM. TRUE CONDITIONALS CAN'T BE KNOWN. MINIMIZES THEIR EXPECTATION → POPULATION MINIMIZER

BOOSTING APPROXES LOGODDS RATIO, SEQUENTIALLY

$$f_0(x) = \underset{f}{\operatorname{ARGMIN}} \sum L(y, f(x, \gamma))$$

$$(f_m, \gamma_m) = \underset{f, \gamma}{\operatorname{ARGMIN}} \sum L(y, f_{m-1}(x_i) + \beta \phi(x, \gamma))$$

$$f_m(x) = f_{m-1}(x) + \beta_m \phi(x, \gamma_m)$$

• NO GOING BACK TO OPTIMIZE EARLIER PARAMS FORWARD STAGewise ADDITIVE MODELING
• M ITERATIONS, MONITOR PROF ON ~~TRAINING~~ TESTSET AND STOP EARLY, AIC OR BIC

• IN PRACTICE: $f_m(x) = f_{m-1}(x) + \eta \beta_m \phi(x, \gamma_m)$, η STEPSIZE 0.001, η SHRINKAGE

L2 BOOSTING: SQUARE ERROR LOSS, $L(\dots) = ((y_i - f_{m-1}(x_i)) - \phi(x, \gamma))^2$, $\beta = 1$ WLOG. WEAK LEARNER PREDICTS R AND FINDS NEW BASIS FCN
RESIDUAL

ADABOOST BINARY CLASSIF, EXPONENTIAL LOSS. $L_m(\phi) = \sum w_m \exp(-y_i \phi(x_i))$, $w_m = \exp(-\bar{y}_i f_{m-1}(x_i))$ WEIGHT 4 DATAPoint, $\phi_m = \underset{\phi}{\operatorname{ARGMIN}} \sum w_m (y \neq \phi)$
WEAK LEARNER TO WEIGHTED DATASET. SOLVE FOR $\beta_m = \frac{1}{2} \log \frac{1 - \text{ERR}}{\text{ERR}} \rightarrow f_m(x) = f_{m-1}(x) + \beta_m \phi_m(x)$
COMPLEX DECISION BOUNDARIES SUPER RESILIENT TO OVERFITTING.

LOGIT BOOST! EXP LOSS IS TOO HEAVY ON MISCLASSIFICATIONS. SENSITIVE TO OUTLIERS. USE LOGLOSS → MISTAKES HAVE LINEAR PENALTY. CAN EXTRACT PROBS.
 $L(\phi_m) = \sum \log[1 + \exp(-2\bar{y}_i (f_{m-1}(x) + \phi(x_i)))]$ NEWTON UPDATES

GRADIENT BOOSTING: 'GENERIC' MODEL 4 BOOSTING $\hat{f} = \underset{f}{\operatorname{ARGMIN}} L(f)$. USUAL STAGewise. $\gamma_m = \left[\frac{\partial L(y, f(x_i))}{\partial f(x_i)} \right]$. $f = f_{m-1} - \epsilon_m \gamma_m$. ϵ_m STEPSIZE SO TO $\phi = \underset{\phi}{\operatorname{ARGMIN}} L(f)$
→ FIT A WEAK LEARNER $\gamma_m = \underset{\gamma}{\operatorname{ARGMIN}} \sum (-\gamma_{i,m} - \phi(x, \gamma))^2$. UPDATE $f_m(x) = f_{m-1}(x) + \eta \phi(x, \gamma_m)$. FLEXIBLE

SPARSE BOOSTING: PICK WL THAT BEST PREDICTS RESIDUAL VECTOR. FORWARD STAGewise LINEAR REGRESSION. $\beta_m = \beta_{m-1} + \eta (0, 0, \dots, \beta_{JM}, \dots, 0, 0)$
 $\eta \rightarrow 0$ IS LAR. INCREASING NO STEPS = REDUCING REGULARIZATION λ , + VARIABLE DELETION → LARS

MARS: MULTIVARIATE ADAPTIVE REGRESSION TREES. GRADIENT BOOSTING + SHALLOW DECISION TREES/STUMPS. SHALLOW DE VARIANCE IS LOW AND BIAS IS FIXED WITH MORE ROWS

• BOOSTING CAN BE SEEN AS ℓ_1 REGULARIZATION. ELIMINATES 'IRRELEVANT' FEATS. ADABOOST + $\ell_1 = \text{L1-ADABOOST}$: GREEDILY ADDS BEST FEATS, THEN PRUNES WITH ℓ_1
→ ALSO BOOSTING MAXIMIZES MARGIN ON TRAINING DATA

• BAYESIAN BOOSTING IS Akin TO MIXTURE OF EXPERTS $p(y|x, \theta) = \sum \pi_m p(y|x, \gamma_m)$, EACH EXPERT IS WEAK LEARNER.

FEED FORWARD NEURAL NETWORKS - MLPs

SEQUES OF LOGISTIC MODELS STACKED ON TOP OF EACH OTHER. FINAL LAYER IS LOGISTIC OR LINEAR.

$$P(y|x, \theta) = N(y|w^T z(x), \sigma^2)$$

$$z(x) = g(Vx) = [g(V_1^T x) \dots g(V_H^T x)]$$

g = NONLINEARITY, ACTIVATION FCN

W = WEIGHT VECTOR HIDDEN \rightarrow OUTPUT

$z(x)$ = HIDDEN LAYER

IF g LINEAR \rightarrow COLLAPSE TO $y = W^T(Vx)$

H = NO. OF HIDDEN

UNIVERSAL APPROXIMATORS: ANY FCN GIVEN ENOUGH

V = WEIGHT MATRIX INPUT \rightarrow HIDDEN

NONLINEARITIES = SIGM (OUTPUT, BIN) SOFTMAX (OUTPUT MULTICLASS)
TANH (HIDDEN)

$|N|$ TRAINING $\approx 10 \times |W|$ NO WEIGHTS

- BINARY CLASSIF: $P(y|x, \theta) = \text{BER}(y|\text{SIGM}(W^T z(x)))$ - LINE GLM

- MULTI CLASSIF: VIA MUTUAL INFORMATION / SUM TO ONE / SOFTMAX
 $P(y|x, \theta) = \text{CAT}(y|\text{SIG}(Wz(x)))$

- REGRESSION: $P(y|x, \theta) = N(y|W\phi(x, V), \sigma^2)$

CONVNETS

PURPOSE OF HIDDEN \rightarrow LEARN NONLINEAR COMBINATIONS OF INPUT. FEATURE EXTRACTION. FEATS USED BY FINAL GLM. GOOD WHERE ORIGINAL INPUTS AREN'T VERY INFORMATIVE.
IN CONVNETS HIDDEN HAVE RECEPTIVE FIELDS AND PARAMETERS ARE TIED / SHARED PER LAYER. EXHIBIT TRANSITIONAL INVARIANCE.
CONV LAYERS CREATE FEATURE MAPS. INCREASE RESILIENCE BY AUGMENTING DATASET WITH DISTORTED VERSION OF ORIGINAL DATA.

LENET5 INTRODUCES SUBSAMPLING BY POOLING / AVERAGING / AVERAGE BETWEEN CONV. LAYERS. FINAL LAYER IS RAF NET \rightarrow VS. SIGM OR SOFTMAX
 \rightarrow SHIFT INVARIANCE

OTHERS

• SKIP ARCS \rightarrow DIRECT ARCS BETWEEN IN AND OUT

• RECURRENT NETS \rightarrow FEEDBACK CONNECTIONS, NONLINEAR DYNAMICAL SYSTEM

• HOPFIELD NETWORK / ASSOCIATIVE MEMORY IF WE ALLOW SYMMETRIC CONNECTIONS BETWEEN HIDDEN \approx BOLTZMANN MACHINE (PROBABILISTIC)

BACKPROPAGATION

MLP OF MLP IS NONCONVEX. WE HAVE TO USE GRADIENT-BASED OPTIMIZERS LIKE SGD BECAUSE ONLINE.

MLP FORWARD MODEL: $x_n \xrightarrow{V} a_n \xrightarrow{g} z_n \xrightarrow{W} b_n \xrightarrow{h} \hat{y}_n$; $\theta = (V, W)$

ALGO: $\Delta \theta \rightarrow \Delta H \rightarrow W_{2,1:n}, W_{1,1:n}$

• COMPUTE LOG LIKELIHOODS

REGRESSION = SQ. ERROR, CLASSIF = X-ENTROPY
 $J(\theta) = -\sum_n \sum_k (\hat{y}_{nk}(\theta) - y_{nk})^2$ | $J(\theta) = -\sum_n \sum_k y_{nk} \log \hat{y}_{nk}(\theta)$

• OUTPUT LAYER GRADIENT: $\nabla_{w_{nk}} J = \frac{\partial J}{\partial b_{nk}} z_n = \delta_{nk}^w z_n$; $\delta_{nk}^w = (\hat{y}_{nk} - y_{nk})$

• INPUT LAYER GRADIENT: $\nabla_{v_{nj}} J = \frac{\partial J}{\partial a_n} = \delta_{nj}^v x_n$; $\delta_{nj}^v = \frac{\partial J}{\partial a_n} = \sum_k \frac{\partial J}{\partial b_{nk}} \cdot \frac{\partial b_{nk}}{\partial a_n} = \sum_k \delta_{nk}^w w_{nk} \cdot g'(a_n)$

CAN BE DONE VIA
LAYER 2 ERRORS BACK TO W
 \rightarrow LOCALLY

• FINAL: $\nabla_{\theta} J(\theta) = \sum_n [\delta_{nj}^v x_n, \delta_{nk}^w z_n]$

IDENTIFIABILITY

COLLAPSED. NONLINEARITIES ARE SYMMETRIC, ODD. 2^H SETTINGS FOR SIGN FLIPS + SWITCHING HIDDEN = TOTAL OF $H! 2^H$ EQUIVALENT PERMUTATIONS.
PLUS LOCAL MINIMA DUE TO NLL NONCONVEXITY

REGULARIZATION

EARLY STOPPING (STOP WHEN TEST ERROR STARTS TO RISE) • IMPOSE PRIOR ON PARAMS, THEN MAP. LINE $N(0, \sigma^2 I) \rightarrow l_2$ REG. WEIGHT DECAY
 \rightarrow NLL BECOMES $J(\theta) = -\sum_n \log P(y_n|x_n, \theta) + \frac{\alpha}{2} \left[\sum_n v_n^2 + \sum_n w_n^2 \right]$, $\nabla_{\theta} J(\theta) = \left[\sum_n \delta_{nj}^v x_n + \alpha v_j, \sum_n \delta_{nk}^w z_n + \alpha w_k \right]$

• STRONG REGULARIZATION \rightarrow WHO CARES IF MANY IT.

• STANDARDIZE INPUT FOR NICE PLOT BEHAVIOR

• SET H LARGE, THEN REGULARIZE

• SET α WITH CV OR E.D.

• SAME REG. PARAMS FOR 1 AND 2 LAYER \rightarrow NO INVARIANCE; SET DIFFERENT ONES FOR WEIGHTS AND BIASES PER LAYER

$P(\theta) = N(W|0, \frac{1}{\alpha} I) N(V|0, \frac{1}{\alpha} I) N(b|0, \frac{1}{\alpha} I) N(c|0, \frac{1}{\alpha} I)$ • ENCOURAGE SPARSITY: WITH l_1 OR AND OR AD-HOC

• ENCOURAGE PARAMS TO HAVE SIMILAR STATISTICAL PROPERTIES. $P(\theta)$ AS MIXTURE OF DIAGONAL GAUSSIANS \rightarrow SAME CLUSTER, SAME μ, σ^2
SOFT WEIGHT SHARING

SEMI-SUPERVISED EMBEDDING

ENCOURAGE HIDDEN TO ASSIGN SIMILAR OBJECTS TO SIMILAR REPRESENTATIONS. $L(\mathbf{f}_1, \mathbf{f}_2, s) = \begin{cases} \|\mathbf{f}_1 - \mathbf{f}_2\|^2 & s=1 \\ \max(0, m - \|\mathbf{f}_1 - \mathbf{f}_2\|^2) & s=0 \end{cases}$

$$\bullet \sum \text{NLL}(\mathbf{f}(x_i), y_i) + \lambda \sum L(\mathbf{f}(x_i), \mathbf{f}(x_j), s_{ij})$$

- PICK RANDOM UNLABELLED EXAMPLE, AND GRADIENT ALL; PICK RANDOM SIMILAR UNLABELLED EXAMPLES AND GRADIENT L_1 ; PICK RANDOM UNLABELLED WITH HIGH PROB DISSIMILAR AND OPTIMIZE GRADIENT L_0

BAYESIAN INFERENCE

INTEGRATING VS OPTIMIZING PARAMS IS STRONGER REGULARIZATION. BAYESIAN MODEL SELECTION FOR HYPERPARAMS IS DESIRABLE. UNCERTAINTY IN PARAM DESIRABLE FOR CERTAIN PROBLEMS. ONLINE INFERENCE FOR ONLINE LEARNING. LAPLACE APPROXIMATION, HYBRID MONTECARLO, VARIATIONAL BAYES, AND FASTER

- REGRESSION POSTERIOR: $P(\mathbf{w} | \mathbf{d}, \beta, D) \approx N(\mathbf{w} | \mathbf{w}_{MP}, \mathbf{A}^{-1})$; \mathbf{A} IS HESSIAN OF $E_{\text{DATA}} = \nabla \nabla (\mathbf{w}_{MP}) = \beta \mathbf{I} + \mathbf{A}$

- CLASSIF POSTERIOR: SAME FOR REG BUT $\beta=1$ AND E IS CROSS-ENTROPY

- REGRESSION PRED. POSTERIOR: $P(y | x, \mathbf{d}, \beta, D) \approx N(y | \mathbf{f}(x, \mathbf{w}_{MP}), \sigma^2(x))$; $\sigma^2(x) = \beta^{-1} + \mathbf{g}^T \mathbf{A}^{-1} \mathbf{g}$

- CLASSIF PRED. POSTERIOR: $P(y=1 | x, D) = \int \text{SIGM}(u) P(\mathbf{d} | x, D) d\mathbf{u} \approx \text{SIGM}(u(\sigma_u^2) \mathbf{b}^T \mathbf{w}_{MP})$, $u(\sigma_u^2) \approx (1 + \pi \sigma_u^2 / 8)^{-1/2}$

- ARD ONCE LAPLACE APPROX, CAN OPTIMIZE MARGINAL LIKELIHOOD WRT HYPERPARAMS α . ONE \times WEIGHT VECTOR. EFFECT SIMILAR TO GROUP LASSO
CAN PRUNE OUT IRRELEVANT INPUT OR HIDDEN FEATURES. $P(\theta) = \prod_{d=1}^D N(\mathbf{v} | \mathbf{0}, \frac{1}{\alpha_d} \mathbf{I}) \prod_{w=1}^H N(\mathbf{w} | \mathbf{0}, \frac{1}{\alpha_w} \mathbf{I})$

RADIAL BASIS NETWORKS

MLP w/ RBF AS HIDDEN NODES. OUTPUT AS USUAL

TRAINING: 1- POSITION RBF NODES (RANDOM, K-MEANS, OTHER CLUSTERING) UNSUPERVISED

2 - FWD PASS

3 - UPDATE WEIGHTS (PERCEPTRON OR PSEUDOINVERSE) $G = (G^T G)^{-1} G^T$

• CAN DO REGRESSION. FCN APPROX, TIME SERIES PREDICTION, ...

SELF-ORGANIZING MAPS

IS AN ANN-LIKE THING USED FOR DIMENSIONALITY REDUCTION AND VISUALIZATION. UNSUPERVISED LEARNING VIA VECTOR QUANTIZATION. INPUT SPACE \leftrightarrow MAP SPACE (2D/3D AND VARIOUS /)

• **FEATURE MAPPING:** NEARBY NEURONS MAP/ACTIVATE TO SIMILAR FEATURES. INPUTS FULLY CONNECTED TO EACH MAP NEURON.

• **ALGORITHM:** SELECT BEST-MATCHING-UNIT/NEURON WRT PRESENTED INPUT.

- UPDATE BMU WEIGHTS - UPDATE DAMPED WEIGHTS OF MAP-NEIGHBORING NEURONS

- REDUCE USING RNE/NEIGHBORHOOD RADIUS - UNTIL CONVERGENCE

- WEIGHT INIT IS RANDOM OR VIA PCA OF DATA

• **SELF ORGANIZATION:** SPACE BECS CLUSTALY ORGANIZED BY ONLY LOCAL INTERACTIONS HAPPENING

• VARIES ON MAP TOPOLOGY/CONNECTEDNESS. LASTING THE MAP GROW OVER CYCLES

• **NEURAL GAS:** SIMILAR, BUT FEATURE VECTORS MOVE IN SPACE.

ENSEMBLE LEARNING

LEARNING A WEIGHTED COMBINATION OF MODELS $f(y|x, \pi) = \sum w_m \cdot f_m(y|x)$. COMMITTEE METHOD. NEURAL NETS, BOOSTING, CAN BE SEEN AS ENSEMBLES

- **STACKING!** $\hat{w} = \underset{w}{\text{ARGMIN}} \sum_w L(y_i, \sum_{m=1}^M w_m \cdot \hat{f}_m(x))$
↓
PREDICTOR BY REMOVING (x_i, y_i) .

- **ERROR CORRECTING OUTPUT CODES!** ECC, MULTICLASS CLASSIFICATION. USE MORE BIT THAN $B = \log_2 C$ FOR CLASS LABEL. MAXIMIZE HAMMING DISTANCE
→ RESISTANCE TO BIT-FLIPPING ERROR $\hat{c}(x) = \underset{c}{\text{MIN}} \sum_b |c_b - \hat{p}_b(x)|$

- **BAYES MODEL AVERAGING!** USE WEIGHTED AVERAGE OF PREDICTIONS MADE BY EACH MODEL. $P(y|x, D) = \sum_{m=1}^M P(y|x, m, D) P(m|D)$
USUALLY INTERCHANGE, PICK FEW SAMPLE FROM POSTERIOR
• NOT EQUIVALENT TO ENSEMBLE LEARNING → HERE WE ENLARGE THE MODEL SPACE,
CONVEX COMBINATION OF BASE MODELS: $P(y|x, \pi) = \sum \pi_m P(y|x, m)$

INTERPRETATION, PEAF ANALYSIS

BUNCHBOX MODELS ARE NOT INTUITIVE TO INTERPRET.

- **PARTIAL DEPENDENCE PLOT:** $f(x_j)$ vs x_j , WITH OTHER PREDICTORS AVERAGED OUT. $f(x_j) = \frac{1}{N} \sum_{i=1}^N P(x_j, x_1, \dots, x_N)$. f_j IS ^{PREDICTED} RESPONSE.
- **RELATIVE IMPORTANCE OF PREDICTOR VARS:** 'MODEL FREE' VARIABLE SELECTION. IN TREES COUNT HOW OFTEN VARS ARE USED.