

CS 4260 and CS 5260

Vanderbilt University

Lecture on NBCs

This lecture assumes that you have

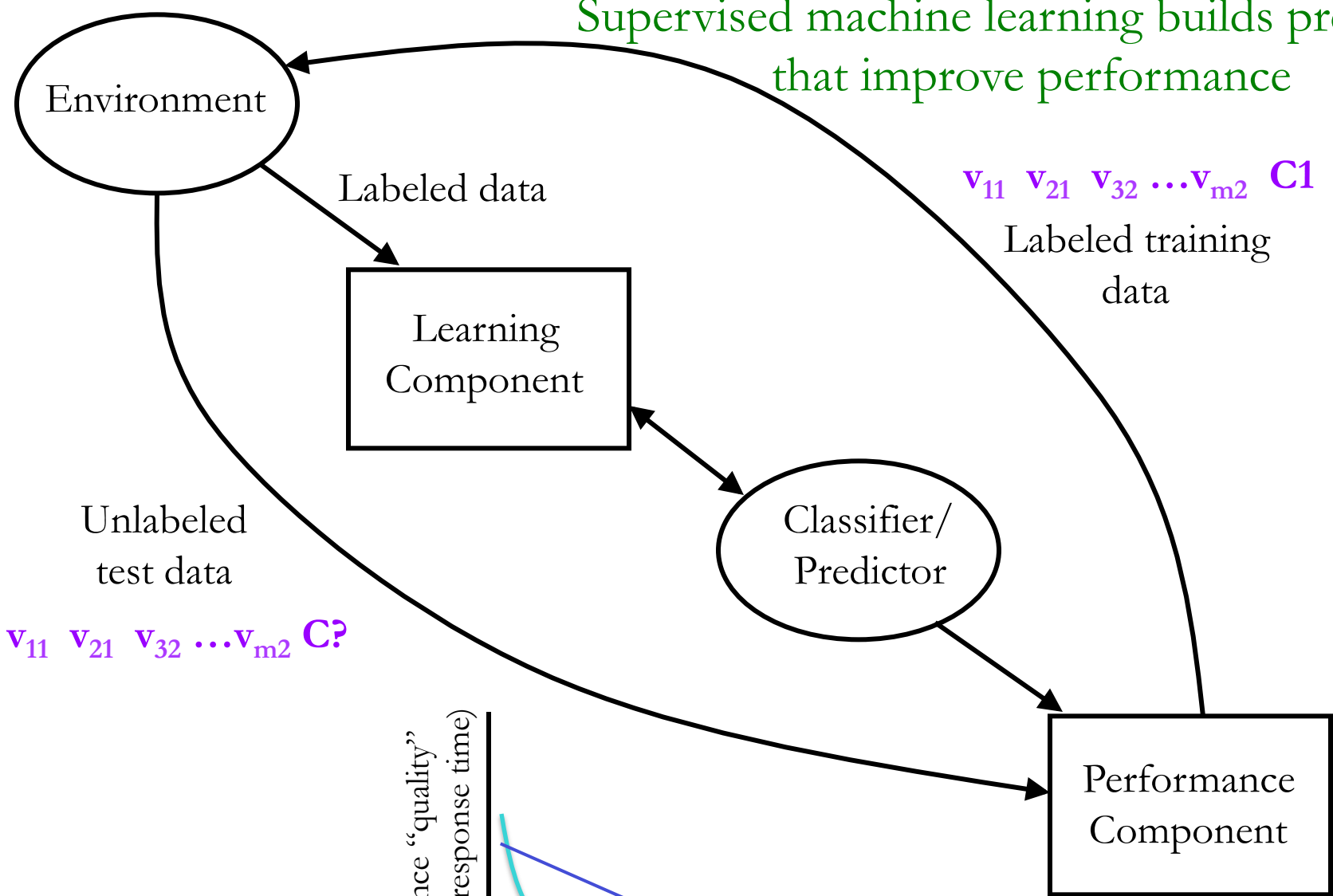
- Read Section 7.1 through 7.2 of *ArtInt* and

ArtInt: Poole and Mackworth, *Artificial Intelligence 2E*

at <http://artint.info/2e/html/ArtInt2e.html>

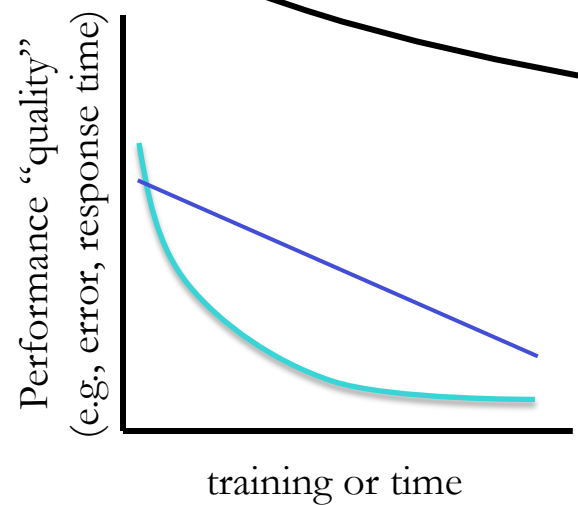
to include slides at <http://artint.info/2e/slides/ch04/lect1.pdf>

Supervised machine learning builds predictors that improve performance



v_{11} v_{21} v_{32} $\dots v_{m2}$ $C1$
Labeled training data

Unlabeled test data
 v_{11} v_{21} v_{32} $\dots v_{m2}$ $C?$



Empirical (aka data-driven) Supervised Learning

Given: a set of classified objects (a *training* data set)

Find: a classifier or predictor (for predicting class membership of unclassified data – a *test* set)

An example training set:

Index	V_1	V_2	V_3	\cdot	\cdot	\cdot	\cdot	V_m	C
1	v_{11}	v_{21}	v_{31}					v_{m1}	c_1
2	v_{12}	v_{22}	v_{32}					v_{m2}	c_2
\cdot	\cdot	\cdot	\cdot						
n	v_{1n}	v_{2n}	v_{3n}					v_{mn}	c_n

Subsequent example (but NOT general) assumptions:
two values per variable, two classes

Bayesian Classifier

Given a vector $V = \{v_{11}, v_{22}, v_{31}, \dots, v_{m2}, c_?\}$

Compute:

$$P(\mathbf{c}_1 | v_{11}, v_{22}, v_{31}, \dots, v_{m2}) = P(\mathbf{c}_1, v_{11}, v_{22}, v_{31}, \dots, v_{m2}) / P(v_{11}, v_{22}, v_{31}, \dots, v_{m2})$$

Compute:

$$P(\mathbf{c}_2 | v_{11}, v_{22}, v_{31}, \dots, v_{m2}) = P(\mathbf{c}_2, v_{11}, v_{22}, v_{31}, \dots, v_{m2}) / P(v_{11}, v_{22}, v_{31}, \dots, v_{m2})$$

Classify V as in \mathbf{c}_1 or \mathbf{c}_2 , whichever yields higher *conditional* probability

Background: $P(a | y) = P(a, y)/P(y)$, where (a, y) is same as $(a \wedge y)$ and $(a \text{ and } y)$

Bayesian Classifier cont

Note that denominators are equal, so if we are only interested in most probable, we need only compute the numerators

Compute:

$$P(\mathbf{c}_1 | v_{11}, v_{22}, v_{31}, \dots, v_{m2}) \text{ propto } P(\mathbf{c}_1, v_{11}, v_{22}, v_{31}, \dots, v_{m2}) / \cancel{P(v_{11}, v_{22}, v_{31}, \dots, v_{m2})}$$

Compute:

$$P(\mathbf{c}_2 | v_{11}, v_{22}, v_{31}, \dots, v_{m2}) \text{ propto } P(\mathbf{c}_2, v_{11}, v_{22}, v_{31}, \dots, v_{m2}) / \cancel{P(v_{11}, v_{22}, v_{31}, \dots, v_{m2})}$$

Classify V as in \mathbf{c}_1 or \mathbf{c}_2 , whichever yields higher *joint* probability

Background: the chain rule:

$$P(c, v1, v2, v3, v4)$$

A factorization ordering

$$= P(c) * P(v1 | c) * P(v2 | c, v1) * P(v3 | c, v1, v2) * P(v4 | c, v1, v2, v3)$$

$P(c, v1)$

$P(c, v1, v2)$

$P(c, v1, v2, v3)$

$P(c, v1, v2, v3, v4)$

$$P(c, v1, v2, v3, v4)$$

An alternative ordering

$$= P(v4) * P(v2 | v4) * P(v3 | v4, v2) * P(v1 | v4, v2, v3) * P(c | v4, v2, v3, v1)$$

Background: conditional independence

$$P(c, v_1, v_2, v_3, v_4)$$

$$= P(c) * P(v_1 | c) * P(v_2 | c, v_1) * P(v_3 | c, v_1, v_2) * P(v_4 | c, v_1, v_2, v_3)$$

$$= P(c) * P(v_1 | c) * P(v_2 | c) * P(v_3 | c) * P(v_4 | c)$$

if v_i s are independent conditioned on c (or conditionally independent given c)

Example:

Skin-cover in { hair, feathers, scales }

Heart in {4-chambers, 3-chambers, 2 chambers}

Transport in {walk, fly, swim}

Class in {mammal, bird, fish}

$$P(\text{mammal, hair, 4-chamber, walk})$$

$$\approx P(\text{mammal}) * P(\text{hair} | \text{mammal}) * P(\text{4-chamber} | \text{mammal}) * P(\text{walk} | \text{mammal})$$

$$P(\text{fish, scale, 2-chamber, swim}) \approx P(\text{fish}) * P(\text{scale} | \text{fish}) * P(\text{2-chamber} | \text{fish}) * P(\text{swim} | \text{fish})$$

$$P(\text{bird, hair, 3-chamber, swim}) \approx P(\text{bird}) * P(\text{hair} | \text{bird}) * P(\text{3-chamber} | \text{bird}) * P(\text{swim} | \text{bird})$$

Example:

Language in {English, Mandarin, Hindi, Spanish, French, ..., Tsalagi, Esperanto}

Country in {United States, China, India, Spain, Germany, Mexico, Canada, ...}

$$P(\text{US, English}) \approx P(\text{US}) * P(\text{English} | \text{US})$$

$$P(\text{US, English, Mandarin, French}) \approx P(\text{US}) * P(\text{English} | \text{US}) * P(\text{Mandarin} | \text{US}) * P(\text{French} | \text{US})$$

Naive Bayesian Classifier

Given a vector $V = \{v_{11}, v_{22}, v_{31}, \dots, v_{m2}, c_?\}$

$P(c_1 | v_{11}, v_{22}, v_{31}, \dots, v_{m2})$ proportional to $P(c_1, v_{11}, v_{22}, v_{31}, \dots, v_{m2})$

$$= P(v_{11} | v_{22}, v_{31}, \dots, v_{m2}, c_1) P(v_{22} | v_{31}, \dots, v_{m2}, c_1), \dots, P(v_{m2} | c_1) P(c_1)$$

$$= P(v_{11} | c_1) P(v_{22} | c_1) P(v_{31} | c_1) \dots P(v_{m2} | c_1) P(c_1)$$

under assumption that V_i 's are independent conditioned on C

$P(c_2 | v_{11}, v_{22}, v_{31}, \dots, v_{m2})$ proportional to $P(c_2, v_{11}, v_{22}, v_{31}, \dots, v_{m2})$

$$= P(v_{11} | v_{22}, v_{31}, \dots, v_{m2}, c_2) P(v_{22} | v_{31}, \dots, v_{m2}, c_2) \dots P(v_{m2} | c_2) P(c_2)$$

$$= P(v_{11} | c_2) P(v_{22} | c_2) P(v_{31} | c_2) \dots P(v_{m2} | c_2) P(c_2)$$

under assumption that V_i 's are independent conditioned on C

Classify V as in c_1 or c_2 , whichever yields higher **joint probability under assumption of conditional independence**

Learning a Naive Bayesian Classifier

View probabilities as proportions computed over training set.

$$\begin{aligned} & P(v_{11} | c_1) * P(v_{22} | c_1) * P(v_{31} | c_1) * \dots * P(v_{m2} | c_1) * P(c_1) \\ &= [v_{11}, c_1] / [c_1] * [v_{22}, c_1] / [c_1] * [v_{31}, c_1] / [c_1] * \dots * [v_{m2}, c_1] / [c_1] * [c_1] / [] \end{aligned}$$

where `[conditions]` is the number of objects/rows in the training set that satisfy all the `conditions`. So `[v11, c1]` is the number of training data that are members of `c1` and have `V1 = v11`, `[c1]` is the number of training objects in `c1`, `[]` is the total number of training objects.

Learning in this case, is a matter of counting the number of rows in training data in which various conditions satisfied:

- each class/variable-value pair
- each class
- total number of rows

$$\begin{aligned}
& P(v_{11}, v_{22}, v_{31}, c_1) \\
&= P(v_{11} \mid v_{22}, v_{31}, c_1) * P(v_{22} \mid v_{31}, c_1) * P(v_{31} \mid c_1) * P(c_1) \\
&= [v_{11}, v_{22}, v_{31}, c_1] / [v_{22}, v_{31}, c_1] * [v_{22}, v_{31}, c_1] / [v_{31}, c_1] * [v_{31}, c_1] / [c_1] * [c_1] / [] \\
&= [v_{11}, v_{22}, v_{31}, c_1] / []
\end{aligned}$$

Given a vector $V = \{1, -1, 0, \dots, 1\}$

Compute: There are often unknown values during training (and test)

$P(-1 | 1, -1, 0, \dots, 1)$ as $P(1|-1)P(-1|-1)P(0|-1) \dots P(1|-1)P(-1)$

$P(1 | 1, -1, 0, \dots, 1)$ as $P(1|1)P(-1|1)P(0|1) \dots P(1|1)P(1)$

Classify V as in $c1$ or $c2$, whichever yields higher probability

$$\begin{array}{cccccc}
 & V1 & V2 & V3 & \dots\dots\dots & Vm \\
 c1 & \left[\begin{array}{ccccc} [v11,c1] & [v21,c1] & [v31,c1] & \dots\dots\dots & [vm1,c1] \\ [v12,c1] & [v22,c1] & [v32,c1] & \dots\dots\dots & [vm2,c1] \end{array} \right] & & & & [c1] \\
 c2 & \left[\begin{array}{ccccc} [v11,c2] & [v21,c2] & [v31,c2] & \dots\dots\dots & [vm1,c2] \\ [v12,c2] & [v22,c2] & [v32,c2] & \dots\dots\dots & [vm2,c2] \end{array} \right] & & & & [c2] \\
 & [v11] & [v21] & [v31] & & [vm1] & [\\
 & [v12] & [v22] & [v32] & & [vm2] &]
 \end{array}$$

Consider an (multidimensional) array implementation of int, and estimate $P(v_{ij}|c_k)$ as $([v_{ij},c_k]+1) / ([c_k]+2)$, and $P(c_k)$ as $([c_k]+1) / ([]+2)$

Number of classes

Number of V_i values

Pseudo-counts to avoid the zero-probability problem and to mitigate overfitting