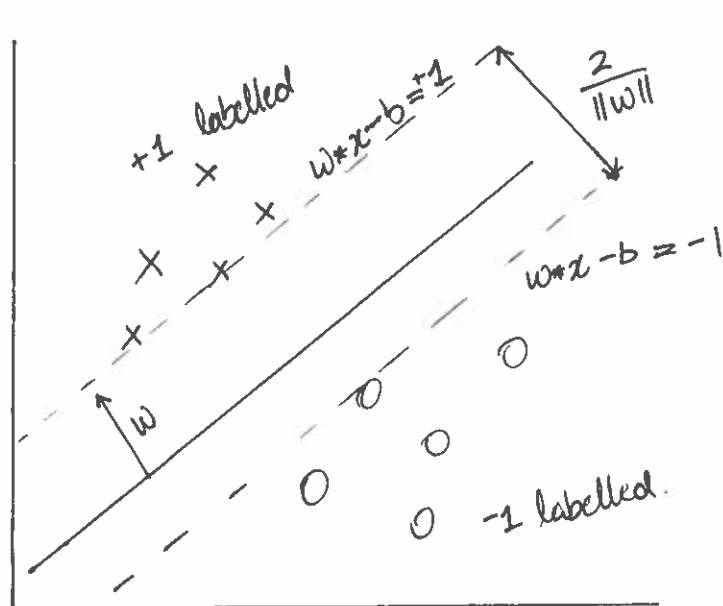


SVM.

Training data: $(\vec{x}_i, y_i) ; i=1 \dots n$

①



linearly separable data

Anything above $\vec{w}^T \vec{x} - b = 1$ is a "x" datapoint

$$\text{Maximize } \frac{2}{\|\vec{w}\|} \Rightarrow \text{Minimize } \|\vec{w}\|_2$$

\Rightarrow If $y_i = 1, \vec{w}^T \vec{x}_i - b \geq 1$
 $y_i = -1, \vec{w}^T \vec{x}_i - b \leq -1$ } Points must be correctly classified

$$\Rightarrow y_i (\vec{w}^T \vec{x}_i - b) \geq 1 \quad \forall 1 \leq i \leq n \quad \text{Constraint}$$

Optimization problem:

$$\min_{\vec{w}, b} \|\vec{w}\|_2^2$$

Quadratic optimization,
solve w/ Lagrange multipliers

If we remove support vectors,
hyperplane changes!

sub. $y_i (\vec{w}^T \vec{x}_i - b) \geq 1 \quad \forall i \in \{1 \dots n\}$

Quadratic \Rightarrow single global minima.

\uparrow hyperplane is determined completely by nearest \vec{x}_i on either side!

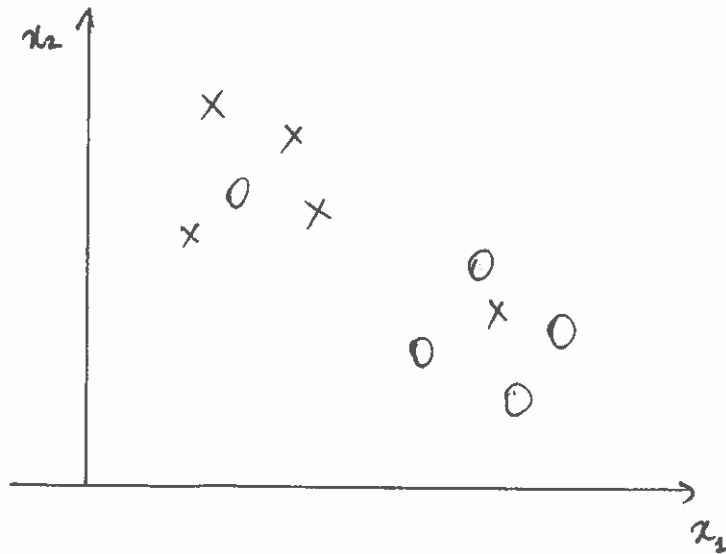
\Rightarrow Most difficult pts to classify

Use the optimization of maximizing margin to reduce number of nonzero weights ~~that~~ \Rightarrow only look at weights that matter \Rightarrow i.e. correspond to the support vectors!

Key diff b/w SVM & NNs/linear reg \Rightarrow use all data!

Ultimately, classifier: $\text{sign}(\vec{w}^T \vec{z} - b)$ for some new \vec{z} , using learned / ②
optimized \vec{w} , b .

"Fuzzy" data



Some points' labels are
 "misclassified".

Modify minimization: (using hinge loss)

$$\max(0, 1 - y_i(\vec{w}^T \vec{z}_i - b))$$

returns 0 if point is correctly classified/labelled.
 (0,1), proportional to distance from margin if
 incorrectly classified.

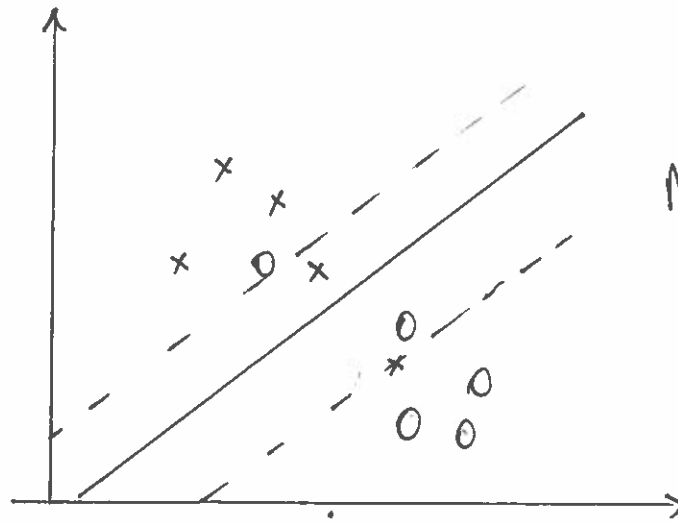
$$\Rightarrow \text{minimize } \lambda \|\vec{w}\|_2^2 + \left[\frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(\vec{w}^T \vec{z}_i - b)) \right]$$

Rewrite: $\min_{\vec{w}, b, \xi} \|\vec{w}\|_2^2 + C \sum_{i=1}^n \xi_i$

Riemann zeta
 function

sub. $y_i(\vec{w}^T \vec{z}_i - b) \geq 1 - \xi_i, \xi_i \geq 0 \quad \forall i \in \{1 \dots n\}$

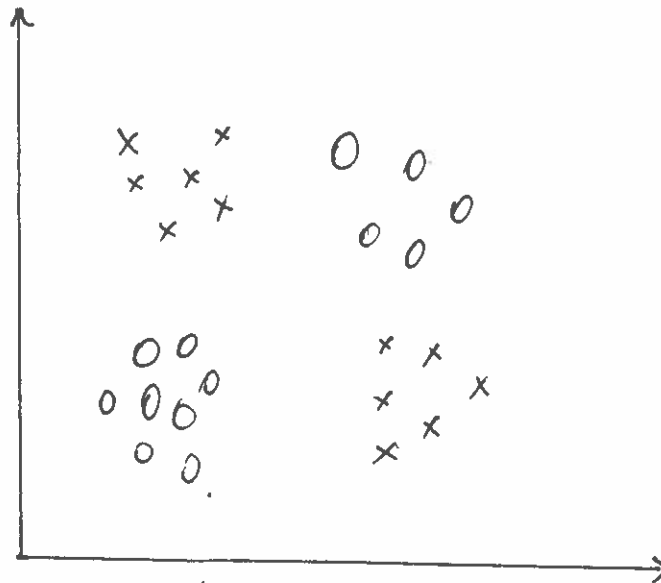
⇒ Our learned classifier now looks like:



Number of
support vectors = 4

Very similar to linear
discriminant analysis (LDA)

Non-linearly Separable Data:



XOR dataset

Gain linear separation by mapping data into a high-dimensional space

We actually solve SVM optimization via the dual formulation (removes dependence on \vec{w} , b).

Dual problem -

Interchange variables & constraints:

$$\begin{array}{ll} \min & c \cdot x \\ \text{sub.} & A \cdot x = b \\ & x \geq 0. \end{array}$$

$$\begin{array}{ll} \max & b \cdot y \\ \text{sub.} & A^T y \leq c. \\ & y \text{ Free. (Unconst).} \end{array}$$

Here:

$$\min \|w\|_2^2 \rightarrow \left. \begin{array}{l} \min. f(w) \\ \text{sub. } y_i (\vec{w}^T \vec{x}_i - b) \\ \quad \text{sub. } g_j(\vec{w}) \leq 0. \\ \quad h_k(\vec{w}) = 0. \end{array} \right\}$$

General form (lin sep. / fuzzy).

↓ Dual formulation

$$\mathcal{L} \equiv \frac{1}{2} f(\vec{w}) + \sum_j \alpha_j g_j(\vec{w}) + \sum_k \beta_k h_k(\vec{w})$$

0 ∴ no equality const.
(Method of Lagrange multipliers)

⇒ For a min/max:

$$\frac{\partial \mathcal{L}}{\partial w_i} = 0 \quad ; \quad \frac{\partial \mathcal{L}}{\partial \beta_k} = 0$$

From strong duality & Kuhn-Tucker thms.,

$$\max. \mathcal{L}(\alpha, \beta).$$

$$\text{sub. } \alpha_j \geq 0 \quad \forall j. \equiv \text{original primal form.}$$

Generally true \forall convex fns.
 $f(w)$

Further, if \hat{w} is optimal soln of primal & $\hat{\alpha}$ & $\hat{\beta}$ are optimal solns of dual:

$$\left. \begin{array}{l} f(\hat{w}) = \mathcal{L}(\hat{\alpha}, \hat{\beta}) \\ \hat{\alpha}_j g_j(\hat{w}) = 0 \quad \forall j \end{array} \right\} \text{Karush-Kuhn-Tucker (KKT) complementarity condition.}$$

If we sub. our original constraints for the SVM problem into the Lagrangian,

$$\mathcal{L} = \frac{1}{2} \vec{w} \cdot \vec{w} + \sum_i \alpha_i [1 - y_i (\vec{w} \cdot \vec{x}_i + b)]$$

For max/min:

$$0 = \frac{\partial \mathcal{L}}{\partial \vec{w}} = \vec{w} - \sum_i \alpha_i y_i \vec{x}_i \quad \text{and} \quad 0 = \frac{\partial \mathcal{L}}{\partial b} = \sum_i \alpha_i y_i$$

$$\Rightarrow \hat{w} = \sum_i \hat{\alpha}_i y_i \vec{x}_i$$

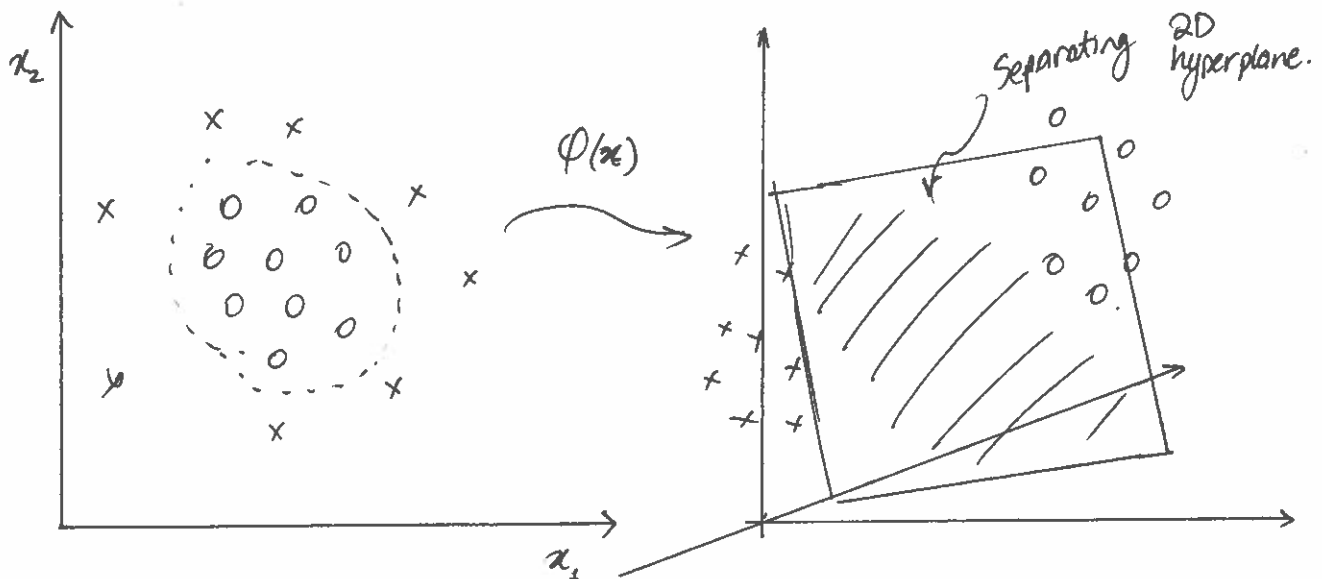
$$\Rightarrow \mathcal{L}(\alpha) = \sum_j \alpha_j - \frac{1}{2} \sum_{j,k} \alpha_j y_j (\vec{x}_j \cdot \vec{x}_k) y_k \alpha_k$$

The only place where \vec{x} shows up is here!

This is the only thing that scales w/ number of features. \Rightarrow

all other parts scale w/ num. pts. \Rightarrow Primal, to dual \rightarrow shift scaling from number of features to number of pts (dual!) \Rightarrow Favors data w/ \downarrow numbers of data pts., but ~~not~~ huge features

Now, let's talk about non linearity. \Rightarrow "Kernel Trick".



Primal form: $\min. \frac{1}{2} \vec{w} \cdot \vec{w} + \lambda \sum_i \xi_i$

sub. $y_i (\vec{w} \cdot \phi(\vec{x}_i) + b) \geq 1 - \xi_i \quad i = 1 \dots m.$

But, embedding dim may be $\mathbb{R}^{\text{Billions}}$!

\Rightarrow Intractable!

Dual problem:

~~$\min. \frac{1}{2} \alpha \cdot \text{diag}(y)$~~

$$\min. -\sum_j \alpha_j + \frac{1}{2} \sum_{j,k} \alpha_j y_j \alpha_k y_k K_{jk}$$

$$K_{jk} = K(x_j, x_k) = \phi(x_j) \cdot \phi(x_k) \quad \text{dot prod.}$$

But, we don't need to explicitly know mapping $\phi(x)$!

All we need is some way to compute K_{jk} that could have come from some $\phi(x)$! \Rightarrow $m \times m$ matrix w/ mathematical properties of inner product !

- $K_{ij} = K(x_i, x_j)$ must be symm. in i, j , and nonneg. eigenval.

- K can be composed by addition, mult., scaling by const.

Popular kernels:

linear: $K(x_i, x_j) = x_i \cdot x_j.$

power: $K(x_i, x_j) = (x_i \cdot x_j)^d \quad 2 \leq d \leq 20 \text{ (usually).}$

polynom: $K(x_i, x_j) = (a x_i \cdot x_j + b)^d$

sigmoid: $K(x_i, x_j) = \tanh(a x_i \cdot x_j + b).$

Gaussian rbf : $K(x_i, x_j) = e^{-\frac{1}{2} \frac{|x_i - x_j|^2}{\sigma^2}}$

Tips for using kernels:

- Gaussian rbf is very popular, \therefore only 1 hyperparam.

Guess good initial σ via $\frac{\text{avg}}{1}$ distance b/w points in feature space.

- Polynomial kernels:

Start by choosing a, b s.t. $a x_i = x_j + b$ b/w $-1, 1 \neq i, j$.

$d \Rightarrow$ roughly interpret as number of features to be mixed ~~to~~ during partitioning.

$\Rightarrow d=1 \Rightarrow$ partition space by 1 feature at a time

$2 \Rightarrow$ 2 features ...

Diff b/w power & polynomial: power considers only d features at once, polynomial considers all combos of d or fewer features.

Decision Tree Models

①

Start w/ training data, labels

$$x_i \in \mathbb{R}^n, i=1 \dots N \quad \text{data}$$

$$y_i \in \mathbb{R}^N \quad \text{labels.}$$

No index, just concat.

all labels into one big vector.

split

I.e., each ~~set~~ element of \mathbb{R}^N

We want to partition the space of features st. the x_i 's w/ the same y_i are grouped together.

Let's represent the data @ some node m w/ Q_m , w/ n_m samples.
 made up of a set of x_i, y labels

Then, represent each split as $\theta = (j, t_m)$
 \uparrow Feature (I.e., index of each x_i element)
 \uparrow Threshold value.

\Rightarrow Partition data into left & right subsets:

$$Q_m^{\text{left}}(\theta) = \{(x, y) \mid x_j \leq t_m\} \quad (\text{Below or equal to threshold})$$

$$Q_m^{\text{right}}(\theta) = Q_m \setminus Q_m^{\text{left}}(\theta) \quad (\text{everything else})$$

Then, define a quality score based on a loss function:

$$G(Q_m, \theta) = \frac{n_m^{\text{left}}}{n_m} H(Q_m^{\text{left}}(\theta)) + \frac{n_m^{\text{right}}}{n_m} H(Q_m^{\text{right}}(\theta))$$

Then, self select θ^* that min. G : $\theta^* = \underset{\theta}{\operatorname{argmin}} G(Q_m, \theta)$

To define H , define a measure for the proportion of class k observations in node m :

②

$$p_{mk} = \frac{1}{n_m} \sum_{y \in Q_m} \mathbb{I}(y=k)$$

\Rightarrow Gini loss: $H(Q_m) = \sum_k p_{mk} (1 - p_{mk})$

log loss: $H(Q_m) = -\sum_k p_{mk} \log(p_{mk})$

Random Forest: & XG Boost.

General idea: combine influence of several different decision tree models.

RF :

Build a forest where each tree uses a random subset of features \rightarrow low correlation b/w trees.

Sklearn implementation - use all features or random subset of features. "Random subspace method"

2) Each tree is built from a sample drawn w/ replacement from the training set.

Net effect: decrease variance of ensemble, reduce overfitting.
Combine diverse trees, at cost of increasing model bias.

\rightarrow If training data is $[x_1, x_2, x_3, x_4, x_5, x_6]$,
one tree gets $[x_1, x_1, x_3, x_4, x_4, x_4, \text{...}]$. \swarrow same size
 \uparrow \uparrow \uparrow
Sampled with replacement
 "Bagging" / "Bootstrapping"

Form a prediction w/ a consensus of trees.

Additive trees: classifier & regressor trees
& CART:

Group Members

Louis

Rob.

Sabari

Yeonjoon

Shree.

Interests

Likes GPUs

Yes

No.

Sabari
Yeonjoon
Shree

Louis
Rob.

Pred. score
+2

Pred score
0.5

Mh.

Organic

Rugs

Yes

XTB regularly

No.

Louis
Rob
Yeonjoon

Sabari
Shree

Pred score
+1.5

Pred
score
0.3

Organic

Mh

$$f(\text{Sabari}) = 2 + 0.3 = 2.3.$$

Mathematically, $\hat{y}_i = \sum_{k=1}^K f_k(x_i)$, $f_k \in \mathcal{F} \Rightarrow \mathcal{F}$ is set of all possible CARTS.

$K \Rightarrow$ total number of trees.

$$\Rightarrow \min_{\theta} \text{obj} = \sum_{i=1}^N l(y_i, \hat{y}_i) + \sum_{k=1}^K \omega(f_k)$$

Loss fn.
over data.

Regularization
over trees

\rightarrow Controls
tree
complexity

③.

What are parameters of trees \Rightarrow what's "f"?

\hookrightarrow Structure of tree + leaf scores!

This isn't a gradient optimization problem.

Learning on space of all possible trees is totally intractable!

Instead, additive predictions: Fix what we have learned, add one tree at a time.

$$\hat{y}_i^0 = 0$$

$$\hat{y}_i^1 = f_1(x_i) := \hat{y}_i^0 + f_1(x_i) \quad \leftarrow \text{First CART tree.}$$

$$\hat{y}_i^2 = \frac{f_1(x_i) + f_2(x_i)}{2} := \hat{y}_i^{(2)} + f_2(x_i)$$

$$\hat{y}_i^{(t)} = \sum_{k=1}^t f_k(x_i) := \hat{y}_i^{(t-1)} + f_t(x_i)$$

Into objective fn:

$$\text{obj}^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t)}) + \sum_{i=1}^t \omega(f_i)$$

$$= \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \omega(f_t) + \text{const}$$

\leftarrow We can only optimize structure of t^{th} tree!
Roll other ω 's into const.

If we apply MSE as the loss fn:

$$\text{obj}^{(t)} = \sum_{i=1}^n (y_i - (\hat{y}_i^{(t-1)} + f_t(x_i)))^2 + \omega(f_t) + \text{const}$$

$$= \sum_{i=1}^n \left(\underbrace{y_i^2}_{\text{const}} - 2y_i(\underbrace{\hat{y}_i^{(t-1)}}_{\text{const}} + f_t(x_i)) + \underbrace{(\hat{y}_i^{(t-1)})^2}_{\text{const}} + 2\hat{y}_i^{(t-1)}f_t(x_i) + (f_t(x_i))^2 \right) + \omega(f_t) + \text{const}$$

$$= \sum_{i=1}^n [2(\hat{y}_i^{(t-1)} - y_i)f_t(x_i) + (f_t(x_i))^2] + \omega(f_t) + \text{const}$$

$$\text{obj}^{(t)} = \sum_{i=1}^n \left[2(\hat{y}^{(t-1)} - y_i) f_x(x_i) + (f_x(x_i))^2 \right] + w(f_t) + c.$$

\uparrow 1st order term (residual) \uparrow 2nd order term

We got a nice form in the case of MSE; but for other loss fns this isn't the case! So, in general, use Taylor expansion:

$$\text{obj}^{(t)} = \sum_{i=1}^n \left[l(y_i, \hat{y}_i^{(t-1)}) + g_i f_x(x_i) + \frac{1}{2} h_i f_x^2(x_i) \right] + w(f_t) + c$$

where $g_i = \frac{\partial l(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}}$; $h_i = \frac{\partial^2 l(y_i, \hat{y}_i^{(t-1)})}{\partial^2 \hat{y}_i^{(t-1)}}$

} "gradient boosting"

Constants don't affect minimization procedure:

$$\Rightarrow \text{obj}^{(t)} = \sum_{i=1}^n \left[g_i f_x(x_i) + \frac{1}{2} h_i f_x^2(x_i) \right] + w(f_t)$$

XGBoost defines $w(f)$ as:

$$w(f) = \gamma T + \frac{1}{2} \gamma \sum_{j=1}^T \omega_j^2$$

γ Gamma T Total num of trees ω_j Vector of scores on each leaf
 ω_j^2 hyperparam

Rewrite & compress: $\text{obj}^{(t)} = \sum_{j=1}^T \left[G_j \omega_j + \frac{1}{2} (H_j + \lambda) \omega_j^2 \right] + \gamma T.$

SHAP Analysis

①

SHAP \Rightarrow "Shapley Additive Explanations"

Comes from cooperative game theory:

A number of players cooperate to achieve an objective, that leads to some overall gain.

But, some players may contribute more

One might have more bargaining power

Another might threaten to destroy everything

\Rightarrow Question: what amount of the overall gain should be assigned to each player?

I.e. how important is each player to the overall gain, and how much payoff should each player expect

John: Imagine the group being formed one ^{/feature} person at a time. Each person demands their contribution as "fair compensation".

Then, for each person, average contribution over the different ways we can form the group.

For ML: Assume only some features are present, while others aren't

>

Shapley values are only attribution method that satisfies:

1) Efficiency: Feature contributions sum to the diff of a prediction x and the average:

$$\sum_{j=1}^p \phi_j = \hat{f}(x) - \overset{\text{Expectation value.}}{E_x(\hat{f}(x))}$$

2) Symmetry: Contributions of features j & k are the same if they contribute equally to all possible groups:

$$\text{If } \text{val.}(S \cup \{j\}) = \text{val.}(S \cup \{k\}).$$

$$\forall S \subseteq \{1, \dots, p\} \setminus \{j, k\}$$

↑ Group w/o j, k .

$$\text{Then } \phi_j = \phi_k.$$

3) Dummy: A feature j that doesn't change predicted value, regardless of group, gets a value of 0.

$$\text{If } \text{val}(S \cup \{j\}) = \text{val}(S) \quad \forall S \subseteq \{1, \dots, p\}.$$

$$\phi_j = 0$$

Additivity: For a objective w/ combined payouts, ~~the~~ SHAP values are additive. (3)

$$\phi_j + \phi_j^+ \Rightarrow \text{combined value.}$$

This is important, since for RF/ensemble models, additivity guarantees that if we calc. SHAP val. for each indiv. tree & average \rightarrow SHAP for forest.